

Random Features for Dot Product Kernels and Beyond

Dissertation

submitted to the

Doctoral School 130: Computer Science,
Telecommunications and Electronics of Paris
of the Sorbonne Université

*in partial fulfillment of the requirements for the degree of
Doctor of Philosophy*

Author:

Jonas WACKER

Presented on the 12th July 2022 to the committee in charge:

Reviewer	Edwin V. BONILLA	<i>CSIRO's Data61 and Australian National University, Australia</i>
Reviewer/ Examiner	Zoltán SZABÓ	<i>London School of Economics and Political Science, United Kingdom</i>
Examiner	Alessandro RUDI	<i>INRIA and École Normale Supérieure, Paris, France</i>
Examiner	Maria A. ZULUAGA	<i>EURECOM, France</i>
Advisor	Maurizio FILIPPONE	<i>EURECOM, France</i>
Invited	Damien GARREAU	<i>Université Côte d'Azur and INRIA Sophia Antipolis, France</i>



Caractéristiques aléatoires pour les noyaux de produit scalaire et les noyaux similaires

Thèse

soumise à

l'École Doctorale 130: Informatique, Télécommunications
& Électronique (EDITE) de Paris
de la Sorbonne Université

pour l'obtention du grade de docteur

présentée par:

Jonas WACKER

et soutenue publiquement le 12 juillet devant un jury composé de:

Rapporteur	Edwin V. BONILLA	<i>CSIRO's Data61 and Australian National University, Australia</i>
Rapporteur/ Examineur	Zoltán SZABÓ	<i>London School of Economics and Political Science, United Kingdom</i>
Examineur	Alessandro RUDI	<i>INRIA and École Normale Supérieure, Paris, France</i>
Examineur	Maria A. ZULUAGA	<i>EURECOM, France</i>
Directeur de Thèse	Maurizio FILIPPONE	<i>EURECOM, France</i>
Invité	Damien GARREAU	<i>Université Côte d'Azur and INRIA Sophia Antipolis, France</i>



Abstract

Dot product kernels, such as polynomial and exponential (softmax) kernels, are among the most widely used kernels in machine learning, as they enable modeling the interactions between input features, which is crucial in applications like computer vision, natural language processing, and recommender systems. However, a fundamental drawback of kernel-based statistical models is their limited scalability to a large number of inputs, which requires resorting to approximations.

In this thesis, we study techniques to linearize kernel-based methods by means of random feature approximations and we focus on the approximation of polynomial kernels and more general dot product kernels to make these kernels more useful in large scale learning.

We begin by analyzing a class of random feature approximations for polynomial kernels from the literature that we refer to as *Product-Sketches*. In particular, we focus on a variance analysis as a main tool to study and improve the statistical efficiency of such sketches. We thus develop structured and complex generalizations of Product-Sketches, and elucidate conditions under which they can yield significant variance improvements. These generalizations allow us to construct a structured *Complex-to-Real (CtR) TensorSRHT* sketch that achieves state-of-the-art kernel approximation errors as well as regression and classification performance, while being faster than comparable sketches in the literature.

Moreover, we combine the Product-Sketches discussed in our thesis to approximate more general dot product kernels and Gaussian kernels. The aforementioned variance analysis hereby becomes the foundation to develop a data-driven optimization approach to allocate the number of random features between Product-Sketches of different degrees. We describe the improvements brought by these contributions with extensive experiments on a variety of tasks and data sets.

Lastly, we show that an optical hardware (the Optical Processing Unit - OPU), computes random feature approximations that are closely related to the ones studied in this thesis. We derive the associated limiting kernel, and we demonstrate its usefulness for a set of computer vision tasks. We show that the OPU can potentially be faster than a recent GPU while using less energy, making such devices interesting from a sustainability viewpoint.

Résumé

Les noyaux de produit scalaire, tels que les noyaux polynomiaux et exponentiels (softmax), sont parmi les noyaux les plus utilisés en apprentissage automatique, car ils permettent de modéliser les interactions entre les composantes des vecteurs d'entrée, ce qui est crucial dans des applications telles que la vision par ordinateur, le traitement du langage naturel et les systèmes de recommandation. Cependant, un inconvénient fondamental des modèles statistiques basés sur les noyaux est leur évolutivité limitée à un grand nombre de données d'entrée, ce qui nécessite de recourir à des approximations.

Dans cette thèse, nous étudions des techniques pour linéariser les méthodes à base de noyaux au moyen d'approximations de *caractéristiques aléatoires* et nous nous concentrons sur l'approximation de noyaux polynomiaux et de noyaux plus généraux de produit scalaire pour rendre ces noyaux plus utiles dans l'apprentissage à grande échelle.

Nous commençons par analyser une classe d'approximations de caractéristiques aléatoires pour les noyaux polynomiaux de la littérature que nous appelons *croquis de produit* (*Product-Sketches*). En particulier, nous nous concentrons sur une analyse de variance comme outil principal pour étudier et améliorer l'efficacité statistique de ces croquis (sketches). Nous développons donc des généralisations structurées et complexes de croquis de produit, et nous élucidons les conditions dans lesquelles elles peuvent produire des améliorations significatives de la variance. Ces généralisations nous permettent de construire un croquis de produit structuré, nommé *Complex-to-Real (CtR) TensorSRHT*, qui atteint des erreurs d'approximation de noyau faibles ainsi que des performances de régression et de classification de pointe, tout en étant plus rapide que les caractéristiques aléatoires comparables dans la littérature.

De plus, nous combinons les croquis de produit discutés dans notre thèse pour approximer des noyaux de produit scalaire plus généraux et des noyaux gaussiens. L'analyse de variance susmentionnée devient ainsi la base pour développer une approche d'optimisation basée sur les données pour allouer le nombre de caractéristiques aléatoires entre les croquis de produit de différents degrés. Nous décrivons les améliorations apportées par ces contributions à l'aide d'expériences approfondies sur une variété de tâches et d'ensembles de données.

Enfin, nous montrons qu'un matériel optique (nommé *Optical Processing Unit - OPU*), calcule des approximations de caractéristiques aléatoires qui sont étroitement liées à celles étudiées dans cette thèse. Nous dérivons le noyau limite associé, et nous démontrons son utilité pour un ensemble de tâches de vision par ordinateur. Nous montrons que l'OPU peut potentiellement être plus rapide qu'un GPU récent tout en utilisant moins d'énergie, ce qui rend ces dispositifs intéressants du point de vue de la durabilité.

Acknowledgements

During my PhD, I have been granted a great amount of freedom to explore a scientific topic of my choice during more than three years. This experience has taught me how to become an independent researcher and to immerse myself in a less well-known, but fundamentally important research topic.

First and foremost, I would like to thank my supervisor Maurizio Filippone who provided me with a safe research environment in which I could freely develop the contents of this thesis. I am especially grateful for his trust and his optimism that I believe to be a rare gift given to a PhD student.

Exploring fundamental research ideas may not yield immediate rewards to financial stakeholders, which is why I would like to thank the French *Agence Nationale de Recherche* and thereby the French government for their trust and financial support for my PhD. Being given such a support without being bound to commercial use cases allowed me to work on a topic that may shape research in the long term.

Next, I would like to thank Motonobu Kanagawa for his incredible endurance and support in reviewing the contents of this thesis. Thanks to him, I could sharpen the main motivation and contributions of this work, and present results with the scientific rigor that is expected in this field.

I would also like to thank the academic jury that voluntarily accepted to review my thesis and attend my PhD defense. I have the privilege to count true experts in the field among the jury members and their feedback is of incredible value to this work.

Last but not least, I would like to thank my parents and all my friends who accompanied me during my PhD experience at the Côte d'Azur. It is thanks to them that this time was also an incredibly enriching social and cross-cultural experience from which I probably drew as much personal development as from the PhD itself. A special thanks hereby goes to Sagar, Marine and Julien, with whom I shared my flat, and who supported me during the good and bad times with helpful conversations and first-class culinary experiences. Another special thanks goes to my friends from the FROST group with whom I have shared some of the most joyful moments of my life.

Jonas Wacker

Contents

Preface	vi
Mathematical Notation	ix
Acronyms and Abbreviations	x
1 Introduction	1
1.1 Positive Definite Kernels	1
1.2 Gaussian Process Regression	3
1.3 Applications for Dot Product Kernels	5
1.4 The Optical Processing Unit	13
2 Random Feature Approximations of Shift-Invariant Kernels and Polynomial Kernels	15
2.1 Random Fourier Features for Shift-Invariant Kernels	16
2.2 Random Features for Polynomial Kernels	21
3 Analysis and Extension of Product-Sketches	35
3.1 Variance Analysis for Real-Valued Product-Sketches	36
3.2 Complex-Valued Product-Sketches	38
3.3 Variance of Complex-Valued Product-Sketches	40
3.4 Probabilistic Error Bounds for Rademacher Sketches	42
3.5 Structured Product-Sketches	43
4 Complex-to-Real Product-Sketches	58
4.1 Transforming Complex into Real Sketches	58
4.2 Variances of CtR Product-Sketches	59
4.3 Comparison against the State-of-the-Art	67
5 Approximating Dot Product Kernels	76
5.1 Maclaurin Expansion of Dot Product Kernels	77
5.2 Random Sketch based on the Maclaurin Expansion	78
5.3 Optimization for a Truncated Maclaurin Approximation	79
5.4 Empirical Evaluation of the Optimized Maclaurin Method	90

5.5	Pathology when Approximating the Gaussian Kernel	102
5.6	Discussion	109
6	The Optical Processing Unit	111
6.1	The Optical Processing Unit (OPU)	111
6.2	Computing the Kernel	112
6.3	Classification Experiments	115
6.4	Projection Time and Energy Consumption	117
6.5	Discussion	118
7	Conclusions and Future Work	120
A	Appendix of Chapter 2	125
A.1	Sum of Independent and Uniformly Distributed Variables	125
B	Appendix of Chapter 3	126
B.1	Proofs for Sections 3.1-3.4	126
B.2	Proofs for Section 3.5	128
B.3	Gaussian Processes with Complex Random Features	134
C	Appendix of Chapter 4	141
C.1	Pseudo-Variations of Gaussian and Rademacher Polynomial Sketches	141
C.2	Gaussian and Rademacher CtR Variance Advantage over their Real- Valued Analogs	142
C.3	Pseudo-Variance of TensorSRHT	145
D	Appendix of Chapter 5	147
D.1	Convex Surrogate Functions for Stacked TensorSRHT Variances . .	147
D.2	Additional Experiments	152
E	Appendix of Chapter 6	157
E.1	Proof of Theorem 6.2.2	157

Preface

This thesis is divided into 7 main chapters. Chapter 1 is an introductory chapter and Chapter 2 is a literature review that provides a context for the contributions made in this thesis. Chapters 3-5 should be read consecutively since intermediate contributions build the basis for the contributions that follow. Chapter 6 can be read as an independent chapter, but some technical aspects are closely related to Chapters 3-5. Chapter 7 is a concluding chapter that discusses future work. We briefly summarize each chapter including the contributions made therein in the following.

Chapter 1 is an introductory chapter that motivates the use of random features, in particular for dot product kernels. We define positive definite kernels, discuss associated scalability bottlenecks for Gaussian processes regression, and provide an intuitive introduction to modern example applications that make use of random features for dot product kernels.

Chapter 2 is a literature review that discusses random feature approximations for shift-invariant kernels and for polynomial kernels. We further draw connections between the two. The main purpose of this chapter is to provide a context for Chapter 3. Some of the methods discussed also serve as baselines for our experiments in subsequent chapters.

In **Chapter 3**, we discuss and extend random feature approximations for polynomial kernels that we refer to as *polynomial sketches*. We make several contributions to the literature in this chapter.

Our first contribution is the derivation of closed form variance formulas for a group of polynomial sketches complementing existing theoretical results. Prior to our work, only upper bounds for such variances had been derived in the literature. Our novel variance formulas make it possible to study conditions under which a certain polynomial sketch yields lower mean squared errors than another. In particular, we show for the first time that the variance of structured polynomial sketches is upper-bounded by the variance of their unstructured counterparts when the degree of the polynomial is odd.

Our second contribution in this chapter is to construct complex-valued generalizations of the aforementioned polynomial sketches. We show that complex polynomial sketches can yield much lower variances than their real-valued coun-

terparts, especially when the degree of the polynomial is large.

In **Chapter 4**, we convert the complex polynomial sketches introduced in Chapter 3 into real-valued ones, and discuss the implications of such complex-to-real (CtR) transformations on their variances. In particular, we prove that CtR Gaussian and Rademacher polynomial sketches yield lower variances than their real-valued analogs under the same conditions as it is the case for the complex polynomial sketches introduced before. We demonstrate empirically that our novel structured *CtR-TensorSRHT* sketch achieves state-of-the-art performance in terms of kernel approximation errors and feature construction time.

In **Chapter 5**, we show how any dot product kernel can be represented as a positively weighted sum of polynomial kernels using a Maclaurin series. This observation allows us to construct random features for dot product kernels and even for the Gaussian kernel using the polynomial sketches introduced in the chapters before. Using our derived variance formulas from Chapter 3, we develop a data-driven optimization approach to determine the number of random features assigned to each polynomial sketch used in this construction. In an extensive empirical evaluation, we show that our optimized sketches achieve state-of-the-art performance for the tasks of polynomial and Gaussian kernel approximation. Lastly, we identify a pathology of our approach when approximating a Gaussian kernel with short length scales and propose a cure for this pathology.

Chapter 6 deals with the *Optical Processing Unit (OPU)*, an optical hardware used to compute random projections. We show that such random projections naturally give rise to a limiting kernel that is closely related to a degree-2 polynomial kernel. In addition, we numerically modify these random feature maps to obtain higher-degree polynomials and derive associated variances in closed form. Finally, we benchmark projection times and energy consumption of the OPU against a recent GPU.

In **Chapter 7**, we conclude our work and discuss improvements that could be tackled in future work. In particular, we discuss the attention mechanism in Transformer neural network architectures and show that current random feature methods do not yield sufficiently good approximations in this case.

Chapters 3 and 5 are based on

- J. Wacker, M. Kanagawa, and M. Filippone. Improved random features for dot product kernels. *arXiv preprint arXiv:2201.08712*, 2022a
submitted to the *Journal of Machine Learning Research*.
- J. Wacker and M. Filippone. Local random feature approximations of the gaussian kernel. In *Proceedings of the 26th International Conference on Knowledge-Based and Intelligent Information & Engineering Systems*, Procedia Computer Science. Elsevier, 2021.

Chapter 4 is based on

- J. Wacker, R. Ohana, and M. Filippone. Complex-to-real random features for polynomial kernels. *arXiv preprint arXiv:2202.02031*, 2022b.

Chapter 6 is based on the following conference paper:

- R. Ohana, J. Wacker, J. Dong, S. Marmin, F. Krzakala, M. Filippone, and L. Daudet. Kernel computations from large-scale random features obtained by optical processing units. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 9294–9298, 2020.

The following additional conference paper was published during the course of writing this thesis:

- J. Wacker, M. Ladeira, and J. E. V. Nascimento. Transfer learning for brain tumor segmentation. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries - 6th International Workshop, BrainLes 2020, Held in Conjunction with MICCAI 2020.*, volume 12658 of *Lecture Notes in Computer Science*, pages 241–251. Springer, 2020.

However, this contribution was made completely independently from the other works and we therefore consider it to be out of scope for this thesis.

Mathematical Notation

We reserve italic letters (e.g., a) for scalars. We use lower case boldface letters (e.g., \mathbf{a}) for column vectors, and upper case boldface letters (e.g., \mathbf{A}) for matrices.

\mathbf{A}^H	Hermitian transpose of \mathbf{A}
\mathbf{A}^\top	transpose of \mathbf{A}
$\ \mathbf{a}\ $	Euclidean norm of \mathbf{a}
$\ \mathbf{A}\ _F$	Frobenius norm of \mathbf{A}
$(\mathbf{A})_{i,j}$	(i, j) -th entry of \mathbf{A}
$(\mathbf{a}_1, \dots, \mathbf{a}_n)$	horizontal concatenation of the column vectors $\mathbf{a}_1, \dots, \mathbf{a}_n$ into a matrix
$\{\mathbf{A}_i\}_{i=1}^N$	set consisting of the matrices $\mathbf{A}_1, \dots, \mathbf{A}_N$
\mathbf{I}_n	identity matrix of dimension n ; the subscript may be omitted when no confusion arises
$\text{tr}(\mathbf{A})$	trace of \mathbf{A}
$\text{Re}\{\mathbf{A}\}$	real part of \mathbf{A}
$\text{Im}\{\mathbf{A}\}$	imaginary part of \mathbf{A}
$\overline{\mathbf{A}}$	conjugate of a matrix \mathbf{A}
\mathbf{e}_i	i -th standard basis vector, i.e., the i -th column of \mathbf{I}
$\mathbf{0}$	zero vector with proper dimension
$:=, =:$	definitions from the right and from the left
$\text{diag}(\mathbf{a})$	diagonal matrix with the vector \mathbf{a} being its diagonal
$\text{vec}(\mathbf{A})$	column vector obtained by stacking column-wise the elements of \mathbf{A}
$\mathbf{A} \otimes \mathbf{B}$	Kronecker product of two matrices \mathbf{A} and \mathbf{B}
$\mathbf{A} \odot \mathbf{B}$	element-wise product of two matrices \mathbf{A} and \mathbf{B}
$\mathbb{E}[\mathbf{A}]$	expectation of \mathbf{A}
$\mathbb{V}[a]$	variance of a
$\text{Pr}\{\mathcal{E}\}$	probability of an event \mathcal{E}
$\mathcal{CN}(\boldsymbol{\mu}, \mathbf{K})$	Circularly symmetric complex Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix \mathbf{K}
\mathbb{R}	set of real numbers
\mathbb{C}	set of complex numbers
\mathbb{N}	set of natural numbers
\mathcal{S}^{d-1}	$(d - 1)$ -dimensional unit sphere in d -dimensional space

Acronyms and Abbreviations

CPU	Central Processing Unit
GPU	Graphical Processing Unit
OPU	Optical Processing Unit
NLP	Natural Language Processing
RF	Random Features
w.r.t.	with respect to
w.l.o.g.	without loss of generality
a.s.	almost surely
i.i.d.	independent and identically distributed
e.g.	exempli gratia, for example
i.e.	id est, that is

Chapter 1

Introduction

The contents of this thesis are all related to random features for polynomial kernels and more general dot product kernels. The purpose of random features is to accelerate statistical models by constructing approximations of such kernels. Before delving into the details about random features, this introductory chapter serves to provide the reader with an intuitive understanding about the *benefits* of such approximations.

In the following, we define positive definite kernels including the subclasses of dot product kernels and shift-invariant kernels that are relevant to this thesis. We then discuss several example applications that suffer from a scalability bottleneck when exact kernel evaluations are used, and show how this bottleneck is resolved through the use of random features. Several of these examples are related to dot product kernels motivating the need for random feature approximations for this particular class of kernels. Lastly, we provide a brief introduction to the Optical Processing Unit (OPU), an optical hardware that we worked with, and from which we drew the inspiration for the main contents of this thesis.

For now, we treat random feature maps as black-box functions and postpone the introduction of first concrete examples to Chapter 2.

1.1 Positive Definite Kernels

Positive definite kernels are ubiquitous in machine learning. They are traditionally used in kernel methods (Scholkopf and Smola, 2002) and Gaussian processes (Rasmussen and Williams, 2006) to model non-linear phenomena in a theoretically principled way. Associated learning algorithms are often simple to solve and can therefore be equipped with theoretical guarantees. This is why traditional learning algorithms using positive definite kernels are often preferred in applications where the learning machine should behave in an anticipated manner, e.g., when high-stake decision-making is involved or when safety is required.

At the same time, many successful modern deep learning approaches embed kernel methods into neural network architectures to either yield richer feature rep-

representations or to incorporate prior domain knowledge. We discuss such examples in the remainder of this chapter.

We now give a brief formal introduction to positive definite kernels. A more complete introduction can be found in, e.g., [Scholkopf and Smola \(2002, Chapter 2\)](#) and [Steinwart and Christmann \(2008, Chapter 4\)](#). The kernels we discuss in this work are generally real-valued, but their approximations can be complex-valued as we will see in [Chapters 2 and 3](#). Therefore, we provide a more general definition for complex positive definite kernels in the following.

Positive definite kernel. Let \mathcal{X} be a nonempty set. A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ is a positive definite kernel if for every $x_1, \dots, x_m \in \mathcal{X}$ and $c_1, \dots, c_m \in \mathbb{C}$

$$\sum_{i=1}^m \sum_{j=1}^m c_i \bar{c}_j k(x_i, x_j) \geq 0 \quad (1.1)$$

holds, where $\bar{\cdot}$ denotes the complex conjugate.

We call the (positive semi-definite) matrix $\mathbf{K} \in \mathbb{C}^{m \times m}$ with elements $(\mathbf{K})_{i,j} = k(x_i, x_j)$ the Gram (or kernel) matrix of k with respect to x_1, \dots, x_m . \mathbf{K} is Hermitian, i.e., we have $k(x_i, x_j) = \overline{k(x_j, x_i)}$. For a real-valued kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, this simplifies to $k(x_i, x_j) = k(x_j, x_i)$. In this case, we use $c_1, \dots, c_m \in \mathbb{R}$ in [Eq. \(1.1\)](#).

Feature map. Kernels can equivalently be defined through feature maps. Let \mathcal{X} be a nonempty set. A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ is a positive definite kernel if there exists a \mathbb{C} -Hilbert space \mathcal{H} and a map $\phi : \mathcal{X} \rightarrow \mathcal{H}$ such that for all $x, y \in \mathcal{X}$

$$k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$$

holds, where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is an inner product in \mathcal{H} . We call ϕ a feature map and \mathcal{H} a feature space of k .

In this work, we generally assume $\mathcal{X} \subseteq \mathbb{R}^d$. We denote vector-valued inputs by bold-faced letters. Hence, we write $k(\mathbf{x}, \mathbf{y})$ for some $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Depending on the kernel, \mathcal{H} can be infinite dimensional.

The goal of this work is therefore to construct low-dimensional *random* feature maps yielding an *approximate* kernel \hat{k} . We denote these feature maps by $\Phi_{\mathcal{R}} : \mathbb{R}^d \rightarrow \mathbb{R}^D$ if they are real-valued and by $\Phi_{\mathcal{C}} : \mathbb{R}^d \rightarrow \mathbb{C}^D$ if they are complex-valued, and we drop the subscript when the context is clear. For unbiased random feature maps, the approximate kernel $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{R}}(\mathbf{x})^{\top} \Phi_{\mathcal{R}}(\mathbf{y})$ or $\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{C}}(\mathbf{x})^{\top} \overline{\Phi_{\mathcal{C}}(\mathbf{y})}$ converges in probability to $k(\mathbf{x}, \mathbf{y})$ as D increases. This leads to a trade-off between approximation quality and computational costs that both increase with D .

Classes of kernel functions. For the purpose of this work, we distinguish two major classes of kernels, *shift-invariant* or *stationary* kernels, and *dot product* kernels (see [Rasmussen and Williams, 2006](#), Chapter 4). Shift-invariant kernels depend only on the lag vector $\boldsymbol{\tau} := \mathbf{x} - \mathbf{y} \in \mathbb{R}^d$, i.e., we can overload the notation and write $k(\mathbf{x}, \mathbf{y}) = k(\boldsymbol{\tau})$ for two inputs $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. A well-known example is the Gaussian kernel $\exp(-\|\boldsymbol{\tau}\|^2/(2l^2))$ with length scale parameter $l > 0$. Dot product kernels on the other hand depend only on the inner product $\mathbf{x}^\top \mathbf{y} \in \mathbb{R}$, so we write $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x}^\top \mathbf{y})$. Well-known examples include the polynomial kernel $(\gamma \mathbf{x}^\top \mathbf{y} + \nu)^p$ for some $\gamma, \nu \geq 0, p \in \mathbb{N}$ and the exponential dot product kernel $\exp(\mathbf{x}^\top \mathbf{y}/l^2)$ for some $l > 0$. We will cover random feature approximations for shift-invariant kernels and polynomial kernels in detail in Chapters 2, 3 and 4. General dot product kernels are covered later on in Chapter 5.

Scalability bottleneck. Statistical models such as kernel methods ([Scholkopf and Smola, 2002](#)) and Gaussian processes ([Rasmussen and Williams, 2006](#)) suffer from a scalability bottleneck when direct kernel evaluations are used. This is because the kernel needs to be evaluated for all pairs of inputs, which implies a time and memory complexity that is at least quadratic in the number of input points.

The scalability issue has been a focus of research since the earliest literature ([Wahba, 1990](#), Chapter 7), and many approximation methods for reducing the computational costs have been developed (e.g., [Williams and Seeger 2000](#); [Rahimi and Recht 2007](#); [Titsias 2009](#); [Hensman et al. 2018](#)). The *random features* that we discuss in this work, are a particularly popular approximation method that was originally introduced for *shift-invariant* kernels in the seminal work by [Rahimi and Recht \(2007\)](#). We follow this line of research and show in the following how such approximations can be used to improve the scalability of various statistical models that make use of kernels.

We start with the example of Gaussian process regression that is used in many of the numerical experiments that we carried out for this thesis.

1.2 Gaussian Process Regression

Gaussian process regression (GPR) ([Rasmussen and Williams, 2006](#), Chapter 2) is a Bayesian non-parametric approach to supervised learning that makes use of positive-definite kernels. It is an attractive modelling choice as it provides uncertainty estimates through predictive variances next to the actual prediction. At the same time, the hyperparameters of the kernel function k can be obtained through a gradient-based optimization of the *log marginal likelihood* ([Rasmussen and Williams, 2006](#), Chapter 5.4) avoiding time-consuming cross-validation.

However, GPR does not scale well to a large number of training points due to the necessity for algebraic manipulations of the Gram matrix. We introduce GPR

here and show how its scalability bottleneck can be overcome by means of random feature approximations.

Suppose a training data set $\mathcal{D} := \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ with $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ that we summarize in matrix notation as $\mathbf{X} := (\mathbf{x}_1, \dots, \mathbf{x}_N)^\top \in \mathbb{R}^{N \times d}$ and $\mathbf{y} := (y_1, \dots, y_N)^\top \in \mathbb{R}^N$. We assume that the observations \mathbf{y} have been generated from \mathbf{X} by an unknown latent function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that has been corrupted by independent Gaussian noise, i.e., $y_i = f(\mathbf{x}_i) + \epsilon_i$ with $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_{\text{noise}}^2)$ and $\sigma_{\text{noise}}^2 > 0$.

In GPR, the vector of function evaluations $\mathbf{f} := (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^\top \in \mathbb{R}^N$ is assumed to have a joint Gaussian distribution with mean $\boldsymbol{\mu} \in \mathbb{R}^N$ and covariance matrix $\mathbf{K}_{\mathbf{ff}} \in \mathbb{R}^{N \times N}$. We follow the standard approach and set $\boldsymbol{\mu}$ to zero here although more complex models exist (Rasmussen and Williams, 2006, Chapter 2.7). For now we assume to be working with real-valued GPs and introduce complex-valued GPs later on in this work. So the entries of $\mathbf{K}_{\mathbf{ff}}$ correspond to the evaluations of a real-valued positive definite kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ with $(\mathbf{K}_{\mathbf{ff}})_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ that determines the covariance of a pair of function values $f(\mathbf{x}_i)$ and $f(\mathbf{x}_j)$.

The task of GPR is to predict the latent function value at a new test input $\mathbf{x}_* \in \mathbb{R}^d$ given a training set \mathcal{D} . The predictive distribution of $f(\mathbf{x}_*) | \mathcal{D} \in \mathbb{R}$ can be computed in closed form, and it is $\mathcal{N}(\mu_*, \sigma_*^2)$ with

$$\mu_* := \mathbf{k}_{\mathbf{f}*}^\top (\mathbf{K}_{\mathbf{ff}} + \sigma_{\text{noise}}^2 \mathbf{I})^{-1} \mathbf{y} \quad \text{and} \quad \sigma_*^2 := k_{**} - \mathbf{k}_{\mathbf{f}*}^\top (\mathbf{K}_{\mathbf{ff}} + \sigma_{\text{noise}}^2 \mathbf{I})^{-1} \mathbf{k}_{\mathbf{f}*}, \quad (1.2)$$

where $\mathbf{k}_{\mathbf{f}*} = (k(\mathbf{x}_1, \mathbf{x}_*), \dots, k(\mathbf{x}_N, \mathbf{x}_*))^\top \in \mathbb{R}^N$, $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$ and $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the identity matrix.

Although being available in closed form, the computation of the GPR predictor can be expensive in practice. The computational bottleneck here is to solve the linear systems in Eq. (1.2), which costs $\mathcal{O}(N^3)$ time. Even storing the matrix $\mathbf{K}_{\mathbf{ff}}$ requires $\mathcal{O}(N^2)$ memory and becomes infeasible when N is large, typically greater than 10,000 for a regular machine, so that approximations become necessary. We now present an alternative GPR formulation that is amenable to the random feature approximations discussed in this thesis.

Explicit feature space formulation. If there exists a finite-dimensional feature map $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ such that $k(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$, it can be shown (Rasmussen and Williams, 2006, Chapter 2) that Eq. (1.2) can be reformulated as

$$\mu_* := \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{X})^\top \mathbf{y} / \sigma_{\text{noise}}^2 \quad \text{and} \quad \sigma_*^2 := \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{x}_*), \quad (1.3)$$

with $\mathbf{A} := \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) / \sigma_{\text{noise}}^2 + \mathbf{I}$ and $\Phi(\mathbf{X}) = (\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_N))^\top \in \mathbb{R}^{N \times D}$.

The feature space representation (1.3) changes the computational cost from $\mathcal{O}(N^3)$ to $\mathcal{O}(ND^2)$, and thus improves the scaling of GPR drastically if $D \ll N$. Unfortunately, exact feature maps can be high-dimensional, as is the case for, e.g., polynomial kernels, or even infinite dimensional rendering an explicit feature space

formulation infeasible. In such cases, we can use the random feature maps studied in this thesis as a drop-in replacement for Φ in Eq. (1.3), where the dimension D of Φ becomes a parameter determining the trade-off between computational costs and approximation quality of the exact GP predictor.

There exist further applications for positive definite kernels that, similar to GPR, suffer from a scalability bottleneck due to pairwise kernel evaluations. In the following, we introduce some example applications from the literature, which make use of dot product kernels and their random feature approximations in particular. Therefore, these applications motivate the necessity for random features for polynomial kernels and more general dot product kernels that are the focus of this manuscript.

1.3 Applications for Dot Product Kernels

While random features for shift-invariant kernels have been studied extensively in the past (see Liu et al. (2020) for a recent overview), there are relatively few works studying random features for the general class of dot product kernels. The purpose of this section is therefore to provide recent examples from the literature that benefit from such approximations.

As we will show in Chapter 5, any dot product kernel and even the Gaussian kernel can be formulated as a positively weighted sum of polynomial kernels. This turns polynomial kernels into very general function approximators to which we assign a great importance in this work. Polynomial kernels in turn model multiplicative interactions between the input features, which is why dot product kernels such as polynomial kernels are used in applications that make use of such interactions. Such applications include genomic data analysis (Aschard, 2016; Weissbrod et al., 2016), recommender systems (Rendle, 2010; Blondel et al., 2016), computer vision (Lin et al., 2018; Gao et al., 2016; Fukui et al., 2016), natural language processing (Yamada and Matsumoto, 2003; Chang et al., 2010; Vaswani et al., 2017), and evaluation metrics like the Kernel Inception Distance (KID) (Bińkowski et al., 2018)¹.

In the following, we focus on recent example applications in deep learning, namely bilinear pooling for fine-grained visual recognition (Lin et al., 2018), multi-modal bilinear pooling for visual question answering (Fukui et al., 2016), and the dot-product attention mechanism used in Transformers (Vaswani et al., 2017).

1.3.1 Bilinear Pooling

We begin by discussing applications for the so-called *bilinear pooling* operation that is used in fine-grained visual recognition (Lin et al., 2018) and has been extended to a multi-modal context for visual question answering (Fukui et al., 2016). We

¹We cover a random feature approximation of the KID in more detail in Chapter 4.

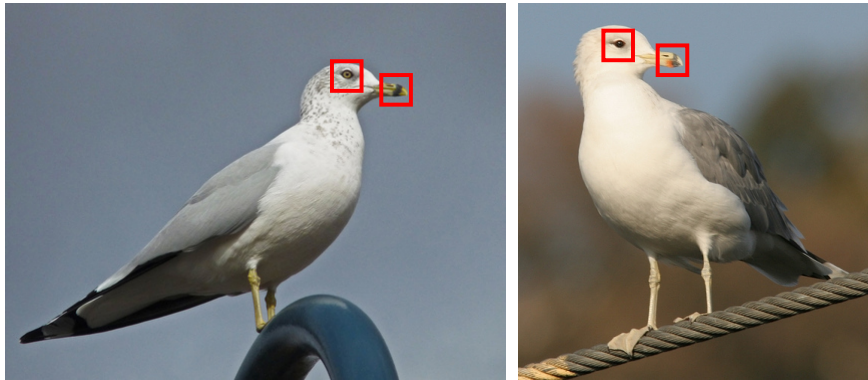


Figure 1.1: Comparing a Ring Billed Gull (left) to a California Gull (right) from the CUB-200 data set. The Ring Billed Gull has a yellow bill with a black ring around it, and a pale eye. The California Gull has a red spot on the lower bill and a dark eye.

cover the example of fine-grained visual recognition first, and show therein how bilinear pooling relates to explicit feature maps of a degree-2 polynomial kernel.

Fine-grained visual recognition. The task of fine-grained visual recognition is about the classification of pictures *within* a subordinate category. This classification usually requires the detection of fine-grained details. An example for such a task is the Caltech-UCSD Birds (CUB) 200 data set (Welinder et al., 2010) in which birds are categorized into 200 different species. Fig. 1.1 shows two example images of two similar bird species that can be distinguished by detecting details on their bills and eyes, respectively².

One way to solve this task would be to rely on anatomic annotations, locating different body parts of the bird on each picture. In this way, features could be extracted from each body part allowing for a detailed comparison. Although being provided in the CUB-200 data set, such annotations are expensive to obtain for a large number of images. This makes *holistic* models that globally process an image without annotations an attractive choice.

Lin et al. (2018) achieve remarkable improvements of holistic models on CUB-200 as well as other fine-grained visual recognition data sets, by applying a bilinear pooling operation to the outputs of a convolutional neural network. We define this operation and explain its scalability bottleneck in the following.

²See https://www.allaboutbirds.org/guide/Ring-billed_Gull/species-compare for an explanation. Accessed: 26/04/2022.

Bilinear Pooling (Lin et al., 2018). We consider a set of images \mathcal{I} , where an image $\mathbf{A} \in \mathcal{I}$ can be of arbitrary size. We then define a feature descriptor as a function $f : \mathcal{I} \times \mathcal{L} \rightarrow \mathbb{R}^d$ that describes an image \mathbf{A} at a location $l \in \mathcal{L}$, where the set of locations \mathcal{L} is defined by the feature descriptor and generally depends on \mathbf{A} .

Let $\mathbf{a} \otimes \mathbf{b} := \text{vec}(\mathbf{a}\mathbf{b}^\top) \in \mathbb{R}^{d_1 \cdot d_2}$ be the vectorized outer product of two vectors $\mathbf{a} \in \mathbb{R}^{d_1}$ and $\mathbf{b} \in \mathbb{R}^{d_2}$. The bilinear pooling operation is then defined as

$$\text{Bil}(f, \mathbf{A}) := \sum_{l \in \mathcal{L}} f(\mathbf{A}, l) \otimes f(\mathbf{A}, l) \in \mathbb{R}^{d^2}, \quad (1.4)$$

which forms a global image descriptor, since the sum over \mathcal{L} , referred to as a *pooling* operation, aggregates local image descriptions over image locations. f could correspond to feature extractors such as SIFT (Lowe, 1999), HOG (Dalal and Triggs, 2005), or a convolutional neural network (Lecun et al., 1998).

The term $f(\mathbf{A}, l) \otimes f(\mathbf{A}, l)$ in Eq. (1.4) contains second-order multiplicative interactions of the coordinates of $f(\mathbf{A}, l)$. These interactions correspond to the explicit feature map of a degree-2 polynomial kernel as shown in Scholkopf and Smola (2002, Proposition 2.1). Taking the inner product of the bilinearly pooled features of two different images $\mathbf{A}, \mathbf{B} \in \mathcal{I}$ thus yields a sum of degree-2 polynomial kernels:

$$\text{Bil}(f, \mathbf{A})^\top \text{Bil}(f, \mathbf{B}) = \sum_{l \in \mathcal{L}} \sum_{l' \in \mathcal{L}} (f(\mathbf{A}, l)^\top f(\mathbf{B}, l'))^2$$

The polynomial kernel here is defined as $k(\mathbf{x}, \mathbf{y}) = (\gamma \mathbf{x}^\top \mathbf{y} + \nu)^2$, where we set $\mathbf{x} = f(\mathbf{A}, l)$, $\mathbf{y} = f(\mathbf{B}, l')$, $\gamma = 1$ and $\nu = 0$.

Computational bottleneck. The dimension of the bilinear features $\text{Bil}(f, \cdot)$ (1.4) is d^2 , where d is the dimension of the features extracted by f . The quadratic scaling in d can render the bilinear pooling operation expensive. E.g., Lin et al. (2018) extract *relu5* outputs from a VGG-M network (Chatfield et al., 2014), which yields $d = 512$ dimensional outputs with $|\mathcal{L}| = 27^2$ for an RGB input image of size $448 \times 448 \times 3$ from the CUB-200 data set.

As there are 200 classes in the CUB-200 data set, a linear classifier used on top of the bilinear pooling features has $512^2 \cdot 200 \approx 52$ million parameters imposing high computational and memory costs. At the same time, storing 512^2 dimensional features for subsequent image comparison and retrieval tasks requires one megabyte of storage per image, which is around 10 times as much as storing the original image in JPEG format.

Hence, Gao et al. (2016) propose *Compact Bilinear Pooling* to reduce the dimensionality of $\text{Bil}(f, \cdot)$ (1.4) using TensorSketch (Pham and Pagh, 2013). More precisely, they replace $f(\mathbf{A}, l) \otimes f(\mathbf{A}, l)$ in Eq. (1.4) by $\Phi(f(\mathbf{A}, l))$, where $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ is a random feature map of the polynomial kernel. This can drastically reduce computational costs when $D \ll d^2$ holds. We will introduce such random feature maps as well as TensorSketch in Chapter 2.

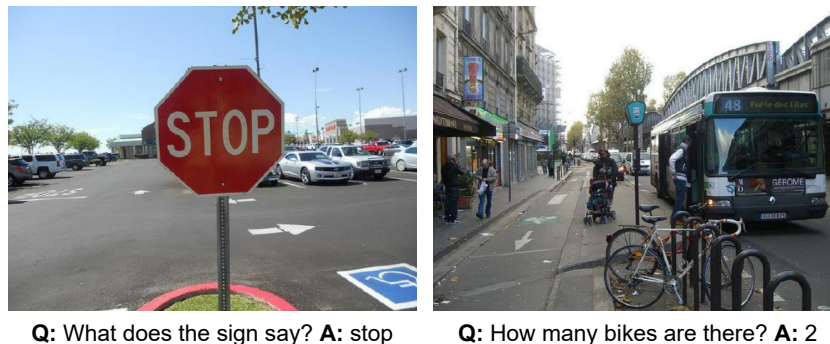


Figure 1.2: Two question-answer pairs from the VQA data set (Antol et al., 2015) along with the corresponding images.

In the following, we provide an additional example application, where bilinear pooling is used to combine the features of two different modalities, namely image and text.

Multi-modal bilinear pooling for visual question answering. The task of Visual Question Answering (VQA) is to provide an answer to a question that relates to an image (see Fig. 1.2 for an example). The answers to be predicted can be multiple choice or open-ended, meaning that they need to be provided in the form of free text. Models solving the VQA task well could be used for instance to describe images to visually impaired people.

Answers in the VQA benchmark (Antol et al., 2015) take the form of short descriptive words instead of full sentences. They may contain only a single word, and many answers such as numbers, colors and objects appear frequently. Therefore, Fukui et al. (2016) model the VQA problem as a classification task, where the input to the model is an image and a question, and the output is the most probable answer among a provided list of frequently occurring answers in the data set.

The challenge in modelling this task is to combine the feature representations of two different modalities (image and text) for further processing. These representations then need to be classified into a large number of categories, e.g., Fukui et al. (2016) use 3000 of the most common answers as classes. They introduce the technique of *multi-modal compact bilinear pooling* to combine the features of both modalities in a low-dimensional representation that can be projected onto a 3000-dimensional class representation under a reasonable computational budget. This feature representation is an instance of random feature approximations for polynomial kernels. We therefore introduce multi-modal bilinear pooling in the following.

Let $f_{\text{img}} : \mathcal{I} \times \mathcal{L} \rightarrow \mathbb{R}^{d_1}$ be an image feature extraction function as introduced before for bilinear pooling. To convert a given question into a numerical represen-

tation, we further define a function $f_{\text{txt}} : \mathcal{T} \rightarrow \mathbb{R}^{d_2}$ that converts a text $t \in \mathcal{T}$ into a d_2 -dimensional feature vector. For instance, this function could be a recurrent neural network or a Transformer (Vaswani et al., 2017).

We then extend Eq. (1.4) to a multi-modal context as follows:

$$\text{MulBil}(f_{\text{img}}, f_{\text{txt}}, \mathbf{A}, t) := \sum_{l \in \mathcal{L}} f_{\text{img}}(\mathbf{A}, l) \otimes f_{\text{txt}}(t) \in \mathbb{R}^{d_1 \cdot d_2} \quad (1.5)$$

Thus, the vectorized outer product is applied to the features of two different modalities instead of a feature map with itself, giving rise to a multi-modal feature vector.

This way of combining the features of different modalities provides a richer representation than taking their element-wise sum/product or simply concatenating them, since we consider multiplicative interactions between all *pairs* of elements. As for standard bilinear pooling, the caveat is that this feature representation is high-dimensional if d_1 and/or d_2 are large. E.g., $d_1 = d_2 = 2048$ in Fukui et al. (2016) lead to a subsequent classification layer with $2048^2 \cdot 3000 \approx 12.6$ billion parameters occupying around 48 gigabytes of memory when using 32 bits precision.

Similar to Gao et al. (2016), Fukui et al. (2016) propose to employ TensorSketch to obtain a low-dimensional sketch of $f_{\text{img}}(\mathbf{A}, l) \otimes f_{\text{txt}}(t)$ without ever constructing this vector explicitly. The sketching methods for polynomial kernels discussed in this thesis such as TensorSketch can be applied here despite the fact that the outer product \otimes in Eq. (1.5) is not applied to a vector with itself. This is because random feature maps for polynomial kernels are generally constructed by combining independent linear projections of individual inputs as we will show in Chapter 2. In this way, Fukui et al. (2016) obtain a 16,000-dimensional representation, which corresponds to around 0.4% of the original number of elements that would be obtained through Eq. (1.5). The storage of the classification layer is thus reduced to around 183 megabytes rendering the training of the model feasible. Efficiency improvements of random features for polynomial kernels that we discuss in this thesis could yield further computational savings for applications of bilinear pooling.

The applications discussed so far are connected to the use of polynomial kernels. We now provide an example application for the exponential dot product kernel.

1.3.2 Dot Product Attention for Transformers

Transformers (Vaswani et al., 2017) are neural network architectures that have been improving the state-of-the-art performance in numerous machine learning benchmarks over the last years. Examples include natural language processing (NLP) tasks like language generation (Brown et al., 2020) and machine translation (Banar et al., 2020), but also other tasks like image classification (Liu et al., 2021). Transformers gained particular popularity in NLP as they can be pretrained on large amounts of unstructured text available on the web, and be fine-tuned to the task at hand on a regular GPU affordable to the larger public (Devlin et al., 2019).

These architectures heavily rely on the attention mechanism (Vaswani et al., 2017) that models pairwise dependencies between the elements of an input sequence. E.g., in machine translation, the elements of an input sequence would be the tokens (word fragments like word stems or syllables) of a sentence to be translated. However, modelling pairwise interactions inside an input sequence introduces a quadratic time and memory complexity in the number L of elements in the sequence, making it impossible to train and fine-tune such models for large $L \gg 512$. This is also because transformers commonly use multiple attention modules, e.g., BERT-large (Devlin et al., 2019) uses a total of 24 layers equipped with 16 attention heads leading to a total number of 384 attention modules. This computational bottleneck makes it difficult to capture long-range dependencies between different tokens. Such dependencies need to be captured in tasks like text summarization, where an input sequence could consist of thousands of tokens contained in a long text.

The computational bottleneck for the modelling of long input sequences has been addressed in numerous works (e.g. Xiong et al., 2021; Choromanski et al., 2021; Zaheer et al., 2020). We provide a brief introduction to the kernelized viewpoint here for which random feature approximations have been proposed (Choromanski et al., 2021). This viewpoint is of particular interest to this thesis as it shows that the attention matrix is a normalized Gram matrix of a dot product kernel that can be approximated through random features, enabling a linear instead of a quadratic time and memory complexity with respect to L .

Dot-product attention. The dot product self-attention mechanism (Vaswani et al., 2017) can be described as follows. Consider an input sequence of objects $\mathbf{x}_1, \dots, \mathbf{x}_L$ such as text tokens or image fragments, where each object $\{\mathbf{x}_i\}_{i=1}^L$ is represented numerically by a d_{emb} -dimensional feature vector, such that it can be processed by a neural network. We summarize the sequence in matrix form as $\mathbf{X} := (\mathbf{x}_1, \dots, \mathbf{x}_L)^\top \in \mathbb{R}^{L \times d_{\text{emb}}}$.

The self-attention mechanism is a *learnable* transformation that converts the input sequence \mathbf{X} into an output sequence $\mathbf{Z} := (\mathbf{z}_1, \dots, \mathbf{z}_L)^\top \in \mathbb{R}^{L \times d}$, where the output dimension d is determined by a set of trainable weight matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_{\text{emb}} \times d}$. We define a set of linear projections $\mathbf{Q} := \mathbf{X}\mathbf{W}_Q, \mathbf{K} := \mathbf{X}\mathbf{W}_K, \mathbf{V} := \mathbf{X}\mathbf{W}_V$ and summarize the self-attention transformation in matrix-form as follows:

$$\mathbf{Z} = \text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) := \mathbf{D}^{-1} \mathbf{A} \mathbf{V}$$

$$\text{with } \mathbf{A} := \exp(\mathbf{Q}\mathbf{K}^\top / \sqrt{d}) \quad \text{and} \quad \mathbf{D} := \text{diag}\left(\sum_{i=1}^L A_{1i}, \dots, \sum_{i=1}^L A_{Li}\right)$$

Each element of $\mathbf{D}^{-1} \mathbf{A}$ hence corresponds to the evaluation of the *softmax*³ func-

³Softmax is defined as a function $f : \mathbb{R}^L \rightarrow \mathbb{R}^L$ with $(f(\mathbf{x}))_j = \exp(x_j) / \sum_{j=1}^L \exp(x_j)$.

tion:

$$\alpha_{i,j} := (\mathbf{D}^{-1}\mathbf{A})_{i,j} = \frac{\exp((\mathbf{W}_Q\mathbf{x}_i)^\top(\mathbf{W}_K\mathbf{x}_j)/\sqrt{d})}{\sum_{j=1}^L \exp((\mathbf{W}_Q\mathbf{x}_i)^\top(\mathbf{W}_K\mathbf{x}_j)/\sqrt{d})}, \quad (1.6)$$

where \mathbf{D}^{-1} takes the role of normalizing the rows of \mathbf{A} . $\alpha_{i,j} \in [0, 1]$ for $i, j = 1, \dots, L$ are called *attention weights*. This means that each individual output $\{\mathbf{z}_i\}_{i=1}^L$ is computed as $\mathbf{z}_i = \sum_{j=1}^L (\mathbf{W}_V\mathbf{x}_j)\alpha_{i,j}$, and is hence a positively weighted average of the vectors $\{\mathbf{W}_V\mathbf{x}_j\}_{j=1}^L$.

The learnable attention weights $\alpha_{i,j}$ should be large for related embeddings and small for unrelated ones. For instance, in NLP tasks, these embeddings could encode semantic information, and a weighted average of semantically related tokens in an input sequence could yield a contextualized representation avoiding ambiguities. To give a concrete example, the word ‘‘Jaguar’’ can either relate to an animal or a car manufacturer and its numerical embedding may therefore be ambiguous. Obtaining a novel numerical representation by averaging out its embedding vector with the one for the word ‘‘car’’ that appears in the same input sequence would disambiguate its numerical representation. The attention weight between the embeddings of ‘‘Jaguar’’ and ‘‘car’’ should thus be large.

Kernelized view. We can see \mathbf{A} as a Gram matrix consisting of evaluations of an exponential dot product kernel $k(\mathbf{x}, \mathbf{y}) = \exp(\mathbf{x}^\top\mathbf{y})$ for some $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Setting $(\mathbf{A})_{i,j} = k(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d}, \mathbf{W}_K\mathbf{x}_j/\sqrt{d})$ shows this correspondence. We will show in the following how random feature maps can accelerate the computation of the attention mechanism explained before.

There exist random feature approximations for the exponential dot product kernel (Kar and Karnick, 2012; Choromanski et al., 2021) and we derive novel random feature maps for it in Chapter 5. We can thus define the approximate kernel $\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x})^\top\Phi(\mathbf{y})$ for some feature map $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ such that $\mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})] = k(\mathbf{x}, \mathbf{y})$. Substituting \hat{k} into Eq. (1.6) yields the approximate attention weights

$$\hat{\alpha}_{ij} = \frac{\Phi(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d})^\top\Phi(\mathbf{W}_K\mathbf{x}_j/\sqrt{d})}{\sum_{j=1}^L \Phi(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d})^\top\Phi(\mathbf{W}_K\mathbf{x}_j/\sqrt{d})}.$$

We can now efficiently compute each $\{\mathbf{z}_i\}_{i=1}^L$ as

$$\begin{aligned} \mathbf{z}_i &= \sum_{j=1}^L (\mathbf{W}_V\mathbf{x}_j)\hat{\alpha}_{ij} = \frac{\sum_{j=1}^L (\mathbf{W}_V\mathbf{x}_j)\Phi(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d})^\top\Phi(\mathbf{W}_K\mathbf{x}_j/\sqrt{d})}{\Phi(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d})^\top \sum_{j=1}^L \Phi(\mathbf{W}_K\mathbf{x}_j/\sqrt{d})} \\ &= \frac{\left(\sum_{j=1}^L (\mathbf{W}_V\mathbf{x}_j)\Phi(\mathbf{W}_K\mathbf{x}_j/\sqrt{d})^\top\right) \Phi(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d})}{\Phi(\mathbf{W}_Q\mathbf{x}_i/\sqrt{d})^\top \left(\sum_{j=1}^L \Phi(\mathbf{W}_K\mathbf{x}_j/\sqrt{d})\right)}, \end{aligned} \quad (1.7)$$

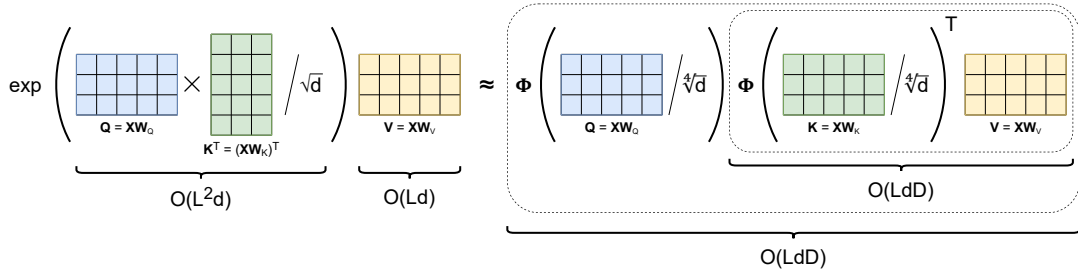


Figure 1.3: Exploiting the associativity of the matrix product to enable linear time complexity when computing random feature attention. The original unnormalized attention matrix $\exp(\mathbf{Q}\mathbf{K}^\top/\sqrt{d})$ (left) is factorized into the approximate matrix $\Phi(\mathbf{Q}/\sqrt[4]{d})\Phi(\mathbf{K}/\sqrt[4]{d})^\top$ (right) that can be efficiently multiplied with \mathbf{V} .

where the terms in parenthesis in Eq. (1.7) only need to be computed once in $\mathcal{O}(LdD)$ time as the terms are independent of i . Computing all $\{z_i\}_{i=1}^L$ then takes $\mathcal{O}(LdD)$ again. Fig. 1.3 summarizes this computation in matrix form. We thus achieve $\mathcal{O}(LdD)$ instead of $\mathcal{O}(L^2d)$ computational complexity, where an advantage is achieved with random features if $D < L$.

A similar random feature approximation has been proposed to sample efficiently from the softmax distribution for classification problems with a large number of classes (Rawat et al., 2019).

The efficiency gains in Transformer architectures brought by random feature approximations could have a significant impact on computational as well as energy savings in the future. In fact, it was shown in Strubell et al. (2019) that training a simple and commonly used BERT-base (Devlin et al., 2019) language model on a GPU cloud platform creates the same amount of CO_2 emissions as a flight from New York to San Francisco for one passenger, while the price paid to the cloud platform lies between 3,700 and 12,500 US dollars. These amounts are orders of magnitude larger when including the necessary hyperparameter tuning and architecture search. Due to the widespread use of such models, efficient random feature approximations of the exponential kernel could thus contribute to more sustainable machine learning infrastructures in this domain.

Having motivated random feature approximations for dot product kernels through a set of example *applications*, we describe in the remainder of this chapter how such approximations *naturally occur* as the result of a physical process. The Optical Processing Unit is an optical hardware that makes use of this phenomenon to compute random features efficiently, both in terms of speed and power consumption. This device could potentially yield additional energy savings and working with it inspired the line of research that we pursue in this thesis.

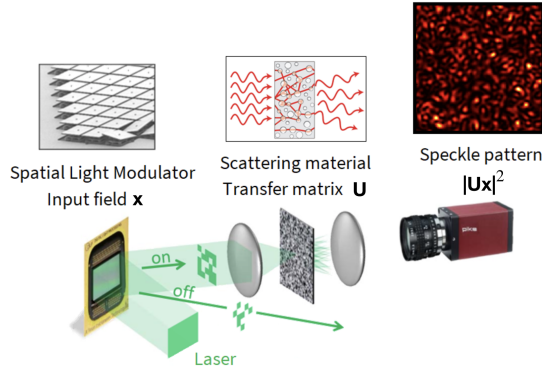


Figure 1.4: Experimental setup of the Optical Processing Unit (modified with permission from [Saade et al. \(2016\)](#)). The data vector is encoded in the coherent light from a laser using a DMD. Light then goes through a scattering medium and a speckle pattern is measured by a camera.

1.4 The Optical Processing Unit

The *Optical Processing Unit (OPU)* ([Saade et al., 2016](#); [Ohana et al., 2020](#)) is an optical hardware that performs random projections literally at the speed of light, since these projections are the result of an optical scattering process. The principle of the random projections performed by the OPU is based on the use of a semi-transparent “scattering” material that refracts the laser light passing through it. As discussed in [Liutkus et al. \(2014\)](#), light going through such a scattering medium follows many complex paths that depend on the refractive index inhomogeneities at random positions. For a fixed scattering medium, the resulting process is still linear, deterministic, and reproducible making the features obtained by the OPU useful for machine learning.

[Fig. 1.4](#) shows a schematic drawing of the physical apparatus of the OPU. We transform an arbitrary data point $\mathbf{x} \in \mathbb{R}^d$ into an analog representation by filtering a laser light beam at different positions using a *digital micromirror device (DMD)*. This encoded light then passes through the heterogeneous medium, implicitly performing a random matrix multiplication, where the random matrix is inherently determined by the scattering medium. The scattered light results in a so-called *speckle pattern* that is recorded by a camera, where the light intensity at each measuring point is the result of a superposition of the light beams filtered by the DMD. This superposition is mathematically modelled by a sum of the components of \mathbf{x} weighted by random coefficients. Measuring the intensity of light at an arbitrary pixel location recorded by the camera induces a non-linear transformation of this sum. D arbitrary pixel values taken together then constitute the random feature map

$$\Phi(\mathbf{x}) = |\mathbf{U}\mathbf{x}|^2/\sqrt{D}, \quad \text{with } \mathbf{U} \in \mathbb{C}^{D \times d} \quad \text{and} \quad (\mathbf{U})_{i,j} \stackrel{i.i.d.}{\sim} \mathcal{CN}(0,1),$$

where $|\cdot|^2$ is the squared absolute value taken element-wise.

We showed in our work (Ohana et al., 2020) that the random projections computed and recorded by the OPU correspond to the random features of a kernel that is closely related to a degree-2 polynomial kernel. This is a remarkable result because it shows that the random features obtained by the aforementioned light scattering process are useful in machine learning applications where polynomial kernels are used.

We derive this result and further extensions in Chapter 6, where we also compare the energy consumption of the OPU against a recent GPU. Since the OPU has a lower power consumption (30 Watts) than a typical GPU (200-300 Watts) while potentially achieving faster projection times, these findings are interesting both from an ecological as well as from a performance perspective. Moreover, recent work (Gupta et al., 2019) has shown that the OPU can also be used to obtain purely *linear* random projections offering a wide range of additional applications in random feature approximations and randomized linear algebra.

Chapter 2

Random Feature Approximations of Shift-Invariant Kernels and Polynomial Kernels

In the previous chapter, we motivated the use of random feature approximations through a set of example applications, while treating random features as black-box functions. The purpose of this chapter is to provide a literature overview on random feature methods with a focus on random features for polynomial kernels, since these form the basis for the random feature methods developed in the subsequent chapters of this thesis.

At the same time, we cover random features for shift-invariant kernels, as they are the most widely used random feature approximations (see [Liu et al. \(2020\)](#) for a recent overview). Random feature constructions for shift-invariant kernels and polynomial kernels are conceptually different, but [Pennington et al. \(2015\)](#) established a connection between the two. As we compare against their method as well as against random Fourier features for the approximation of the Gaussian kernel in the empirical evaluations carried out for this thesis, we consider it important to discuss both lines of research in this chapter.

Random features were originally introduced as *random Fourier features* for shift-invariant kernels by [Rahimi and Recht \(2007\)](#), and several other works followed later on that studied random feature approximations for the polynomial kernel, in particular [Kar and Karnick \(2012\)](#); [Pham and Pagh \(2013\)](#); [Hamid et al. \(2014\)](#); [Avron et al. \(2014\)](#); [Pennington et al. \(2015\)](#); [Ahle et al. \(2020\)](#).

We therefore cover random Fourier features for shift-invariant kernels first.

2.1 Random Fourier Features for Shift-Invariant Kernels

In this section, we give a general introduction to random features for shift-invariant kernels. In addition, we discuss the subclass of isotropic kernels for which these random features take on a special form.

A shift-invariant (or stationary) kernel $k(\mathbf{x}, \mathbf{y}) \in \mathbb{C}$ is a positive definite function of $\boldsymbol{\tau} = \mathbf{x} - \mathbf{y} \in \mathbb{R}^d$, making it invariant to translations of the inputs. We will thus overload the notation and write $k(\boldsymbol{\tau})$ instead.

Shift-invariant kernels are characterized through Bochner's theorem (see [Stein, 1999](#), Chapter 2.5):

Theorem 2.1.1 (Bochner's theorem). *A complex-valued shift-invariant kernel $k(\boldsymbol{\tau})$ on \mathbb{R}^d is positive definite if and only if it can be written in the following form:*

$$k(\boldsymbol{\tau}) = \int_{\mathbb{R}^d} \exp(i\boldsymbol{\omega}^\top \boldsymbol{\tau}) \mu(\boldsymbol{\omega}) d\boldsymbol{\omega}, \quad (2.1)$$

where μ is a non-negative finite measure and $i := \sqrt{-1}$.

The construction of random Fourier features relies entirely on the above theorem and we introduce them in the following.

Random Fourier features. *Random Fourier features* were introduced in the seminal work by [Rahimi and Recht \(2007\)](#) who exploit the fact that the measure μ is proportional to a probability density p . To see this, let $\sigma^2 := k(\mathbf{0}) = \int_{\mathbb{R}^d} \mu(\boldsymbol{\omega}) d\boldsymbol{\omega} \geq 0$. We then obtain a density $p(\boldsymbol{\omega}) := \mu(\boldsymbol{\omega})/\sigma^2$ because $p(\boldsymbol{\omega}) \geq 0$ and $\int_{\mathbb{R}^d} p(\boldsymbol{\omega}) d\boldsymbol{\omega} = 1$ hold. W.l.o.g., we assume $\sigma^2 = 1$ from now on. ([Stein, 1999](#), Chapter 2.5) call $p(\boldsymbol{\omega})$ the *spectral density*, which can be found for a given kernel through the following inversion formula:

$$p(\boldsymbol{\omega}) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \exp(-i\boldsymbol{\omega}^\top \boldsymbol{\tau}) k(\boldsymbol{\tau}) d\boldsymbol{\tau} \quad (2.2)$$

Therefore, there is a one-to-one correspondence between spectral densities and shift-invariant kernels.

Given a properly scaled shift-invariant kernel and its spectral density, we can rewrite [Eq. \(2.1\)](#) as:

$$\begin{aligned} k(\boldsymbol{\tau}) &= \int_{\mathbb{R}^d} \exp(i\boldsymbol{\omega}^\top \boldsymbol{\tau}) p(\boldsymbol{\omega}) d\boldsymbol{\omega} = \mathbb{E}_{\boldsymbol{\omega} \sim p(\boldsymbol{\omega})} [\exp(i\boldsymbol{\omega}^\top \boldsymbol{\tau})] \\ &= \mathbb{E}_{\boldsymbol{\omega} \sim p(\boldsymbol{\omega})} \left[\exp(i\boldsymbol{\omega}^\top \mathbf{x}) \overline{\exp(i\boldsymbol{\omega}^\top \mathbf{y})} \right]. \end{aligned} \quad (2.3)$$

We can then obtain a Monte-Carlo estimate $\hat{k}_c(\boldsymbol{\tau}) := \Phi_c(\mathbf{x})^\top \overline{\Phi_c(\mathbf{y})}$ of $k(\boldsymbol{\tau})$ with

$$\Phi_c(\mathbf{x}) = (\exp(i\boldsymbol{\omega}_1^\top \mathbf{x}), \dots, \exp(i\boldsymbol{\omega}_D^\top \mathbf{x}))^\top / \sqrt{D} \in \mathbb{C}^D,$$

where $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_D \stackrel{i.i.d.}{\sim} p(\boldsymbol{\omega})$. $\Phi_{\mathcal{C}} : \mathbb{R}^d \rightarrow \mathbb{C}^D$ is called a *random feature map* and we use the subscript \mathcal{C} to emphasize the use of complex-valued basis functions. From now on, we drop the subscript $\boldsymbol{\omega} \sim p(\boldsymbol{\omega})$ of the expectation for brevity.

The dimension D of the random feature map determines a trade-off between computational cost and approximation quality. We briefly derive the variances of random Fourier features to show this effect.

Variance of the complex estimator. The variance of $\hat{k}_{\mathcal{C}}(\boldsymbol{\tau})$ is:

$$\begin{aligned} \mathbb{V}[\hat{k}_{\mathcal{C}}(\boldsymbol{\tau})] &= \frac{1}{D^2} \sum_{\ell=1}^D \mathbb{V}[\exp(i\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})] = \frac{1}{D^2} \sum_{\ell=1}^D \mathbb{E} \left[|\exp(i\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})|^2 \right] - |\mathbb{E}[\exp(i\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})]|^2 \\ &= \frac{1}{D^2} \sum_{\ell=1}^D \mathbb{E} [\cos^2(\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau}) + \sin^2(\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})] - k(\boldsymbol{\tau})^2 = \frac{1}{D} (1 - k(\boldsymbol{\tau})^2) \end{aligned}$$

As the $\{\boldsymbol{\omega}_{\ell}\}_{\ell=1}^D$ are i.i.d., we have $\mathbb{V}[\hat{k}(\boldsymbol{\tau})] \propto 1/D$ giving better kernel approximations for larger D . At the same time, the input $\boldsymbol{\tau}$ also affects the variance. E.g., if $\boldsymbol{\tau} = \mathbf{0}$, the variance is always zero regardless of D .

When the kernel of interest is real-valued, we can derive a real-valued feature map as follows.

Real-valued feature maps. For kernels $k(\boldsymbol{\tau}) \in \mathbb{R}$, we can rewrite Eq. (2.3) as:

$$\begin{aligned} k_{\mathcal{R}}(\boldsymbol{\tau}) &= \mathbb{E} [\exp(i\boldsymbol{\omega}^{\top} \boldsymbol{\tau})] = \mathbb{E} [\cos(\boldsymbol{\omega}^{\top} \boldsymbol{\tau})] + i\mathbb{E} [\sin(\boldsymbol{\omega}^{\top} \boldsymbol{\tau})] \\ &= \mathbb{E} [\cos(\boldsymbol{\omega}^{\top} \boldsymbol{\tau})] = \mathbb{E} [\cos(\boldsymbol{\omega}^{\top} \boldsymbol{x} - \boldsymbol{\omega}^{\top} \boldsymbol{y})] \\ &= \mathbb{E} [\cos(\boldsymbol{\omega}^{\top} \boldsymbol{x}) \cos(\boldsymbol{\omega}^{\top} \boldsymbol{y}) + \sin(\boldsymbol{\omega}^{\top} \boldsymbol{x}) \sin(\boldsymbol{\omega}^{\top} \boldsymbol{y})], \end{aligned}$$

where we use the trigonometric identity $\cos(\alpha - \beta) = \cos(\alpha)\cos(\beta) + \sin(\alpha)\sin(\beta)$. Thus, we can use an alternative random feature map $\Phi_{\mathcal{R}} : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$ for real-valued kernels $\hat{k}_{\mathcal{R}}(\boldsymbol{\tau}) := \Phi_{\mathcal{R}}(\boldsymbol{x})^{\top} \Phi_{\mathcal{R}}(\boldsymbol{y})$ with

$$\Phi_{\mathcal{R}}(\boldsymbol{x}) := [\cos(\boldsymbol{\omega}_1^{\top} \boldsymbol{x}), \sin(\boldsymbol{\omega}_1^{\top} \boldsymbol{x}), \dots, \cos(\boldsymbol{\omega}_D^{\top} \boldsymbol{x}), \sin(\boldsymbol{\omega}_D^{\top} \boldsymbol{x})]^{\top} / \sqrt{D} \in \mathbb{R}^{2D}. \quad (2.4)$$

Variance of the real estimator. In this case the variance is

$$\begin{aligned} \mathbb{V}[\hat{k}_{\mathcal{R}}(\boldsymbol{\tau})] &= \frac{1}{D^2} \sum_{\ell=1}^D \mathbb{V}[\cos(\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})] = \frac{1}{D^2} \sum_{\ell=1}^D \mathbb{E} [\cos^2(\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})] - \mathbb{E} [\cos(\boldsymbol{\omega}_{\ell}^{\top} \boldsymbol{\tau})]^2 \\ &= \frac{1}{D} \left(\frac{1}{2} + \frac{1}{2} k(2\boldsymbol{\tau}) - k(\boldsymbol{\tau})^2 \right) \leq \mathbb{V}[\hat{k}_{\mathcal{C}}(\boldsymbol{\tau})], \end{aligned}$$

where we use the trigonometric identity $\cos(\alpha^2) = \frac{1}{2} + \frac{1}{2} \cos(2\alpha)$.

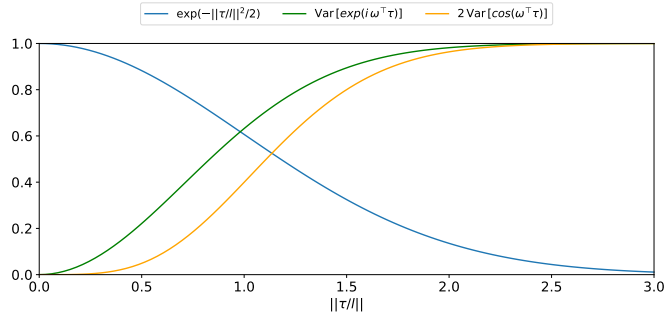


Figure 2.1: The Gaussian kernel (blue) and the variance of complex (green) and real-valued (yellow) kernel estimators. We double the variance of the real-valued estimator to compare feature maps of equal dimensions.

Comparing the variance of real and complex kernel estimates. Although we obtain a variance reduction w.r.t. D when using the real-valued estimator, this comes at the expense of doubling the dimension of the feature map compared to the complex case. We therefore compare $\mathbb{V}[\hat{k}_{\mathcal{R}}(\boldsymbol{\tau})]$ and $\mathbb{V}[\hat{k}_{\mathcal{C}}(\boldsymbol{\tau})]$ on the example of the Gaussian kernel defined as $k(\boldsymbol{\tau}) := \exp(-\|\boldsymbol{\tau}\|^2/(2l^2))$ with length scale parameter $l > 0$ in the following.

Fig. 2.1 shows this comparison across different values of $\|\boldsymbol{\tau}/l\|$. In order to yield a fair comparison, we double the variance of the real-valued estimator since it requires $2D$ features (orange line). It can be seen that the real-valued estimator has a lower variance than its complex-valued counterpart (green line) for all values $\|\boldsymbol{\tau}/l\|$. Thus, for the Gaussian kernel, real-valued random Fourier features perform better despite doubling the feature map dimension. Furthermore, we can generally observe that variances increase as $\|\boldsymbol{\tau}/l\|$ increases, which indicates that more random features are required when l is short compared to the scaling of the data itself. This happens in particular when the function to be modelled oscillates quickly, i.e., it uses high frequencies. We will empirically study this phenomenon in Chapter 5.

Although complex-valued feature maps perform worse in terms of the approximation variance in this example, we will show in Chapter 3 that the opposite can be the case for polynomial kernel approximations.

For further statistical guarantees on random Fourier features, we refer the reader to the following related works. Exponential concentration bounds are derived in the original work by [Rahimi and Recht \(2007\)](#) and are tightened by [Sutherland and Schneider \(2015\)](#). Statistical guarantees for kernel ridge regression are derived in ([Avron et al., 2017](#); [Rudi and Rosasco, 2017](#)).

2.1.1 Spherical Random Features for Isotropic Kernels

A special class of random Fourier features are Spherical Random Features (SRF) introduced in [Pennington et al. \(2015\)](#) for isotropic kernels. They are of particular

importance to this thesis, as they will provide us with a baseline method when approximating high-degree polynomial kernels in Chapter 5.

We begin by characterizing isotropic kernels.

Isotropic kernels. Isotropic kernels are shift-invariant (or stationary) kernels of the form $k(\mathbf{x}, \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|)$ for some $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, where $\|\cdot\|$ denotes the Euclidean norm. This means that these kernels are also rotation-invariant. To see this, consider a rotation matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$ with $\mathbf{R}^\top \mathbf{R} = \mathbf{I}_d$ being applied to both \mathbf{x} and \mathbf{y} . Then we have:

$$\|\mathbf{R}\mathbf{x} - \mathbf{R}\mathbf{y}\|^2 = (\mathbf{R}\mathbf{x})^\top (\mathbf{R}\mathbf{x}) + (\mathbf{R}\mathbf{y})^\top (\mathbf{R}\mathbf{y}) - 2(\mathbf{R}\mathbf{x})^\top (\mathbf{R}\mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2.$$

We now show that all isotropic kernels are real-valued. Recall from the previous chapter that $\overline{k(\mathbf{x}, \mathbf{y})} = k(\mathbf{y}, \mathbf{x})$ holds for any kernel k . Therefore, we can write

$$\overline{k(\mathbf{x}, \mathbf{y})} = k(\mathbf{y}, \mathbf{x}) = k(\|\mathbf{y} - \mathbf{x}\|) = k(\|\mathbf{x} - \mathbf{y}\|) = k(\mathbf{x}, \mathbf{y}),$$

which implies that $\text{Im}\{k(\boldsymbol{\tau})\} = 0$, where $\text{Im}\{\cdot\}$ denotes the imaginary part of a complex number.

Popular examples of isotropic kernels include Gaussian and Matérn kernels (see [Rasmussen and Williams, 2006](#), Chapter 4.2). In the following, we derive random feature approximations for the specific case of isotropic kernels. For this purpose, we rely on a transformation of Bochner's theorem.

Basis functions for isotropic kernels. Bochner's theorem (2.1) can be reformulated to provide a specific spectral representation for the class of isotropic kernels. Let $k(\|\boldsymbol{\tau}\|)$ be an isotropic kernel with spectral density $p(\boldsymbol{\omega})$. As shown in [Stein \(1999, Chapter 2.10\)](#), the spectral density of isotropic kernels is radial, i.e., it only depends on $\|\boldsymbol{\omega}\|$. This allows to change the integration over the density $p(\boldsymbol{\omega})$ in Eq. (2.1) into an integration over the one-dimensional density $\hat{p}(\|\boldsymbol{\omega}\|)$, i.e., a density over the norm of $\boldsymbol{\omega}$, as follows:

$$k(\|\boldsymbol{\tau}\|) = \int_0^\infty d\|\boldsymbol{\omega}\| \hat{p}(\|\boldsymbol{\omega}\|) \underbrace{\Gamma\left(\frac{d}{2}\right) \left(\frac{2}{\|\boldsymbol{\omega}\| \|\boldsymbol{\tau}\|}\right)^{\frac{d}{2}-1} J_{\frac{d}{2}-1}(\|\boldsymbol{\omega}\| \|\boldsymbol{\tau}\|)}_{:= \Lambda_d(\|\boldsymbol{\omega}\| \|\boldsymbol{\tau}\|)}, \quad (2.5)$$

where J_ν with $\nu \in \mathbb{R}$ is the Bessel function of the first kind ([Abramowitz and Stegun, 1974](#), Chapter 10) and $\Gamma(\cdot)$ is the Gamma function.

The one-dimensional density $\hat{p}(\|\boldsymbol{\omega}\|)$ can be derived from $p(\boldsymbol{\omega})$ by using an integration over spherical coordinates and the fact that both densities need to integrate to one:

$$\int_{\mathbb{R}^d} p(\boldsymbol{\omega}) d\boldsymbol{\omega} = \int_0^\infty \underbrace{p(\boldsymbol{\omega}) \|\boldsymbol{\omega}\|^{d-1} \frac{(2\pi)^{d/2}}{\Gamma(d/2)}}_{:= \hat{p}(\|\boldsymbol{\omega}\|)} d\|\boldsymbol{\omega}\| = 1 \quad (2.6)$$

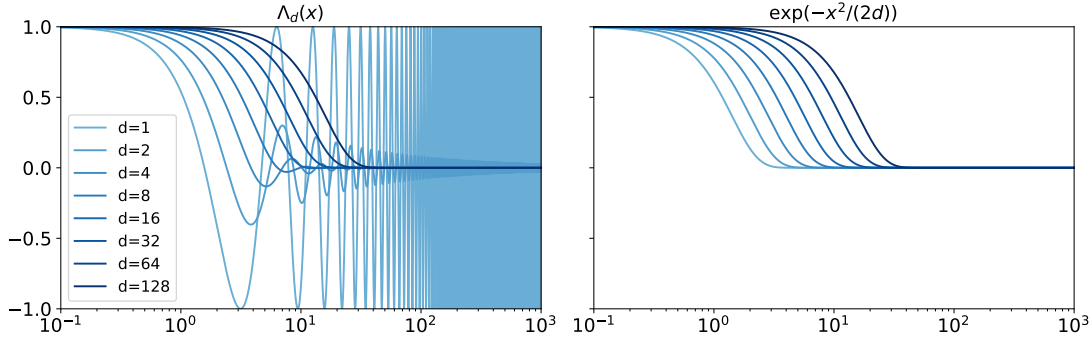


Figure 2.2: Comparing $\Lambda_d(x)$ on the left to $\exp(-x^2/(2d))$ on the right for different values of d . The functions are indistinguishable for large enough d .

In fact, the density $\hat{p}(\|\boldsymbol{\omega}\|)$ depends on the dimension d of $\boldsymbol{\omega}$ and only exists if k is a valid kernel, which is the case when $\int_0^\infty \|\boldsymbol{\tau}\|^{d-1} |k(\|\boldsymbol{\tau}\|)| d\|\boldsymbol{\tau}\| < \infty$ holds (see [Stein, 1999](#), Chapter 2.10).

As noted in [Genton \(2002\)](#), the basis functions $\Lambda_d(\cdot)$ in [Eq. \(2.5\)](#) become less expressive and have an increasing lower bound for increasing d . Therefore, the class of valid isotropic kernels becomes more restrictive as d increases and certain isotropic kernels are positive definite only up to a certain input dimension. [Fig. 2.2](#) shows the basis functions $\Lambda_d(\cdot)$ in [Eq. \(2.5\)](#) for different values of d . It is interesting to note that $\Lambda_d(\|\boldsymbol{\omega}\| \|\boldsymbol{\tau}\|)$ approaches $\exp(-(\|\boldsymbol{\omega}\| \|\boldsymbol{\tau}\|)^2/(2d))$ as d increases. This has the immediate consequence that any isotropic kernel that is positive definite for *all* d , can be modeled as a scale-mixture of Gaussians. This idea is formalized in Schoenberg's theorem:

Theorem 2.1.2 (Schoenberg 1938). *A continuous function $f : [0, \infty) \rightarrow \mathbb{R}$ is positive definite and radial on \mathbb{R}^d if and only if it is of the form $f(r) = \int_0^\infty \exp(-r^2 t^2) d\mu(t)$, where μ is a finite non-negative Borel measure on $[0, \infty)$.*

Having learned that any positive definite isotropic kernel is directly defined by the one-dimensional density $\hat{p}(\|\boldsymbol{\omega}\|)$ allows us to adapt the sampling procedure of random Fourier features for the class of isotropic kernels. We refer to this novel sampling procedure as Spherical Random Features and explain it in the following.

Spherical Random Features (SRF). As noted earlier, the original spectral density $p(\boldsymbol{\omega})$ of isotropic kernels is radial and thus depends only on $\|\boldsymbol{\omega}\|$. This density must therefore be equal for all $\boldsymbol{\omega}$ with an equal radius $r := \|\boldsymbol{\omega}\| > 0$. This observation allows to sample from $p(\boldsymbol{\omega})$ by drawing uniform samples from the unit sphere $\mathcal{S}^{d-1} := \{\boldsymbol{\omega} \in \mathbb{R}^d : \|\boldsymbol{\omega}\| = 1\}$ and to scale them using samples from the norm density $\hat{p}(\|\boldsymbol{\omega}\|)$. [Algorithm 1](#) summarizes this procedure.

[Pennington et al. \(2015\)](#) refer to the resulting random features as *Spherical Random Features*. Since isotropic kernels are always real-valued, we can replace

Algorithm 1: Spherical Random Features (SRF)

Result: An SRF map $\Phi_C(\mathbf{x})$ Draw D i.i.d. samples $\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_D$ from $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$;Draw D i.i.d. samples $\omega_1, \omega_2, \dots, \omega_D$ from $\hat{p}(\|\boldsymbol{\omega}\|)$;Compute $\Phi_C(\mathbf{x}) = \left(\exp\left(i \frac{\omega_1}{\|\mathbf{s}_1\|} \mathbf{s}_1^\top \mathbf{x}\right), \dots, \exp\left(i \frac{\omega_D}{\|\mathbf{s}_D\|} \mathbf{s}_D^\top \mathbf{x}\right) \right) / \sqrt{D} \in \mathbb{C}^D$;

the complex-valued basis functions in [Algorithm 1](#) with sines and cosines to obtain a real-valued feature map as we did in [Eq. \(2.4\)](#).

A crucial advantage of [Algorithm 1](#) is that we can substitute i.i.d. samples from the unit-sphere \mathcal{S}^{d-1} with orthogonal quasi-Monte-Carlo samples as described by [Yu et al. \(2016\)](#) to approximate the integral in [Eq. \(2.1\)](#) more efficiently, i.e., requiring less random features for the same kernel approximation error. An in-depth discussion when orthogonal samples are superior approximators for isotropic kernels can be found in [Choromanski et al. \(2018\)](#).

Spherical Random Features can also be used to approximate polynomial kernels on the unit sphere. However, such random feature approximations are biased. We therefore delay their introduction to the end of this chapter and cover unbiased approximations first.

2.2 Random Features for Polynomial Kernels

We study here random feature approximations of *polynomial kernels*, defined as

$$k(\mathbf{x}, \mathbf{y}) = (\gamma \mathbf{x}^\top \mathbf{y} + \nu)^p, \quad (2.7)$$

where $\gamma, \nu \geq 0$ and $p \in \mathbb{N}$. We call such random features *polynomial sketches*¹. Since polynomial kernels are not shift-invariant, widely known random Fourier features ([Rahimi and Recht, 2007](#)) cannot be applied for their unbiased approximation. Polynomial sketches are a fundamentally different approach, and can be understood as implicit linear random projections, or sketches ([Woodruff, 2014](#)), of the explicit high dimensional feature maps of polynomial kernels.

For simplicity, we focus on *homogeneous* polynomial kernels of the form

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y})^p, \quad (2.8)$$

i.e., we set $\gamma = 1, \nu = 0$ in [Eq. \(2.7\)](#). All other cases can be reduced to the homogeneous case, by appending $\sqrt{\nu}$ to the input vectors and rescaling them, i.e., by setting $\tilde{\mathbf{x}} := (\sqrt{\gamma} \mathbf{x}^\top, \sqrt{\nu})^\top \in \mathbb{R}^{d+1}$ and $\tilde{\mathbf{y}} := (\sqrt{\gamma} \mathbf{y}^\top, \sqrt{\nu})^\top \in \mathbb{R}^{d+1}$, we have

$$(\gamma \mathbf{x}^\top \mathbf{y} + \nu)^p = (\tilde{\mathbf{x}}^\top \tilde{\mathbf{y}})^p.$$

¹We exclude Spherical Random Features ([Pennington et al., 2015](#)) from this terminology as they are not implicit high-dimensional sketches.

In this way, polynomial sketches for the homogeneous case can also be applied to the inhomogeneous case.

There exists an explicit feature map for polynomial kernels. To see this, we introduce some new notation. Recall from the previous chapter that we define the vectorized outer product between two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ as $\mathbf{a} \otimes \mathbf{b} := \text{vec}(\mathbf{a}\mathbf{b}^\top) \in \mathbb{R}^{d^2}$. We further define

$$\mathbf{a}^{(k)} := \underbrace{\mathbf{a} \otimes \cdots \otimes \mathbf{a}}_{k \text{ terms}} \in \mathbb{R}^{d^k} \quad (2.9)$$

for some $k \in \mathbb{N}$. Then the homogeneous polynomial kernel can be written as (Scholkopf and Smola, 2002, Proposition 2.1):

$$(\mathbf{x}^\top \mathbf{y})^p = (\mathbf{x}^{(p)})^\top \mathbf{y}^{(p)} \quad (2.10)$$

Thus, $\mathbf{x}^{(p)}$ and $\mathbf{y}^{(p)}$ are the exact feature maps of the input vectors \mathbf{x} and \mathbf{y} , respectively. This exact feature expansion leads to d^p dimensional vectors² that are infeasible to construct when d and/or p are moderately large. This justifies the need for randomized approximations of the polynomial kernel.

There are several different approaches for the randomized approximation of polynomial kernels in the literature. We identify four different construction types of such polynomial sketches that we call: (A) Product-Sketches, (B) Convolutions of CountSketches (a.k.a. TensorSketch), (C) Kronecker Product of SRHT sketches (a.k.a. TensorSRHT), and (D) Spherical Random Features for polynomial kernels. They all have in common that they can be constructed *without* forming the exact d^p -dimensional feature map of the polynomial kernel.

There exist further algorithms that can be used on top of the aforementioned random feature approximations to improve their efficiency. E.g., Hamid et al. (2014) propose to first create a high-dimensional projection that is down-projected in a second step and Ahle et al. (2020) propose a hierarchical random feature construction. Such algorithms are not the focus of this work and we will not discuss them further. The random features discussed in this thesis should be seen as base sketches that can be used inside such algorithms. The goal of this thesis is to improve the efficiency of these base sketches.

The contributions in this thesis are all made regarding Product-Sketches that were the first ones to be proposed in the literature. We cover them in the following and introduce random feature types (B)-(D) in the remainder of this chapter to give a context over the literature.

²By grouping up equal terms, we can reduce the dimensionality to $\binom{d+p-1}{p}$, which still leads to unrealistic dimensional feature vectors as soon as d is large. For example, working with MNIST (Lecun et al., 1998) images of size 28x28 ($d = 784$) leads to 307,720 features for $p = 2$ and to 80,622,640 features for $p = 3$.

2.2.1 (A) Product-Sketches

We first study polynomial sketches proposed by Kar and Karnick (2012), which are also discussed in Hamid et al. (2014).

Let $D \in \mathbb{N}$ be the number of random features, and $p \in \mathbb{N}$ be the degree of the polynomial kernel (2.8). Suppose we generate $p \times D$ i.i.d. random vectors

$$\mathbf{w}_{i,\ell} \in \mathbb{R}^d \quad \text{satisfying} \quad \mathbb{E}[\mathbf{w}_{i,\ell} \mathbf{w}_{i,\ell}^\top] = \mathbf{I}_d, \quad i \in \{1, \dots, p\}, \quad \ell \in \{1, \dots, D\}, \quad (2.11)$$

where $\mathbf{I}_d \in \mathbb{R}^{d \times d}$ denotes the identity matrix.

Then we define a random feature map as

$$\Phi(\mathbf{x}) := \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{w}_{i,1}^\top \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{w}_{i,D}^\top \mathbf{x} \right) \right]^\top \in \mathbb{R}^D. \quad (2.12)$$

We call the specific form of the feature map in Eq. (2.12) a *Product-Sketch* since they are an element-wise product of individual linear random projections. The resulting approximation of the polynomial kernel (2.8) is given by

$$\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x})^\top \Phi(\mathbf{y}), \quad (2.13)$$

which is unbiased, as the expectation with respect to the random vectors (2.11) gives

$$\mathbb{E} [\Phi(\mathbf{x})^\top \Phi(\mathbf{y})] = \frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p \mathbf{x}^\top \mathbb{E}[\mathbf{w}_{i,\ell} \mathbf{w}_{i,\ell}^\top] \mathbf{y} = (\mathbf{x}^\top \mathbf{y})^p.$$

Kar and Karnick (2012) suggest to define random vectors in (2.11) using the Rademacher distribution: each element of $\mathbf{w}_{i,\ell}$ is independently drawn from $\{-1, 1\}$ with equal probability. Hamid et al. (2014) suggest to use Gaussian as well as structured Rademacher weights based on the subsampled randomized Hadamard transform (SRHT) (Tropp, 2011) instead. We study in Chapter 3 how the distribution of the random vectors affects the quality of the kernel approximation.

Implicit sketching of high-dimensional features. The random feature map (2.12) can be interpreted as a linear sketch (projection) of an explicit high-dimensional feature vector for the polynomial kernel. To describe this, consider the case $D = 1$ and let $\mathbf{w}_i = (w_{i,1}, \dots, w_{i,d})^\top \in \mathbb{R}^d$ with $i = 1, \dots, p$ be i.i.d. random vectors satisfying $\mathbb{E}[\mathbf{w}_i \mathbf{w}_i^\top] = \mathbf{I}_d$. Then, the random feature map $\Phi(\mathbf{x})$ (which is one dimensional in this case) for $\mathbf{x} := (x_1, \dots, x_d)^\top$ is given by

$$\Phi(\mathbf{x}) = \prod_{i=1}^p \mathbf{w}_i^\top \mathbf{x} = \prod_{i=1}^p \sum_{j=1}^d w_{i,j} x_j = \sum_{j_1=1, \dots, j_p=1}^d w_{1,j_1} x_{j_1} \cdots w_{p,j_p} x_{j_p} \quad (2.14)$$

$$= (\mathbf{w}_1 \otimes \cdots \otimes \mathbf{w}_p)^\top \mathbf{x}^{(p)}, \quad (2.15)$$

where $\mathbf{x}^{(p)} \in \mathbb{R}^{d^p}$ is defined as in Eq. (2.9). Therefore, for $D = 1$, the approximate kernel is given as

$$\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}) = (\mathbf{w}_1 \otimes \cdots \otimes \mathbf{w}_p)^\top \mathbf{x}^{(p)} \cdot (\mathbf{w}_1 \otimes \cdots \otimes \mathbf{w}_p)^\top \mathbf{y}^{(p)} \quad (2.16)$$

As shown in Eq. (2.10), $\mathbf{x}^{(p)}$ and $\mathbf{y}^{(p)}$ are the exact feature maps of the polynomial kernel for the input vectors \mathbf{x} and \mathbf{y} , respectively. Thus, the random feature map $\Phi(\mathbf{x}) \in \mathbb{R}$ in (2.15) is a projection of the exact feature map $\mathbf{x}^{(p)} \in \mathbb{R}^{d^p}$ onto \mathbb{R} .

Similarly, if $D > 1$, the random feature map $\Phi(\mathbf{x}) \in \mathbb{R}^D$ in (2.12) can be interpreted as a projection of the exact feature map $\mathbf{x}^{(p)}$ onto \mathbb{R}^D . A remarkable point of this random feature map is that it can be obtained without constructing the exact feature vector $\mathbf{x}^{(p)}$. Indeed, the computational complexity of constructing the random feature map $\Phi(\mathbf{x})$ is $\mathcal{O}(pdD)$, while the exact feature map $\mathbf{x}^{(p)}$ requires $\mathcal{O}(d^p)$.

The variances of the sketches introduced in this section have not been explored prior to our work. We will derive detailed variance formulas for several Product-Sketches in Chapters 3 and 4. However, there exist exponential concentration bounds similar to the ones derived in Rahimi and Recht (2007) for random Fourier features. These can be found in Kar and Karnick (2012); Hamid et al. (2014). We want to emphasize however, that these bounds serve to provide theoretical *guarantees* for the number of random features required for a certain approximation error. They should not be used to compare the statistical efficiency of these sketches. I.e., Hamid et al. (2014) obtain a tighter bound when using Gaussian weights in Eq. (2.11) than the one for Rademacher weights obtained by Kar and Karnick (2012). We will show in Chapter 3 however, that Rademacher sketches, and in particular structured ones yield much lower approximation variances than Gaussian sketches, which is also reflected in their empirical performances in downstream tasks in Chapter 4.

2.2.2 (B) Convolutions of CountSketches a.k.a. TensorSketch

TensorSketch is a structured polynomial sketch that was introduced by Pham and Pagh (2013). It can still be considered a state-of-the-art polynomial sketch in terms of kernel approximation quality and computation speed. Hence, it was used in Gao et al. (2016) and Fukui et al. (2016) to accelerate bilinear pooling operations as explained in the previous chapter. Since *TensorSketch* is a convolution of *CountSketches* (Charikar et al., 2002), we introduce *CountSketches* first.

CountSketch (Charikar et al., 2002). *CountSketches* were originally invented to maintain approximate counts of high-frequency elements in large data streams where it is infeasible to store the entire vocabulary of unique elements (Charikar et al., 2002). Weinberger et al. (2009); Woodruff (2014) show that *CountSketches*

can also be used as a dimensionality reduction technique leading to unbiased estimates of the linear kernel.

CountSketches are usually defined in terms of hash functions (Charikar et al., 2002; Pham and Pagh, 2013). However, in order to keep a consistent framework throughout this literature review, we define CountSketches in terms of randomized linear projections.

Let $\mathbf{C} \in \{-1, 0, 1\}^{D \times d}$ be a *sparse* projection matrix with elements

$$(\mathbf{C})_{i,j} = \delta(h_j = i)s_j, \quad \text{with } i \in \{0, \dots, D-1\}, j \in \{1, \dots, d\}, \quad (2.17)$$

where $\mathbf{h} := (h_1, \dots, h_d)^\top \in \mathbb{R}^d$ has elements independently and uniformly drawn from $\{0, \dots, D-1\}$ and $\mathbf{s} := (s_1, \dots, s_d)^\top \in \{-1, 1\}^d$ is a vector with independent Rademacher samples. $\delta(h_j = i) = 1$ if $h_j = i$ and zero otherwise. The rows of \mathbf{C} are indexed starting from zero as this will simplify our notation later on.

\mathbf{C} is sparse because it has only one non-zero entry per column defined through \mathbf{h} . Thus, it has only d non-zero entries in total allowing to compute the CountSketch $\mathbf{C}\mathbf{x} \in \mathbb{R}^D$ for a given $\mathbf{x} \in \mathbb{R}^d$ in $\mathcal{O}(d)$ time and memory. In fact, it is sufficient to only store \mathbf{h} and \mathbf{s} in the place of \mathbf{C} . The i -th element of $\mathbf{C}\mathbf{x}$ can then be efficiently computed as $(\mathbf{C}\mathbf{x})_i = \sum_{j:h_j=i} x_j s_j$.

We now show that $\mathbf{C}\mathbf{x}^\top \mathbf{C}\mathbf{y}$ is an unbiased estimate of the linear kernel:

$$\begin{aligned} \mathbb{E}[\mathbf{C}\mathbf{x}^\top \mathbf{C}\mathbf{y}] &= \mathbb{E} \left[\sum_{i=0}^{D-1} \left(\sum_{j=1}^d \delta(h_j = i) s_j x_j \right) \left(\sum_{j'=1}^d \delta(h_{j'} = i) s_{j'} y_{j'} \right) \right] \quad (2.18) \\ &= \sum_{i=0}^{D-1} \sum_{j=1}^d \sum_{j'=1}^d \mathbb{E}[\delta(h_j = i) \delta(h_{j'} = i)] \mathbb{E}[s_j s_{j'}] x_j y_{j'} \\ &= \sum_{i=0}^{D-1} \sum_{j=1}^d \mathbb{E}[\delta(h_j = i) \delta(h_j = i)] x_j y_j = \sum_{i=0}^{D-1} \sum_{j=1}^d \frac{1}{D} x_j x_j = \mathbf{x}^\top \mathbf{y}, \end{aligned}$$

where we use the fact that $\mathbb{E}[s_j s_{j'}] = 1$ if $j = j'$ and $\mathbb{E}[s_j s_{j'}] = \mathbb{E}[s_j] \mathbb{E}[s_{j'}] = 0$ otherwise. Moreover, $\mathbb{E}[\delta(h_j = i)] = 1/D \cdot 1 + 0 = 1/D$.

In the following, we describe an efficient way to compute a CountSketch of $\mathbf{x}^{(p)}$ known as *TensorSketch* (Pham and Pagh, 2013) without ever constructing $\mathbf{x}^{(p)}$ itself. TensorSketch is thus a polynomial sketch and it is formed through a convolution of p individual CountSketches.

TensorSketch. The convolution operation is equivalent to a product in the frequency domain. We therefore apply the discrete Fourier transform (DFT) to a CountSketch. Let $k \in \{0, \dots, D-1\}$ be the index over the elements of $\mathbf{C}\mathbf{x}$. The

k -th element of its DFT is computed as:

$$\begin{aligned} (\text{DFT}[\mathbf{C}\mathbf{x}])_k &= \sum_{i=0}^{D-1} \exp\left(-i\frac{2\pi}{D}ki\right) \mathbf{c}_i^\top \mathbf{x} = \sum_{j=1}^d \sum_{i=0}^{D-1} \exp\left(-i\frac{2\pi}{D}ki\right) \delta(h_j = i) s_j x_j \\ &= \sum_{j=1}^d \exp\left(-i\frac{2\pi}{D}kh_j\right) s_j x_j. \end{aligned}$$

Now we take the element-wise product of the Fourier transform of two independent CountSketches $\mathbf{C}_1\mathbf{x}$ and $\mathbf{C}_2\mathbf{x}$ defined through $\mathbf{h}_1, \mathbf{h}_2$ and $\mathbf{s}_1, \mathbf{s}_2$, respectively. The k -th element of this product then yields:

$$\begin{aligned} &(\text{DFT}[\mathbf{C}_1\mathbf{x}])_k \cdot (\text{DFT}[\mathbf{C}_2\mathbf{x}])_k \\ &= \left(\sum_{j=1}^d \exp\left(-i\frac{2\pi}{D}kh_{1,j}\right) s_{1,j}x_j \right) \left(\sum_{j'=1}^d \exp\left(-i\frac{2\pi}{D}kh_{2,j'}\right) s_{2,j'}x_{j'} \right) \\ &= \sum_{j=1}^d \sum_{j'=1}^d \exp\left(-i\frac{2\pi}{D}k(h_{1,j} + h_{2,j'})\right) s_{1,j}s_{2,j'}x_jx_{j'} \\ &= \sum_{j=1}^d \sum_{j'=1}^d \exp\left(-i\frac{2\pi}{D}k((h_{1,j} + h_{2,j'}) \bmod D)\right) s_{1,j}s_{2,j'}x_jx_{j'} = (\text{DFT}[\mathbf{C}_{(1*2)}\mathbf{x}^{(2)}])_k, \end{aligned}$$

where $\mathbf{C}_{(1*2)} \in \{-1, 0, 1\}^{D \times d^2}$ is a novel CountSketch with \mathbf{h} and \mathbf{s} in Eq. (2.17) being replaced by

$$\mathbf{h}^{(1+2)} := \text{vec}(\mathbf{h}_1 + \mathbf{h}_2^\top) \bmod D \in \mathbb{R}^{d^2} \quad \text{and} \quad \mathbf{s}^{(1\otimes 2)} := \mathbf{s}_1 \otimes \mathbf{s}_2 \in \mathbb{R}^{d^2}.$$

Note that the arithmetic modulus arises because $f(x) = \exp(-i\frac{2\pi}{D}x)$ has a period of D .

In Appendix A.1, we show that $\mathbf{h}^{(1+2)}$ has elements drawn uniformly from $\{0, \dots, D-1\}$. Moreover, $\mathbf{s}^{(1\otimes 2)}$ is a vector with elements drawn from the Rademacher distribution because the product of two independent Rademacher variables is again Rademacher distributed. This means that the convolution of two independent CountSketches of \mathbf{x} yields a CountSketch of $\mathbf{x}^{(2)}$. Extending the derivation above to p independent CountSketches $\{\mathbf{C}_i\mathbf{x}\}_{i=1}^p$ gives a new CountSketch $\mathbf{C}_{(1*\dots*p)}$ of $\mathbf{x}^{(p)}$ that can be computed as:

$$\Phi(\mathbf{x}) := \text{DFT}^{-1} [\text{DFT}[\mathbf{C}_1\mathbf{x}] \cdots \text{DFT}[\mathbf{C}_p\mathbf{x}]] = \mathbf{C}_{(1*\dots*p)}\mathbf{x}^{(p)}, \quad (2.19)$$

where DFT^{-1} is the inverse discrete Fourier transform. We call $\Phi(\mathbf{x})$ (2.19) a TensorSketch of degree p . Using the Fast Fourier Transform, the total computational cost for obtaining $\Phi(\mathbf{x})$ is $\mathcal{O}(p(d+D \log D))$. This implies that we implicitly obtain a linear sketch in d^p dimensions without ever having to compute $\mathbf{x}^{(p)}$ explicitly.

Using Eq. (2.18), we can show that the feature map $\Phi(\mathbf{x})$ leads to an unbiased estimate $\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$ of the polynomial kernel:

$$\mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})] = \mathbf{x}^{(p)\top} \mathbf{y}^{(p)} = (\mathbf{x}^\top \mathbf{y})^p$$

A caveat of TensorSketch is that the elements of the outer product $\mathbf{s}^{(1 \otimes \dots \otimes p)} = \mathbf{s}_1 \otimes \dots \otimes \mathbf{s}_p$ are statistically dependent despite having pairwise covariances of zero. This is sufficient for TensorSketch to yield unbiased kernel estimates, but it influences higher order moments making it difficult to study its approximation variances for $p > 1$ as shown in Avron et al. (2014, Lemma 2 of the longer version) (the case $p = 1$ is covered in Weinberger et al. (2009)). There is thus no closed-form variance formula available in the literature³.

Avron et al. (2014) derive only the following upper bound on the variance:

Theorem 2.2.1 (Avron et al. (2014), Lemma 2). *For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the variance of the TensorSketch kernel estimate $\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$ with Φ defined as in Eq. (2.19) is upper-bounded as follows:*

$$\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})] \leq \frac{1}{D} (2 + 3^p) \|\mathbf{x}\|^{2p} \|\mathbf{y}\|^{2p}$$

We will show in Chapter 4 that TensorSketch performs better than Rademacher Product-Sketches empirically, in spite of the upper bound in Theorem 2.2.1 being larger than the one we obtain for the Rademacher and Gaussian Product-Sketches in Chapter 3. To the best of our knowledge, there are also no exponential concentration bounds for TensorSketch available in the literature, as it is the case for the Product-Sketches discussed in the previous section. We hypothesize that this is due to the fact that the rows of the matrix $\mathbf{C}_{(1^* \dots * p)}$ are not sampled independently obstructing a simple application of well-known concentration inequalities. Nonetheless, theoretical guarantees for downstream tasks like kernel ridge regression and kernel PCA are given in Avron et al. (2014); Ahle et al. (2020).

2.2.3 (C) Kronecker Product of SRHT Sketches a.k.a. TensorSRHT

In this section, we introduce structured polynomial sketches known as *TensorSRHT* (*Tensor Subsampled Randomized Hadamard Transform*). Tropp (2011) studied TensorSRHT for $p = 1$ (linear case) called simply SRHT, and Hamid et al. (2014); Ahle et al. (2020) extended it⁴ to arbitrary polynomial degrees p . Ahle et al. (2020) introduced the TensorSRHT sketch that we present in this section. The

³The original TensorSketch paper (Pham and Pagh, 2013) contains a variance formula, but makes the simplifying assumption that TensorSketch has the same variance properties as CountSketch (Charikar et al., 2002) applied to tensorized inputs.

⁴The sketches proposed by Hamid et al. (2014) and Ahle et al. (2020) are similar but not exactly the same. Hamid et al. (2014) uses $p \times B$ independent linear SRHT sketches (see Tropp,

prefix “Tensor” is used here to refer to the fact we form the Kronecker or Tensor product of p independent sketches, which yields a new sketch that is an implicit linear projection of the tensorized input $\mathbf{x}^{(p)}$, as it was also the case for the other sketches described before. TensorSRHT is of great relevance to this thesis because we develop a closely related sketch in Chapter 3.

In order to define TensorSRHT, we need to cover some preliminaries first. More specifically, we need to define Kronecker products and Hadamard matrices.

Kronecker Product. We start by extending the notion of the outer product in Eq. (2.9) applied to vectors being applied to matrices instead. For this purpose, we define the Kronecker product of two matrices $\mathbf{A} \in \mathbb{R}^{m_1 \times n_1}$, $\mathbf{B} \in \mathbb{R}^{m_2 \times n_2}$ as

$$\mathbf{A} \otimes \mathbf{B} := \begin{bmatrix} (\mathbf{A})_{11}\mathbf{B} & \dots & (\mathbf{A})_{1n_1}\mathbf{B} \\ \vdots & \ddots & \vdots \\ (\mathbf{A})_{m_11}\mathbf{B} & \dots & (\mathbf{A})_{m_1n_1}\mathbf{B} \end{bmatrix} \in \mathbb{R}^{(m_1 \cdot m_2) \times (n_1 \cdot n_2)}.$$

The Kronecker product has the *mixed-product property* defined as follows. For matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ of such sizes that one can compute \mathbf{AC} and \mathbf{BD} , it holds that

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}). \quad (2.20)$$

Note especially that we have $(\mathbf{A} \otimes \mathbf{B})(\mathbf{x} \otimes \mathbf{y}) = \mathbf{Ax} \otimes \mathbf{By}$ for two vectors \mathbf{x} and \mathbf{y} . We make use of this property for an efficient computation of TensorSRHT later.

We define structured Hadamard matrices next.

Hadamard matrices. Let $n := 2^m$ with $m \in \mathbb{N}$, and define $\mathbf{H}_n \in \{1, -1\}^{n \times n}$ to be the unnormalized Hadamard matrix, which is recursively defined as

$$\mathbf{H}_{2n} := \begin{bmatrix} \mathbf{H}_n & \mathbf{H}_n \\ \mathbf{H}_n & -\mathbf{H}_n \end{bmatrix}, \quad \text{with} \quad \mathbf{H}_2 := \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}. \quad (2.21)$$

From now on, we always use $\mathbf{H}_d \in \{1, -1\}^{d \times d}$ with d being the dimensionality of input vectors, assuming $d = 2^m$ for some $m \in \mathbb{N}$. If $d \neq 2^m$ for any $m \in \mathbb{N}$, we pad the input vectors with zeros until their dimensionality becomes 2^m for some $m \in \mathbb{N}$. For $i = 1, \dots, d$, let $\mathbf{h}_i \in \{1, -1\}^d$ be the i -th column of \mathbf{H}_d , i.e.,

$$\mathbf{H}_d = (\mathbf{h}_1, \dots, \mathbf{h}_d) \in \{1, -1\}^{d \times d}.$$

\mathbf{H}_d has very useful properties that we exploit for the construction of several polynomial sketches in this thesis. The first useful property is $\mathbf{H}_d \mathbf{H}_d^\top = \mathbf{H}_d^\top \mathbf{H}_d =$

2011), where $B := \lceil \frac{D}{d} \rceil$ is the number of SRHT blocks per degree. The elements of these sketches are then shuffled across degrees and blocks and the blocks are multiplied elementwise over p . Ahle et al. (2020) compute only p independent sketches and subsample from their tensor product instead.

$d\mathbf{I}_d$, which implies that distinct columns (and rows) of \mathbf{H}_d are orthogonal to each other, i.e., $\mathbf{h}_i^\top \mathbf{h}_j = 0$ for $i \neq j$. Another useful property is that $\mathbf{H}_{d^2} = \mathbf{H}_d \otimes \mathbf{H}_d$ holds, which will be needed for the construction of the TensorSRHT sketch in the following. Lastly, the recursive definition of \mathbf{H}_d gives rise to the Fast Walsh-Hadamard Transform (FWHT) (Fino and Algazi, 1976) that multiplies \mathbf{H}_d with a vector $\mathbf{a} \in \mathbb{R}^d$ in $\mathcal{O}(d \log d)$ instead of $\mathcal{O}(d^2)$ time while \mathbf{H}_d does not need to be stored in memory.

We now move on to define SRHT, which is a linear random projection that forms the basis for the construction of TensorSRHT.

Subsampled Randomized Hadamard Transform (SRHT). SRHT (Tropp, 2011) is a linear random projection that is particularly fast due to the structured nature of \mathbf{H}_d .

Let $\mathbf{D} := \text{diag}(\mathbf{w}) \in \mathbb{R}^{d \times d}$ be a diagonal matrix whose diagonal entries $\mathbf{w} := (w_1, \dots, w_d)^\top \in \{1, -1\}^d$ are i.i.d. Rademacher random variables. Further consider a random permutation of the indices $\pi : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$, and let $\pi(1), \dots, \pi(d)$ be the permuted indices. We define a sampling matrix $\mathbf{P}_\pi := (\mathbf{e}_{\pi(1)}, \dots, \mathbf{e}_{\pi(D)})^\top \in \mathbb{R}^{D \times d}$ with $D \leq d$, where $\mathbf{e}_{\pi(\ell)} \in \mathbb{R}^d$ is a vector whose $\pi(\ell)$ -th element is 1 and all other elements are 0.

The SRHT is then defined as

$$\Phi(\mathbf{x}) := \mathbf{P}_\pi \mathbf{H}_d \mathbf{D} \mathbf{x} / \sqrt{D} = \mathbf{P}_\pi (\mathbf{H}_d (\mathbf{w} \odot \mathbf{x})) / \sqrt{D} \in \mathbb{R}^D, \quad (2.22)$$

where \odot denotes the element-wise (Hadamard) product.

The parenthesis on the right-hand-side of Eq. (2.22) indicate the efficient order of computation to leverage the Fast Walsh-Hadamard Transform. $\Phi(\mathbf{x})$ can thus be computed in $\mathcal{O}(d \log d)$ instead of $\mathcal{O}(d^2)$ time, and using only $\mathcal{O}(d)$ memory since \mathbf{H}_d does not need to be stored.

The sketch (2.22) leads to an unbiased estimate $\hat{k}(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$ because

$$\begin{aligned} \mathbb{E}[\Phi(\mathbf{x})^\top \Phi(\mathbf{y})] &= \frac{1}{D} \sum_{\ell=1}^D \mathbb{E}[\mathbf{h}_{\pi(\ell)}^\top (\mathbf{w} \odot \mathbf{x}) \cdot \mathbf{h}_{\pi(\ell)}^\top (\mathbf{w} \odot \mathbf{y})] \\ &= \frac{1}{D} \sum_{\ell=1}^D \sum_{j=1}^d \sum_{j'=1}^d \mathbb{E}[h_{\pi(\ell),j} h_{\pi(\ell),j'}] \mathbb{E}[w_j w_{j'}] x_j y_{j'} \\ &= \frac{1}{D} \sum_{\ell=1}^D \sum_{j=1}^d \underbrace{h_{\pi(\ell),j}^2 \mathbb{E}[w_j^2]}_{=1} x_j y_j = \mathbf{x}^\top \mathbf{y}, \end{aligned}$$

where the expectation is taken w.r.t. both the random permutation π as well as \mathbf{w} . We have $h_{\pi(\ell),j}^2 = 1$ for all $\{h_{\pi(\ell),j}\}_{j=1}^d$ because the elements of $\mathbf{h}_{\pi(\ell)}$ are either -1 or $+1$.

The goal of TensorSRHT is to apply SRHT to the tensorized input $\mathbf{x}^{(p)}$ to obtain an unbiased estimate of the polynomial kernel without ever constructing $\mathbf{x}^{(p)}$ explicitly.

TensorSRHT. The sketch defined by Ahle et al. (2020, Definition 15) exploits the mixed-product property (2.20) of the Kronecker product. This property allows for an efficient construction of a d^p -dimensional SRHT sketch, while working only with random projections in \mathbb{R}^d .

Let $\{\mathbf{D}_i\}_{i=1}^p$ be independently sampled diagonal matrices of size $\mathbb{R}^{d \times d}$ with their diagonal elements being defined through the independent Rademacher vectors $\{\mathbf{w}_i\}_{i=1}^p$ of size \mathbb{R}^d , respectively. We further define a random sampling matrix $\mathbf{P}_\pi := (\mathbf{e}_{\pi(1)}, \dots, \mathbf{e}_{\pi(D)})^\top \in \mathbb{R}^{D \times d^p}$ with $D \leq d^p$ as before.

We now show how to transform p individual SRHT sketches $\{\mathbf{H}\mathbf{D}_i\mathbf{x}\}_{i=1}^p$ into a single SRHT sketch of an input point $\mathbf{x}^{(p)}$. For this purpose, we apply the mixed-product property (2.20) twice and use $\mathbf{H}_d \otimes \mathbf{H}_d = \mathbf{H}_{d^2}$ to show the following key relationship used in TensorSRHT:

$$\Phi(\mathbf{x}) := \mathbf{P}_\pi(\mathbf{H}_d\mathbf{D}_1\mathbf{x} \otimes \dots \otimes \mathbf{H}_d\mathbf{D}_p\mathbf{x})/\sqrt{D} \quad (2.23)$$

$$\begin{aligned} &= \mathbf{P}_\pi(\mathbf{H}_d\mathbf{D}_1 \otimes \dots \otimes \mathbf{H}_d\mathbf{D}_p)(\mathbf{x} \otimes \dots \otimes \mathbf{x})/\sqrt{D} \\ &= \mathbf{P}_\pi(\mathbf{H}_d \otimes \dots \otimes \mathbf{H}_d)(\mathbf{D}_1 \otimes \dots \otimes \mathbf{D}_p)\mathbf{x}^{(p)}/\sqrt{D} \\ &= \mathbf{P}_\pi\mathbf{H}_{d^p}(\mathbf{D}_1 \otimes \dots \otimes \mathbf{D}_p)\mathbf{x}^{(p)}/\sqrt{D} \end{aligned} \quad (2.24)$$

$(\mathbf{D}_1 \otimes \dots \otimes \mathbf{D}_p)$ in Eq. (2.24) is a diagonal matrix of size $\mathbb{R}^{d^p \times d^p}$ with its diagonal elements being $\mathbf{w}_1 \otimes \dots \otimes \mathbf{w}_p \in \mathbb{R}^{d^p}$. As noted for TensorSketch before, these elements are Rademacher distributed with zero covariance. $\Phi(\mathbf{x})$ (2.24) thus corresponds to SRHT applied to $\mathbf{x}^{(p)}$.

The equivalent formulation in Eq. (2.23) allows us to efficiently compute this sketch.

Efficient computation. We can efficiently compute $\Phi(\mathbf{x})$ by first computing $\{\mathbf{H}\mathbf{D}_i\mathbf{x}\}_{i=1}^p$ via the FWHT as shown in Eq. (2.23). Then we sample D index tuples $\{(i_{1,\ell}, \dots, i_{p,\ell})\}_{\ell=1}^D$ without replacement, where $i_{1,\ell}, \dots, i_{p,\ell}$ are uniformly sampled from $\{1, \dots, d\}$. The tuples $\{(i_{1,\ell}, \dots, i_{p,\ell})\}_{\ell=1}^D$ can be understood as D single indices sampled without replacement from $\{1, \dots, d^p\}$.

We then compute each element $\{\Phi(\mathbf{x})_\ell\}_{\ell=1}^D$ of the final TensorSRHT sketch as

$$\Phi(\mathbf{x})_\ell = (\mathbf{H}\mathbf{D}_1\mathbf{x})_{i_{1,\ell}} \dots (\mathbf{H}\mathbf{D}_p\mathbf{x})_{i_{p,\ell}}/\sqrt{D}, \quad (2.25)$$

where $\Phi(\mathbf{x})_\ell$ corresponds to a randomly sampled element (without replacement) from the d^p -dimensional vector $\mathbf{H}_{d^p}(\mathbf{D}_1 \otimes \dots \otimes \mathbf{D}_p)\mathbf{x}^{(p)}/\sqrt{D}$.

The computational cost of TensorSRHT is thus $\mathcal{O}(p(d \log d + D))$, since computing $\{\mathbf{H}\mathbf{D}_i\mathbf{x}\}_{i=1}^p$ costs $\mathcal{O}(pd \log d)$ and taking the products in Eq. (2.25) for $\ell = 1, \dots, D$ costs an additional $\mathcal{O}(pD)$. Recall that TensorSketch discussed in the section before has a cost of $\mathcal{O}(p(D \log D + d))$. This makes TensorSRHT faster than TensorSketch if $D > d$.

Statistical guarantees. As for TensorSketch, the variance of the TensorSRHT sketch discussed here is not available in closed form in the literature as soon as $p > 1$ (a result for $p = 1$ is available in Choromanski et al. (2017)). Once again, we hypothesize that this is because the entries of $(\mathbf{D}_1 \otimes \cdots \otimes \mathbf{D}_p)$ in Eq. (2.24) are statistically dependent despite having zero covariance. Moreover, the elements of $\Phi(\mathbf{x})$ (2.24) additionally have *non-zero* covariances since every row of \mathbf{H}_{d^p} is multiplied by the same matrix $(\mathbf{D}_1 \otimes \cdots \otimes \mathbf{D}_p)$.

One of the main contributions of this thesis is to derive a closely related structured sketch in Chapter 3 for which we obtain closed form variances. While Ahle et al. (2020) focus on statistical guarantees that do not take the data distribution into account, we elucidate under which conditions of the input data our structured sketches perform better than non-structured sketches in terms of variance comparisons. In particular, we show that our structured sketches yield lower variances than unstructured analogs for odd degrees of the polynomial kernel, with the additional advantage of being faster. We are not aware of any other work in the literature that theoretically proves such statistical advantages for structured polynomial sketches. Nonetheless, theoretical guarantees for TensorSRHT for downstream tasks can be found in Ahle et al. (2020).

2.2.4 (D) Spherical Random Features for Polynomial Kernels

Pennington et al. (2015) propose a fundamentally different approach to polynomial kernel approximation from the previous ones presented in this chapter. Instead of constructing an unbiased implicit sketch in d^p dimensions, the authors propose to approximate polynomial kernels through isotropic kernels and random Fourier features (see Section 2.1.1). We review this approach here as it can be considered the state-of-the-art for the approximation of *high-degree* polynomial kernels on the unit sphere. It also serves as an important baseline in our work and we found it rather difficult to implement it in practice, which is why this section provides a comprehensive introduction and discusses some limitations of the method.

The approach by Pennington et al. (2015) is motivated by the fact that polynomial kernels with data lying on the unit-sphere can be expressed as isotropic kernels. The authors consider the approximation of a polynomial kernel $k : \mathcal{S}^{d-1} \times \mathcal{S}^{d-1} \rightarrow \mathbb{R}$ being parameterized as follows:

$$k(\mathbf{x}, \mathbf{y}) = \left(1 - \frac{\|\mathbf{x} - \mathbf{y}\|^2}{a^2}\right)^p = (\nu + \gamma(\mathbf{x}^\top \mathbf{y}))^p, \quad (2.26)$$

where $\nu = 1 - 2/a^2$ and $\gamma = 2/a^2$ for some $a \geq 2$.

Eq. (2.26) shows the correspondence between a shift-invariant kernel (left-hand-side) and a polynomial kernel (right-hand-side) valid for inputs lying on the unit-sphere.

However, the condition $a \geq 2$ makes it impossible to model homogeneous polynomial kernels for which we have $\nu = 0$ and thus $a = \sqrt{2}$. Furthermore, the offset ν and the scaling term γ are inter-dependent through a common dependency on a , which limits the class of polynomial kernels being modelled.

Since Eq. (2.26) takes the form of an isotropic kernel, one could hope to approximate k with random Fourier features as shown in Section 2.1. Yet, for $\mathbf{x}, \mathbf{y} \in \mathcal{S}^{d-1}$, we have $\|\boldsymbol{\tau}\| = \|\mathbf{x} - \mathbf{y}\| = \sqrt{2 - 2\cos\theta} \in [0, 2]$, where θ is the angle between \mathbf{x} and \mathbf{y} . Therefore, the support of $\boldsymbol{\tau}$ is not defined over the entire space \mathbb{R}^d and Bochner's theorem (2.1) does not apply. Pennington et al. (2015) further prove that the support of $\boldsymbol{\tau}$ cannot be artificially extended such that the kernel function (2.26) remains positive definite. Hence, there exists no spectral density for this kernel, which obstructs the use of random Fourier features.

The authors therefore consider a *biased* approximation by finding a closely related isotropic kernel defined on \mathbb{R}^d , i.e., one for which random Fourier features can be obtained.

Finding an approximate isotropic kernel. Recall from Section 2.1 that any isotropic kernel defined on \mathbb{R}^d for some $d \in \mathbb{N}$ can be expressed using Eq. (2.5) as

$$\hat{k}(\|\boldsymbol{\tau}\|) = \int_0^\infty d\|\boldsymbol{\omega}\| \hat{p}(\|\boldsymbol{\omega}\|) \Lambda_d(\|\boldsymbol{\omega}\| \|\boldsymbol{\tau}\|), \quad (2.27)$$

where $\hat{p}(\|\boldsymbol{\omega}\|)$ is a one-dimensional density and $\Lambda_d(\cdot)$ is defined through Eq. (2.5). Pennington et al. (2015) propose to optimize $\hat{p}(\|\boldsymbol{\omega}\|)$ so as to minimize the mean squared error

$$L = \frac{1}{2} \int_0^2 d\|\boldsymbol{\tau}\| \left(k(\|\boldsymbol{\tau}\|) - \hat{k}(\|\boldsymbol{\tau}\|) \right)^2 \quad (2.28)$$

with k being defined as in Eq. (2.26). In practice, L is evaluated on a grid over $\|\boldsymbol{\tau}\| \in [0, 2]$. Here, the density is parameterized as

$$\hat{p}(\|\boldsymbol{\omega}\|) \propto \|\boldsymbol{\omega}\|^{d-1} \max \left(0, \sum_{i=1}^N c_i \left(\frac{1}{\sqrt{2}\sigma_i} \right)^d \exp(-\|\boldsymbol{\omega}\|^2 / 4\sigma_i^2) \right) \quad (2.29)$$

for some $\{c_i\}_{i=1}^N, \{\sigma_i\}_{i=1}^N$ with $c_i \in \mathbb{R}$ and $\sigma_i > 0$.

The formulation of $\hat{p}(\|\boldsymbol{\omega}\|)$ is motivated by Theorem 2.1.2 stating that any isotropic kernel defined on arbitrary domains \mathbb{R}^d can be modelled as a scale-mixture of Gaussian functions. Eq. (2.29) is obtained by choosing a scale-mixture over Gaussian distributions for $p(\boldsymbol{\omega})$ (2.2) and transforming it using Eq. (2.6). Pennington et al. (2015) further modify the resulting density by adding $\max(0, \cdot)$ and allowing negative values for $\{c_i\}_{i=1}^N$ in Eq. (2.29) for greater numerical flexibility. The normalization constant in Eq. (2.29) is found through a numerical integration on a grid over $\|\boldsymbol{\omega}\| \in \mathbb{R}_+$, which is also done for the evaluation of \hat{k} (2.27) itself. $\{c_i\}_{i=1}^N, \{\sigma_i\}_{i=1}^N$ in Eq. (2.29) are then found using a gradient-based optimizer. Once they are obtained, one can sample from $\hat{p}(\|\boldsymbol{\omega}\|)$ using inverse-transform-sampling to eventually sample Spherical Random Features using Algorithm 1.

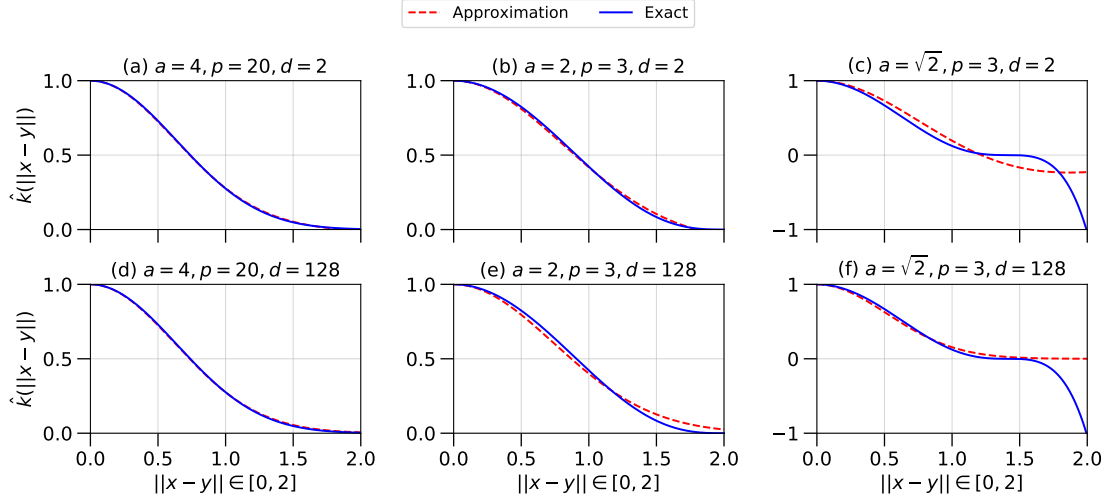


Figure 2.3: Approximating polynomial kernels on the unit-sphere through isotropic kernels. The first row shows approximations for $d = 2$ while the bottom row shows them for $d = 128$.

Limits of the approximation. Unlike the polynomial sketches presented earlier, the approximation (2.27) of k (2.26) is biased. We show some characteristic approximation functions $\hat{k}(\|\tau\|)$ (2.27) in Fig. 2.3 that are obtained through the optimization of Eq. (2.28). Note that these functions are themselves approximated using Spherical Random Features in a second step. Plots (a),(b),(d),(e) show kernel parameterizations with $a \geq 2$ as desired by Pennington et al. (2015). Plots (c) and (f) explicitly show homogeneous polynomial kernels ($a = \sqrt{2}$), and it becomes immediately obvious that they cannot be matched through \hat{k} (2.27) exposing an important limitation of this method.

Furthermore, we note that high-degree ($p = 20$) polynomial kernels are better approximated than low-degree ($p = 3$) ones confirming the findings of the authors, who argue that their method is particularly suitable for large p leading to a lower bias. They further show that Spherical Random Features can outperform previous approaches presented in Sections 2.2.1 and 2.2.2 for data lying on the unit-sphere in this case.

Although performing well for large p , we want to emphasize that the method by Pennington et al. (2015) is challenging to implement in practice. We therefore provide some insights on the numerical stability of the method that can help to simplify the implementation.

Numerical stability. The issue of numerical stability is not discussed in Pennington et al. (2015) but quite relevant. The evaluation of the Bessel function of the first kind $J_{d/2-1}(\cdot)$ as part of $\Lambda_d(\cdot)$ in Eq. (2.27) yields small values and tends to

numerically underflow for large d . The original implementation⁵ proposed by the authors solves this issue by using different series approximations for different values of d . Instead, we propose to use $\exp(-x^2/(2d))$ instead of $\Lambda_d(x)$ from $d > 128$ to simplify the implementation. This gives a reasonably good approximation of $\Lambda_d(x)$ as shown in Fig. 2.2.

⁵<https://github.com/felixyu/SRF>

Chapter 3

Analysis and Extension of Product-Sketches

In this chapter, we focus on the unbiased approximation of polynomial kernels by means of Product-Sketches that we introduced in Chapter 2. As we will learn in Chapter 5, polynomial sketches can be used to construct random feature maps for arbitrary dot product kernels. This gives the study of the efficient approximation of polynomial sketches a large importance. Therefore, we carry out an in-depth variance analysis of Product-Sketches in this chapter and propose structured and complex extensions that improve these variances.

Product-Sketches were used in [Kar and Karnick \(2012\)](#) and [Hamid et al. \(2014\)](#), who proposed to sample corresponding weight vectors from Gaussian and Rademacher distributions, and derived concentration results for both. The bound for Gaussian Product-Sketches derived in [Hamid et al. \(2014, Lemma 3.1\)](#) guarantees a given pointwise kernel approximation error for fewer random features than the bound by [Kar and Karnick \(2012\)](#) for Rademacher sketches. However, this result does not reflect the empirical performance of these two sketches as we show later in this work.

In this chapter, we carry out a variance analysis of Product-sketches and we show that Rademacher sketches achieve a variance lower bound under the assumption that the corresponding weight vectors are i.i.d. We further prove that even lower variances can be attained when lifting the i.i.d. assumption and using structured Rademacher sketches.

To the best of our knowledge, such an analysis is novel to the literature and it elucidates under which conditions of the input data one sketch should be preferred over another. This is a key contribution of this chapter, which complements existing theoretical results that rather focus on the question *how many* random features are needed to provide a certain approximation error at a given probability. Unlike our analysis, such statistical guarantees can make it difficult to compare the *average-case* efficiency of different sketches.

The second major contribution of this chapter is to derive *complex-valued* gen-

eralizations of Product-Sketches and to study under which conditions they achieve lower variances than real-valued analogs. We show that variances generally decrease considerably for positively-valued data and large polynomial degrees.

3.1 Variance Analysis for Real-Valued Product-Sketches

We begin by studying the variance properties of existing real-valued Product-Sketches, namely the Gaussian and Rademacher sketches introduced in [Kar and Karnick \(2012\)](#); [Hamid et al. \(2014\)](#).

Recall from Chapter 2 that a Product-Sketch is defined as

$$\Phi_{\mathcal{R}}(\mathbf{x}) := \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{w}_{i,1}^{\top} \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{w}_{i,D}^{\top} \mathbf{x} \right) \right]^{\top} \in \mathbb{R}^D, \quad (3.1)$$

for some weight vectors $\mathbf{w}_{i,\ell}$ with $i = 1, \dots, p$ and $\ell = 1, \dots, D$, where the subscript \mathcal{R} emphasizes the use of real-valued weight vectors here. We construct complex-valued generalizations later on.

The resulting approximation of the polynomial kernel $(\mathbf{x}^{\top} \mathbf{y})^p$ is then given by

$$\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) := \Phi_{\mathcal{R}}(\mathbf{x})^{\top} \Phi_{\mathcal{R}}(\mathbf{y}). \quad (3.2)$$

For ease of presentation, we focus on the case $D = 1$ (a single random feature), and consider p i.i.d. weight vectors

$$\mathbf{w}_i = (w_{i,1}, \dots, w_{i,d})^{\top} \in \mathbb{R}^d \quad \text{satisfying} \quad \mathbb{E}[\mathbf{w}_i] = \mathbf{0}, \quad \mathbb{E}[\mathbf{w}_i \mathbf{w}_i^{\top}] = \mathbf{I}_d, \quad (3.3)$$

with $i = 1, \dots, p$, where the elements $w_{i,1}, \dots, w_{i,d}$ are themselves i.i.d. Then the resulting approximate kernel (3.2) is given by

$$\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^p \mathbf{w}_i^{\top} \mathbf{x} \mathbf{w}_i^{\top} \mathbf{y}. \quad (3.4)$$

The variances for the case $D > 1$ can be simply obtained by dividing the variance for $D = 1$ by D , since the approximate kernel (3.2) for $D > 1$ is the average of D i.i.d. copies of the approximate kernel (3.4) for $D = 1$.

Since the polynomial sketches discussed here are unbiased, the resulting variance of the approximate kernel $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})$ in (3.2) is the *mean-square error* with respect to the true polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \mathbf{y})^p$, thus serving as a quality metric:

$$\mathbb{V}[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})] := \mathbb{E}[(\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) - \mathbb{E}[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})])^2] = \mathbb{E}[(\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) - k(\mathbf{x}, \mathbf{y}))^2],$$

where the expectation is with respect to the weight vectors \mathbf{w}_i for $i = 1, \dots, p$. Thus, the variance analysis also enables understanding how the distribution of the weight vectors affects the quality of the kernel approximation.

[Theorem 3.1.1](#) below provides a closed form expression of the variance of the approximate kernel and its lower bound. To the best of our knowledge, this result is a novel contribution to the literature. It is useful in understanding how the variance depends on the distribution of the weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_p$ and the properties of input vectors \mathbf{x} and \mathbf{y} . Since this result is a special case of a more general result presented later in [Theorem 3.3.1](#) regarding complex polynomial sketches, we omit its proof.

Theorem 3.1.1. *Let $\mathbf{x} := (x_1, \dots, x_d)^\top \in \mathbb{R}^d$ and $\mathbf{y} := (y_1, \dots, y_d)^\top \in \mathbb{R}^d$ be any input vectors. Let $\mathbf{w}_1, \dots, \mathbf{w}_p \in \mathbb{R}^d$ be i.i.d. random vectors satisfying [\(3.3\)](#), such that elements w_{i1}, \dots, w_{id} of each vector $\mathbf{w}_i = (w_{i1}, \dots, w_{id})^\top$ are themselves i.i.d. Let $\mathbf{w} = (w_1, \dots, w_d)^\top \in \mathbb{R}^d$ be a random vector independently and identically distributed as $\mathbf{w}_1, \dots, \mathbf{w}_p$. Then, for the variance of the approximate kernel [\(3.4\)](#), we have*

$$\mathbb{V} \left[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) \right] = \left(\sum_{k=1}^d \mathbb{E}[w_k^4] x_k^2 y_k^2 + \|\mathbf{x}\|^2 \|\mathbf{y}\|^2 - 3 \sum_{k=1}^d x_k^2 y_k^2 + 2(\mathbf{x}^\top \mathbf{y})^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \quad (3.5)$$

$$\geq \left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + 2 \left[(\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right] \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \quad (3.6)$$

[Eq. \(3.5\)](#) of [Theorem 3.1.1](#) shows that the distribution of the i.i.d. weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_p$ affects the variance of the approximate kernel only through the 4-th moments $\mathbb{E}[w_k^4]$, provided that the distribution satisfies the required conditions in [Theorem 3.1.1](#). In particular, the smaller these 4-th moments are, the smaller the variance becomes. This observation enables discussing the optimal choice for the distribution of the weight vectors. The lower bound in [Eq. \(3.6\)](#) is the lowest possible variance, since we have $\mathbb{E}[w_k^4] \geq (\mathbb{E}[w_k^2])^2 = 1$ by Jensen's inequality, where $\mathbb{E}[w_k^2] = 1$ follows from [Eq. \(3.3\)](#).

Let us discuss two choices for the distribution of the weight vectors in [Eq. \(3.3\)](#): Rademacher and Gaussian distributions. For the Rademacher distribution, we have $\mathbb{E}[w_k^4] = 0.5 \cdot 1^4 + 0.5 \cdot (-1)^4 = 1$. Therefore, the variance in [Eq. \(3.5\)](#) simplifies to:

$$\text{(Var. Rademacher)} \quad \left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + 2 \left[(\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right] \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p}, \quad (3.7)$$

which is equal to the lower bound in [Eq. \(3.6\)](#). Therefore, [Theorem 3.1.1](#) shows the optimality of the Rademacher distribution for the approximate kernel [\(3.4\)](#).

As an alternative, one can consider the standard Gaussian distribution, that is, $w_k \sim \mathcal{N}(0, 1)$. Then we have $\mathbb{E}[w_k^4] = 3$ (the fourth moment of a standard Gaussian random variable), and the variance in Eq. (3.5) simplifies to:

$$\text{(Var. Gaussian)} \quad (\|\mathbf{x}\|^2\|\mathbf{y}\|^2 + 2(\mathbf{x}^\top \mathbf{y})^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p}. \quad (3.8)$$

This is larger than the variance of the Rademacher sketch in Eq. (3.7). Thus, in terms of the variance of the approximate kernel, the Gaussian sketch is less efficient than the Rademacher sketch. These arguments suggest that the Rademacher sketch should be a preferred choice. This is in contrast to (Hamid et al., 2014, Lemma 3.1) who obtain better statistical guarantees for Gaussian sketches than the ones for Rademacher sketches derived in (Kar and Karnick, 2012). This is most likely because the bound by Hamid et al. (2014) exploits upper bounds on higher order moments and is probably tighter than the one given by (Kar and Karnick, 2012) that does not exploit such information.

3.2 Complex-Valued Product-Sketches

We now introduce complex-valued polynomial sketches, one of our novel contributions. We do this by extending the analysis of Choromanski et al. (2017) for linear sketches to polynomial sketches.¹

As before, without loss of generality, we focus on approximating the homogeneous polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y})^p$ of degree $p \in \mathbb{N}$. Let $D \in \mathbb{N}$. Suppose we generate $p \times D$ complex-valued random vectors satisfying

$$\mathbf{z}_{i,\ell} \in \mathbb{C}^d \quad \text{satisfying} \quad \mathbb{E}[\mathbf{z}_{i,\ell} \overline{\mathbf{z}_{i,\ell}^\top}] = \mathbf{I}_d, \quad i \in \{1, \dots, p\}, \quad \ell \in \{1, \dots, D\} \quad (3.9)$$

We then define a complex-valued random feature map as

$$\Phi_C(\mathbf{x}) := \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{z}_{i,1}^\top \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{z}_{i,D}^\top \mathbf{x} \right) \right]^\top \in \mathbb{C}^D, \quad \mathbf{x} \in \mathbb{R}^d, \quad (3.10)$$

and the resulting approximate kernel as

$$\hat{k}_C(\mathbf{x}, \mathbf{y}) := \Phi_C(\mathbf{x})^\top \overline{\Phi_C(\mathbf{y})} = \frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p (\mathbf{z}_{i,\ell}^\top \mathbf{x}) \overline{(\mathbf{z}_{i,\ell}^\top \mathbf{y})}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \quad (3.11)$$

Eq. (3.11) is a generalization of the approximate kernel (3.2) with real-valued features, as Eq. (3.2) can be recovered by defining the complex random vectors $\mathbf{z}_{i,\ell}$

¹More specifically, Choromanski et al. (2017) analyze the variance of the *real part* of the approximate complex-valued kernel in Eq. (3.11) for $p = 1$. In contrast, we study Eq. (3.11) with generic $p \in \mathbb{N}$, and analyze the variance of Eq. (3.11) itself, including both the real and imaginary parts.

in Eq. (3.9) as real random vectors $\mathbf{w}_{i,\ell}$ in Eq. (3.3); in this case the requirement $\mathbb{E}[\mathbf{z}_{i,\ell}\overline{\mathbf{z}_{i,\ell}}^\top] = \mathbb{E}[\mathbf{w}_{i,\ell}\mathbf{w}_{i,\ell}^\top] = \mathbf{I}_d$ is satisfied.

For example, complex-valued random vectors $\mathbf{z}_{i,\ell}$ satisfying Eq. (3.9) can be generated as follows.

Example 1. Suppose we generate $2 \times p \times D$ independent real-valued random vectors

$$\mathbf{v}_{i,\ell}, \mathbf{w}_{i,\ell} \in \mathbb{R}^d \quad \text{satisfying} \quad \mathbb{E}[\mathbf{v}_{i,\ell}] = \mathbb{E}[\mathbf{w}_{i,\ell}] = \mathbf{0}, \quad \mathbb{E}[\mathbf{v}_{i,\ell}\mathbf{v}_{i,\ell}^\top] = \mathbb{E}[\mathbf{w}_{i,\ell}\mathbf{w}_{i,\ell}^\top] = \mathbf{I}_d \quad (3.12)$$

for $i \in \{1, \dots, p\}$, $\ell \in \{1, \dots, D\}$. Then one can define complex-valued random vectors (3.9) as

$$\mathbf{z}_{i,\ell} := \sqrt{\frac{1}{2}}(\mathbf{v}_{i,\ell} + i\mathbf{w}_{i,\ell}) \in \mathbb{C}^d, \quad i \in \{1, \dots, p\}, \quad \ell \in \{1, \dots, D\}. \quad (3.13)$$

The following two examples are specific cases of Example 1 and are complex versions of the real-valued Rademacher and Gaussian sketches discussed previously.

Example 2 (Complex Rademacher Sketch). In Example 1, suppose that elements of random vectors $\mathbf{v}_{i,\ell}$ and $\mathbf{w}_{i,\ell}$ are independently sampled from the Rademacher distribution, i.e., sampled uniformly from $\{1, -1\}$. Then the resulting random vectors $\mathbf{v}_{i,\ell}$, $\mathbf{w}_{i,\ell}$ satisfy the conditions in Eq. (3.12) and thus the complex random vectors in Eq. (3.13) satisfy the condition Eq. (3.9).

Example 3 (Complex Gaussian Sketch). In Example 1, suppose that elements of random vectors $\mathbf{v}_{i,\ell}$ and $\mathbf{w}_{i,\ell}$ are independently sampled from the standard Gaussian distribution, $\mathcal{N}(0, 1)$. Then the resulting random vectors $\mathbf{v}_{i,\ell}$, $\mathbf{w}_{i,\ell}$ satisfy the conditions in Eq. (3.12) and thus the complex random vectors in Eq. (3.13) satisfy the condition Eq. (3.9).

Example 4. Suppose the elements of each random vector $\mathbf{z}_{i,\ell} \in \mathbb{C}^d$ are independently sampled from the uniform distribution on $\{1, -1, i, -i\}$. Then the requirement in Eq. (3.9) is satisfied.

Example 4 is essentially identical to the complex Rademacher sketch in Example 2, in that each element of $\mathbf{z}_{i,\ell}$ in Example 4 can be obtained by multiplying $e^{i\pi/4}$ to an element of $\mathbf{z}_{i,\ell}$ in Example 2, and vice versa. The multiplication by $e^{i\pi/4}$ is equivalent to rotating an element counter-clockwise by 45 degrees. See Fig. 3.1 for an illustration. One can see that this multiplication by $e^{i\pi/4}$ does not change the resulting approximate kernel (3.11). In this sense, the constructions of Example 2 and Example 4 are equivalent. However, the sketch in Example 4 gives a computational advantage over Example 2: Since every element of each random vector $\mathbf{z}_{i,\ell}$ is either real or imaginary, the inner products $\mathbf{z}_{i,\ell}^\top \mathbf{x}$ in Eq. (3.10) can be computed at the same cost as for real polynomial sketches.

We show in the following proposition that the approximate kernel (3.11) is an unbiased estimator of the polynomial kernel $(\mathbf{x}^\top \mathbf{y})^p$.

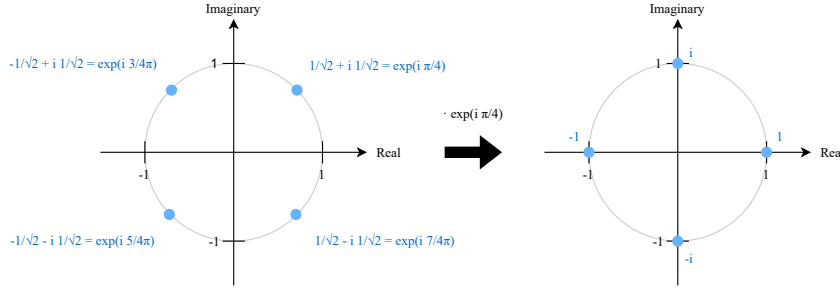


Figure 3.1: Multiplying each element of a random vector $\mathbf{z}_{i,\ell}$ in [Example 2](#) by $\exp(i\frac{\pi}{4})$ corresponds to a counter-clockwise rotation of that element by 45 degrees on the complex plane. The support of the resulting elements is $\{1, -1, i, -i\}$ and the construction of [Example 4](#) is obtained.

Proposition 3.2.1. *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary, and $\hat{k}_c(\mathbf{x}, \mathbf{y})$ be the approximate kernel in (3.11). Then we have*

$$\mathbb{E}[\hat{k}_c(\mathbf{x}, \mathbf{y})] = (\mathbf{x}^\top \mathbf{y})^p$$

Proof. Since [Eq. \(3.11\)](#) is the empirical average of D terms, it is sufficient to show the unbiasedness of each term. To this end, we consider here the case $D = 1$ and drop the index ℓ . We have

$$\mathbb{E} \left[\prod_{i=1}^p \mathbf{z}_i^\top \mathbf{x} \overline{\mathbf{z}_i^\top \mathbf{y}} \right] = \prod_{i=1}^p \mathbb{E} \left[\mathbf{z}_i^\top \mathbf{x} \overline{\mathbf{z}_i^\top \mathbf{y}} \right] = \prod_{i=1}^p \mathbf{x}^\top \mathbb{E} \left[\mathbf{z}_i \overline{\mathbf{z}_i} \right] \mathbf{y} = (\mathbf{x}^\top \mathbf{y})^p.$$

where we used [Eq. \(3.9\)](#) in the last identity. ■

3.3 Variance of Complex-Valued Product-Sketches

We now study the variance of the approximate kernel [\(3.11\)](#) with the complex-valued random feature map [\(3.10\)](#). As for the real-valued case, we consider the case $D = 1$ and drop the index ℓ :

$$\hat{k}_c(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^p \mathbf{z}_i^\top \mathbf{x} \overline{\mathbf{z}_i^\top \mathbf{y}}. \quad (3.14)$$

The variance of the case $D > 1$ can be obtained by dividing the variance of [Eq. \(3.14\)](#) by D , since the approximate kernel [\(3.11\)](#) is the average of D i.i.d. copies of [Eq. \(3.14\)](#). We denote by $z_{i,k}$ the k -th element of \mathbf{z}_i .

Note that the variance of a complex random variable $Z \in \mathbb{C}$ is defined by

$$\mathbb{V}[Z] := \mathbb{E}[|Z - \mathbb{E}[Z]|^2] = \mathbb{E}[(Z - \mathbb{E}[Z])(\overline{Z - \mathbb{E}[Z]})] = \mathbb{E}[|Z|^2] - |\mathbb{E}[Z]|^2$$

Theorem 3.3.1 below characterizes the variance in terms of the input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and the distribution of the complex weight vectors (3.9). The proof is given in Appendix B.1.1.

Theorem 3.3.1. *Let $\mathbf{x} := (x_1, \dots, x_d)^\top \in \mathbb{R}^d$ and $\mathbf{y} := (y_1, \dots, y_d)^\top \in \mathbb{R}^d$ be any input vectors. Let $\mathbf{z}_1, \dots, \mathbf{z}_p \in \mathbb{C}^d$ be i.i.d. random vectors satisfying (3.9), such that elements z_{i1}, \dots, z_{id} of each vector $\mathbf{z}_i = (z_{i1}, \dots, z_{id})^\top$ are themselves i.i.d. Let $\mathbf{z} = (z_1, \dots, z_d)^\top \in \mathbb{C}^d$ be a random vector independently and identically distributed as $\mathbf{z}_1, \dots, \mathbf{z}_p$, and write $z_k = a_k + ib_k$ with $a_k, b_k \in \mathbb{R}$. Suppose*

$$\mathbb{E}[a_k b_k] = 0, \quad \mathbb{E}[a_k^2] = q, \quad \mathbb{E}[b_k^2] = 1 - q \quad \text{where } 0 \leq q \leq 1. \quad (3.15)$$

Then, for the approximate kernel (3.14), we have

$$\begin{aligned} \mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] &= \left(\sum_{k=1}^d \mathbb{E}[|z_k|^4] x_k^2 y_k^2 + \|\mathbf{x}\|^2 \|\mathbf{y}\|^2 - \sum_{k=1}^d x_k^2 y_k^2 \right. \\ &\quad \left. + ((2q - 1)^2 + 1) \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right) \end{aligned} \quad (3.16)$$

$$\geq \left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + ((2q - 1)^2 + 1) \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right). \quad (3.17)$$

Theorem 3.3.1 is a generalization of Theorem 3.1.1 as it applies to a spectrum of complex polynomial sketches in terms of q , where the case $q = 1$ recovers Theorem 3.1.1 for real-valued polynomial sketches. The key condition in Theorem 3.3.1 is Eq. (3.15),² where the constant q is the average length of the real part a_k of each random element $z_k = a_k + ib_k$. Note that $\mathbb{E}[b_k^2] = 1 - q$ follows from $\mathbb{E}[a_k^2] = q$ since $1 = \mathbb{E}[|z_k|^2] = a_k^2 + b_k^2$. Eq. (3.15) is satisfied for Examples 2, 3 and 4 with $q = 1/2$ and for the real-valued Rademacher and Gaussian sketches with $q = 1$. If z_k is sampled uniformly from $\{i, -i\}$, which is eligible as it satisfies Eq. (3.9), then $q = 0$. If z_k is sampled uniformly from $\{1, -1\}$ with probability q and from $\{i, -i\}$ with probability $1 - q$, then Eq. (3.15) is satisfied with this q .

Eq. (3.17) is the smallest possible variance attainable by complex polynomial sketches satisfying the conditions in Theorem 3.3.1. For $q = 1/2$, this lower bound is attained by the complex Rademacher sketch (Example 2) and its equivalent construction (Example 4), for which we have $\mathbb{E}[|z_k|^4] = 1$:

$$\text{(Comp. Rademacher)} \quad \left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \quad (3.18)$$

On the other hand, for the complex Gaussian sketch (Example 3) we have $\mathbb{E}[|z_k|^4] = \mathbb{E}[(a_k^2 + b_k^2)^2] = 2$, and the variance is given by

$$\text{(Comp. Gaussian)} \quad \left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \quad (3.19)$$

²Eq. (3.15) implies that z_k is a *proper* complex random variable (Neeser and Massey, 1993).

Comparing the Real and Complex Polynomial Sketches. Let us now compare the variances of real ($q = 1$) and complex ($q \neq 1$) polynomial sketches. First, it is easy to see that the variance of the complex Gaussian polynomial sketch (Eq. (3.19)) is upper-bounded by the variance of the real Gaussian sketch (Eq. (3.8)). For the lower bound in Eq. (3.17), which is attained by the Rademacher sketches, a more detailed analysis is needed. To this end, consider the term that depends on q :

$$((2q - 1)^2 + 1) ((\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2)$$

The variance in Eq. (3.16) is a monotonically increasing function of this term. Suppose

$$(\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 = \sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d x_i x_j y_i y_j \geq 0 \quad (3.20)$$

Then $q = 1/2$ (e.g., complex sketches in Examples 2, 3 and 4) makes the term the smallest, while $q = 1$ and $q = 0$ (purely real and imaginary polynomial sketches) makes it the largest. In other words, for input vectors \mathbf{x} and \mathbf{y} satisfying Eq. (3.20), complex-valued sketches with $q = 1/2$ result in a lower variance than the real-valued counterparts with $q = 1$. On the other hand, if Eq. (3.20) does not hold, real-valued sketches result in a lower variance than the complex-valued counterparts.

Therefore, whether complex-valued Rademacher sketches ($q = 1/2$) yield a lower variance than real-valued Rademacher sketches ($q = 1$) depends on whether Eq. (3.20) holds. For example, Eq. (3.20) holds true if input vectors $\mathbf{x} = (x_1, \dots, x_d)^\top$ and $\mathbf{y} = (y_1, \dots, y_d)^\top$ are *nonnegative*: $x_1, \dots, x_d \geq 0$ and $y_1, \dots, y_d \geq 0$. Nonnegative input vectors are ubiquitous in real-world applications, e.g., where each input feature represents the amount of a certain quantity, where input vectors are given by bag-of-words representations, one-hot encoding (categorical data), or min-max feature scaling, and where they are outputs of a ReLU neural network³. For such applications with nonnegative input vectors, complex-valued polynomial sketches always yield a smaller variance than the real-valued counterparts.

3.4 Probabilistic Error Bounds for Rademacher Sketches

We present here probabilistic error bounds for the approximate kernel in Eq. (3.11) in terms of the number D of random features, using the variance formula obtained in the previous section and focusing on Rademacher sketches. The proof of the following result is given in Appendix B.1.2.

³c.f. DeepFried Convnets (Yang et al., 2015) and fine-grained image recognition (Gao et al., 2016).

Theorem 3.4.1. *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary input vectors. For $0 \leq q \leq 1$, consider a polynomial sketch in [Theorem 3.3.1](#) such that $\mathbb{E}[|z_k|^4] = 1$ and thus attains the variance in [Eq. \(3.17\)](#). Define a constant $\sigma^2 \geq 0$ by*

$$\sigma^2 := \frac{1}{\|\mathbf{x}\|^{2p}\|\mathbf{y}\|^{2p}} \left[\left(\|\mathbf{x}\|^2\|\mathbf{y}\|^2 + ((2q-1)^2 + 1) ((\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2) \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right]$$

Let $\epsilon, \delta > 0$ be arbitrary, and $D \in \mathbb{N}$ be such that

$$D \geq 2 \left(\frac{2}{3\epsilon} + \frac{\sigma^2}{\epsilon^2} \right) \log \left(\frac{2}{\delta} \right). \quad (3.21)$$

Then, for the approximate kernel $\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})$ in [Eq. \(3.11\)](#), we have

$$\Pr \left[\left| \hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y}) - (\mathbf{x}^\top \mathbf{y})^p \right| \leq \epsilon \|\mathbf{x}\|_1^p \|\mathbf{y}\|_1^p \right] \geq 1 - \delta.$$

[Eq. \(3.21\)](#) shows that the required number D of random features to achieve the relative accuracy of ϵ (where the “relative” is with respect to $\|\mathbf{x}\|_1^p \|\mathbf{y}\|_1^p$) with probability at least $1 - \delta$. For small ϵ , the second term σ^2/ϵ^2 dominates the first term $2/(3\epsilon)$. This second term depends on σ^2 , which is a scaled version of the variance in [Eq. \(3.17\)](#) of the approximate kernel for $D = 1$. Thus, if the variance in [Eq. \(3.17\)](#) is smaller (resp. larger), one needs a smaller (resp. larger) number of random features to achieve the relative accuracy of ϵ .

Let us now compare the real-valued ($q = 1$) and complex-valued ($q = 1/2$) Rademacher sketches. As discussed earlier, the complex Rademacher sketch has a smaller variance than the real Rademacher sketch when the inequality in [Eq. \(3.20\)](#) holds for the two input vectors \mathbf{x}, \mathbf{y} . In particular, this inequality always holds when the input vectors are nonnegative. Therefore, in this case, the complex Rademacher sketch requires a smaller number of random features than the real Rademacher sketch to achieve a given accuracy.

When the inequality in [Eq. \(3.20\)](#) holds, the advantage of the complex Rademacher sketch becomes more significant for larger p . We illustrate this in [Fig. 3.2](#), where $\mathbf{x} = \mathbf{y} = (1, \dots, 1)^\top / \sqrt{d} \in \mathbb{R}^d$. In this case, we have $\sigma^2 = (2 - 1/d)^p - 1$ for the complex Rademacher sketch ($q = 1/2$), while we have $\sigma^2 = (3 - 2/d)^p - 1$ for the real Rademacher sketch ($q = 1$). For larger p , the value of σ^2 for the real Rademacher sketch becomes relatively larger than that for the complex Rademacher sketch, implying that the complex Rademacher sketch is more efficient. [Fig. 3.2](#) empirically supports this observation.

3.5 Structured Product-Sketches

We study here *structured* polynomial sketches and their extensions with complex features. In [Section 3.1](#), we studied polynomial sketches in [Eq. \(3.4\)](#) (or [Eq. \(3.10\)](#))

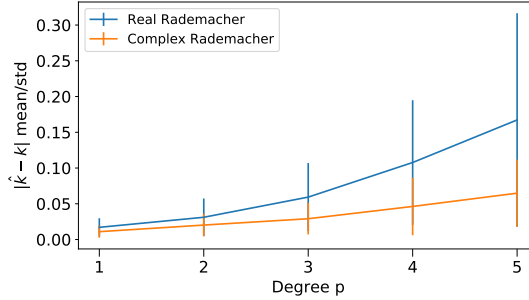


Figure 3.2: This plot shows the mean absolute error $\mathbb{E}[|\hat{k}(\mathbf{x}, \mathbf{y}) - (\mathbf{x}^\top \mathbf{y})^p|]$ for different values of the degree p , where the mean is taken over 100 independent constructions of the approximate kernel $\hat{k}(\mathbf{x}, \mathbf{y})$, for both the real and complex Rademacher sketches. The number of random features is $D = 5,000$ and the dimensionality of input vectors is $d = 1,000$. The input vectors are $\mathbf{x} = \mathbf{y} = (1, \dots, 1)^\top / \sqrt{d} \in \mathbb{R}^d$.

for complex extensions), where the $p \times D$ random vectors $\mathbf{w}_{i,\ell} \in \mathbb{R}^d$ ($i = 1, \dots, p$, $\ell = 1, \dots, D$) are generated in an i.i.d. manner. By putting a structural constraint on these vectors, one can construct more efficient random features with a lower variance. Moreover, such a structural constraint leads to a computational advantage, as the imposed structure may be used for implementing an efficient algorithm for fast matrix multiplication.

We develop two novel structured polynomial sketches that are closely related to *TensorSRHT* (Ahle et al., 2020) introduced in Chapter 2. Just like the original version by Ahle et al. (2020), our sketches apply the Subsampled Randomized Hadamard Transform (SRHT) to an input $\mathbf{x} \in \mathbb{R}^d$ in order to then form an implicit sketch of $\mathbf{x}^{(p)}$. The key difference is that our sketches are embedded into the Product-Sketch framework and permit us to derive closed-form variances for them. Currently there are only upper bounds on the variances for structured polynomial sketches like TensorSRHT and TensorSketch available in the literature, which is why statistical advantages have only been shown empirically up until now. Our contribution is to theoretically characterize statistical advantages over non-structured polynomial sketches by means of closed-form variance derivations and comparisons. Since our two proposed sketches are closely related to TensorSRHT by Ahle et al. (2020), we will stick to the term TensorSRHT and add further precisions when necessary.

Our second main contribution in this section is to develop a complex extension of the aforementioned structured sketches and to study their variance advantage empirically. In Section 3.5.1, we introduce our version of TensorSRHT with real features, and present its extension using complex features in Section 3.5.2. We then make a comparison between the real and complex versions in Section 3.5.3.

3.5.1 Our Version of TensorSRHT

Recall from Chapter 2 that TensorSRHT by Ahle et al. (2020) is defined as

$$\begin{aligned}\Phi(\mathbf{x}) &= \mathbf{P}_\pi(\mathbf{H}_d \mathbf{D}_1 \mathbf{x} \otimes \cdots \otimes \mathbf{H}_d \mathbf{D}_p \mathbf{x}) / \sqrt{D} \\ &= \mathbf{P}_\pi \mathbf{H}_{dp}(\mathbf{D}_1 \otimes \cdots \otimes \mathbf{D}_p) \mathbf{x}^{(p)} / \sqrt{D},\end{aligned}\quad (3.22)$$

where $\mathbf{P}_\pi \in \{0, 1\}^{D \times dp}$ is a sampling matrix, \mathbf{H}_n is the unnormalized Hadamard matrix of size $\{0, 1\}^{n \times n}$ for some $n = 2^m, m \in \mathbb{N}$ and $\{\mathbf{D}_i\}_{i=1}^p$ are diagonal Rademacher matrices of size $\mathbb{R}^{d \times d}$.

By looking at Eq. (3.22), one can see that TensorSRHT is constructed by forming D distinct combinations of the elements of $\{\mathbf{H}_d \mathbf{D}_i \mathbf{x}\}_{i=1}^p$. More precisely, we sample the index tuples $\{(i_{1,\ell}, \dots, i_{p,\ell})\}_{\ell=1}^D$ without replacement, where $i_{1,\ell}, \dots, i_{p,\ell}$ are uniformly sampled from $\{1, \dots, d\}$. Each element $\{\Phi(\mathbf{x})_\ell\}_{\ell=1}^D$ is then computed as

$$\Phi(\mathbf{x})_\ell = (\mathbf{H}_d \mathbf{D}_1 \mathbf{x})_{i_{1,\ell}} \cdots (\mathbf{H}_d \mathbf{D}_p \mathbf{x})_{i_{p,\ell}} / \sqrt{D}.\quad (3.23)$$

In the following, we propose two modifications of TensorSRHT that can be represented in terms of an equivalent Product-Sketch (3.4). The advantage of our proposed sketches is that we can compute their variances in closed form making it possible to study the statistical advantages of our TensorSRHT modifications compared to non-structured Product-Sketches.

We call our modified TensorSRHT sketches *stacked TensorSRHT* and *upsampled TensorSRHT*, respectively. We developed stacked TensorSRHT in our earlier work (Wacker et al., 2022a) since we could easily derive a convex surrogate function of its variance over D , as it is needed for our optimized Maclaurin method that we propose in Chapter 5. In our later work (Wacker et al., 2022b), we developed the upsampled TensorSRHT that reduces the time complexity from $\mathcal{O}(p(D \log d + D))$ of stacked TensorSRHT to $\mathcal{O}(p(d \log d + D))$ to project a single data point, and is thus faster when $D > d$. Upsampled TensorSRHT has the same time complexity as TensorSRHT by Ahle et al. (2020) that can be considered the state-of-the-art in terms of computation time. We start by defining stacked TensorSRHT and introduce upsampled TensorSRHT afterwards.

Stacked TensorSRHT

To simplify the definition of stacked TensorSRHT, we distinguish two cases for D , namely $D \leq d$ and $D > d$. We now cover the first case.

Case $D \leq d$. For now we assume that the number D of random features is less or equal to the dimensionality d of the input vectors: $D \leq d$. For $i = 1, \dots, p$, define $\mathbf{w}_i \in \mathbb{R}^d$ as a random vector whose elements are i.i.d. Rademacher random variables:

$$\mathbf{w}_i := (w_{i,1}, \dots, w_{i,d})^\top \in \mathbb{R}^d, \quad w_{i,j} \stackrel{i.i.d.}{\sim} \text{unif}(\{1, -1\}) \quad (j = 1, \dots, d)$$

Consider a random permutation of the indices $\pi_i : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$ for every $i = 1, \dots, p$, and let

$$\pi_i(1), \dots, \pi_i(d)$$

be the permuted indices. For $i = 1, \dots, p$ and $\ell = 1, \dots, D$, we then define a random vector $\mathbf{s}_{i,\ell} \in \mathbb{R}^d$ as the Hadamard product (i.e., element-wise product) of the Rademacher vector \mathbf{w}_i and the permuted column $\mathbf{h}_{\pi_i(\ell)}$ of the Hadamard matrix:

$$\mathbf{s}_{i,\ell} := \mathbf{w}_i \circ \mathbf{h}_{\pi_i(\ell)} = (w_{i,1}h_{\pi_i(\ell),1}, \dots, w_{i,d}h_{\pi_i(\ell),d})^\top \in \mathbb{R}^d, \quad (3.24)$$

where $h_{\pi_i(\ell),j}$ denotes the j -th element of $\mathbf{h}_{\pi_i(\ell)}$.

Because of the orthogonality of the columns $\mathbf{h}_1, \dots, \mathbf{h}_d$ of the Hadamard matrix \mathbf{H}_d , the random weight vectors $\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,d}$ are orthogonal to each other almost surely: for $\ell \neq m$, we have

$$\mathbf{s}_{i,\ell}^\top \mathbf{s}_{i,m} = (\mathbf{w}_i \circ \mathbf{h}_{\pi_i(\ell)})^\top (\mathbf{w}_i \circ \mathbf{h}_{\pi_i(m)}) = \sum_{j=1}^d w_{i,j}^2 h_{\pi_i(\ell),j} h_{\pi_i(m),j} = h_{\pi_i(\ell)}^\top h_{\pi_i(m)} = 0.$$

Note also that, given the permutation $\pi_i(1), \dots, \pi_i(d)$, the elements of each random vector $\mathbf{s}_{i,\ell}$ in (3.24) are i.i.d. Rademacher variables.

Finally, we define a random feature map $\Phi_{\mathcal{R}}(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^D$ for the case $D \leq d$ as

$$\Phi_{\mathcal{R}}(\mathbf{x}) := \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{s}_{i,1}^\top \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{s}_{i,D}^\top \mathbf{x} \right) \right]^\top \in \mathbb{R}^D, \quad (3.25)$$

which defines an approximate kernel as

$$\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) := \Phi_{\mathcal{R}}(\mathbf{x})^\top \Phi_{\mathcal{R}}(\mathbf{y}) = \frac{1}{D} \sum_{\ell=1}^D \Phi_{\mathcal{R}}(\mathbf{x})_\ell \Phi_{\mathcal{R}}(\mathbf{y})_\ell$$

where $\Phi_{\mathcal{R}}(\cdot)_\ell$ denotes the ℓ -th element of $\Phi_{\mathcal{R}}(\cdot)$.

The orthogonality of the weight vectors in Eq. (3.24) leads to *negative covariances* between the terms $\Phi_{\mathcal{R}}(\mathbf{x})_\ell \Phi_{\mathcal{R}}(\mathbf{y})_\ell$ and $\Phi_{\mathcal{R}}(\mathbf{x})_m \Phi_{\mathcal{R}}(\mathbf{y})_m$ with distinct indices $\ell \neq m$ in the approximate kernel. These negative covariances decrease the overall variance of the approximate kernel, as we will show later in [Theorem 3.5.2](#).

For $D \leq d$, our version of TensorSRHT (3.25) and the one by [Ahle et al. \(2020\)](#) in Eq. (3.22) are almost equivalent. To see this, let $(\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,d})^\top \mathbf{x} = \mathbf{P}_{\pi_i} \mathbf{H}_d \mathbf{D}_i \mathbf{x}$ with $\mathbf{D}_i := \text{diag}(\mathbf{w}_i)$ and $\mathbf{P}_{\pi_i} := (\mathbf{e}_{\pi_i(1)}, \dots, \mathbf{e}_{\pi_i(d)})^\top \in \mathbb{R}^{d \times d}$, where $\mathbf{e}_{\pi_i(\ell)} \in \mathbb{R}^d$ is a vector whose $\pi_i(\ell)$ -th element is 1 and other elements are 0 ($\ell = 1, \dots, d$). \mathbf{P}_{π_i} has the role of shuffling the elements of $\mathbf{H}_d \mathbf{D}_i \mathbf{x}$ at random. As a result, the ℓ -th element of $\Phi_{\mathcal{R}}(\mathbf{x})$ (3.25) can be expressed as

$$\Phi_{\mathcal{R}}(\mathbf{x})_\ell = (\mathbf{H}_d \mathbf{D}_1 \mathbf{x})_{\pi_1(\ell)} \cdots (\mathbf{H}_d \mathbf{D}_p \mathbf{x})_{\pi_p(\ell)} / \sqrt{D}. \quad (3.26)$$

Each permutation index $\{\pi_i(\ell)\}_{\ell=1}^D$ is sampled from $\{1, \dots, d\}$ without replacement for every $i = 1, \dots, p$ due to the shuffling operation. From this it follows that the tuples $\{(\pi_1(\ell), \dots, \pi_p(\ell))\}_{\ell=1}^D$ are sampled without replacement from $\{1, \dots, d\}^p$. However, it is not possible to obtain two tuples that have equal indices $\pi_i(\ell) = \pi_i(\ell')$ for some $\ell \neq \ell'$ and $i \in 1, \dots, p$, which shows a minor difference between our version of TensorSRHT (3.26) and TensorSRHT by Ahle et al. (2020) (3.23).

Case $D > d$. We now explain how to draw $D > d$ random feature samples, which requires an extension of the sketch in Eq. (3.25) described so far. In the case of stacked TensorSRHT, we independently generate the feature map in Eq. (3.25) a total number of $B := \lceil \frac{D}{d} \rceil$ times and concatenate the resulting B vectors to obtain a Bd -dimensional feature map, and then discard the redundant last $Bd - D$ components of it to obtain a D -dimensional feature map. In this way, we can obtain a D -dimensional feature map for arbitrary $D \in \mathbb{N}$, which we can still write as

$$\Phi_{\mathcal{R}}(\mathbf{x}) := \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{s}_{i,1}^{\top} \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{s}_{i,D}^{\top} \mathbf{x} \right) \right]^{\top} \in \mathbb{R}^D. \quad (3.27)$$

as in Eq. (3.25), but with $\{\mathbf{s}_{i,\ell}\}_{\ell=1}^D$ corresponding to our stacked construction. The entire procedure for constructing the structured polynomial sketch in Eq. (3.27) is outlined in Algorithm 2, where we also cover the complex-valued case discussed later.

In Algorithm 2, we use the equivalent matrix formulation introduced in Eq. (3.26), since it enables the Fast Walsh-Hadamard transform by employing the associativity, and thus the feature map can be computed much faster. To explain this more precisely, we can compute

$$(\mathbf{s}_{i,1}^{\top} \mathbf{x}, \dots, \mathbf{s}_{i,d}^{\top} \mathbf{x}) = \mathbf{x}^{\top} (\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,d}) = \mathbf{x}^{\top} (\mathbf{D}_i \mathbf{H}_d \mathbf{P}_{\pi_i}) = ((\mathbf{x}^{\top} \mathbf{D}_i) \mathbf{H}_d) \mathbf{P}_{\pi_i}$$

by 1) first computing $\mathbf{x}^{\top} \mathbf{D}_i$, 2) then multiplying the Hadamard matrix \mathbf{H}_d using the Fast Walsh-Hadamard transform, and 3) lastly multiplying the permutation matrix \mathbf{P}_{π_i} , which is more efficient than first precomputing $\mathbf{D}_i \mathbf{H}_d \mathbf{P}_{\pi_i}$ and then multiplying \mathbf{x}^{\top} . In this way, thanks to the Fast Walsh-Hadamard transform, $(\mathbf{s}_{i,1}^{\top} \mathbf{x}, \dots, \mathbf{s}_{i,d}^{\top} \mathbf{x})$ can be computed in $\mathcal{O}(d \log d)$ instead of $\mathcal{O}(d^2)$ (Fino and Algazi, 1976). The total computational complexity is therefore $\mathcal{O}(p(D \log d + D))$ and the memory requirement is $\mathcal{O}(pD)$, and this is a computational advantage over the non-structured i.i.d. estimator in Eq. (3.4).

The feature map in Eq. (3.27) induces an approximate kernel $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{R}}(\mathbf{x})^{\top} \Phi_{\mathcal{R}}(\mathbf{y})$. The following proposition summarizes that this approximate kernel is unbiased with respect to the target polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^{\top} \mathbf{y})^p$. As mentioned earlier, stacked TensorSRHT discussed here is slightly different from TensorSRHT by Ahle et al. (2020). Therefore, this result is novel in its own right. The result follows from Proposition 3.5.8 in the next subsection, so we omit the proof.

Algorithm 2: Real and Complex stacked TensorSRHT

Result: A feature map $\Phi_{\mathcal{R}/\mathbb{C}}(\mathbf{x})$

Pad \mathbf{x} with zeros so that d becomes a power of 2 ;

Let $B = \lceil \frac{D}{d} \rceil$ be the number of stacked projection blocks ;

forall $b \in \{1, \dots, B\}$ **do**

forall $i \in \{1, \dots, p\}$ **do**

Real case Generate a random vector $\mathbf{w}_i = (w_{i,1}, \dots, w_{i,d})^\top \in \mathbb{R}^d$
as $w_{i,1}, \dots, w_{i,d} \stackrel{i.i.d.}{\sim} \text{unif}(\{1, -1\})$, and define a diagonal matrix
 $\mathbf{D}_i := \text{diag}(\mathbf{w}_i) \in \mathbb{R}^{d \times d}$;

Complex case Generate a random vector
 $\mathbf{z}_i = (z_{i,1}, \dots, z_{i,d})^\top \in \mathbb{C}^d$ as $z_{i,1}, \dots, z_{i,d} \stackrel{i.i.d.}{\sim} \text{unif}(\{1, -1, i, -i\})$,
and define a diagonal matrix $\mathbf{D}_i := \text{diag}(\mathbf{z}_i) \in \mathbb{C}^{d \times d}$;

Randomly permute the indices $1, \dots, d$ to $\pi_i(1), \dots, \pi_i(d)$;

Let $\mathbf{P}_\pi := (\mathbf{e}_{\pi_i(1)}, \dots, \mathbf{e}_{\pi_i(d)}) \in \mathbb{R}^{d \times d}$, where $\mathbf{e}_{\pi_i(\ell)} \in \mathbb{R}^d$ is a vector
whose $\pi_i(\ell)$ -th element is 1 and other elements are 0 ($\ell = 1, \dots, d$)
;

Let $(\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,d}) := \mathbf{D}_i \mathbf{H}_d \mathbf{P}_\pi$;

end

Compute $\Phi_b(\mathbf{x}) := \sqrt{1/D} [(\prod_{i=1}^p \mathbf{s}_{i,1}^\top \mathbf{x}), \dots, (\prod_{i=1}^p \mathbf{s}_{i,d}^\top \mathbf{x})]^\top$;

end

Concatenate the elements of $\Phi_1(\mathbf{x}), \dots, \Phi_B(\mathbf{x})$ to yield a single projection
vector $\Phi_{\mathcal{R}/\mathbb{C}}(\mathbf{x})$ and keep the first D entries ;

Proposition 3.5.1. *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary, and $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{R}}(\mathbf{x})^\top \Phi_{\mathcal{R}}(\mathbf{y})$ be the approximate kernel with $\Phi_{\mathcal{R}}(\mathbf{x}), \Phi_{\mathcal{R}}(\mathbf{y}) \in \mathbb{R}^D$ given by the random feature map in Eq. (3.27). Then we have $\mathbb{E}[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})] = (\mathbf{x}^\top \mathbf{y})^p$.*

We next study the variance of the approximate kernel given by stacked TensorSRHT, which is the mean squared error of the approximate kernel since it is unbiased as shown above. The following theorem provides a closed form expression for the variance, whose proof is given for the more general complex case in Appendix B.2.2. It is a novel result and extends Choromanski et al. (2017, Theorem 3.3) to the setting $p > 1$ and $D > d$.

Theorem 3.5.2 (Variance of Real Stacked TensorSRHT). *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary, and $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{R}}(\mathbf{x})^\top \Phi_{\mathcal{R}}(\mathbf{y})$ be the approximate kernel with $\Phi(\mathbf{x}), \Phi(\mathbf{y}) \in \mathbb{R}^D$ given by the random feature map in Eq. (3.27). Then we have*

$$\mathbb{V} \left[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) \right] = \underbrace{\frac{V_{\text{Rad}}^{(p)}}{D}}_{(A)} - \underbrace{\frac{c(D, d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_{\text{Rad}}^{(1)}}{d-1} \right)^p \right]}_{(B)}, \quad (3.28)$$

where $V_{\text{Rad}}^{(p)} \geq 0$ and $V_{\text{Rad}}^{(1)} \geq 0$ are the variances of the real Rademacher sketch with a single feature in Eq. (3.7) with generic $p \in \mathbb{N}$ and $p = 1$, respectively, and $c(D, d) \in \mathbb{N}$ is defined by

$$c(D, d) := \lfloor D/d \rfloor d(d-1) + \text{mod}(D, d)(\text{mod}(D, d) - 1). \quad (3.29)$$

Remark 3.5.3. The constant $c(D, d)$ in Eq. (3.29) is the number of pairs of indices $\ell, \ell' = 1, \dots, D$ with $\ell \neq \ell'$ for which the covariance of the weight vectors $\mathbf{s}_{i,\ell}$ and $\mathbf{s}_{i,\ell'}$ in Eq. (3.27) is non-zero (see the proof in Appendix B.2.2 for details). If $D = Bd$ for some $B \in \mathbb{N}$, this constant simplifies to $c(D, d) = Bd(d-1)$, and the variance in Eq. (3.28) becomes

$$\mathbb{V} \left[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) \right] = \frac{1}{D} V_{\text{Rad}}^{(p)} - \frac{d-1}{D} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_{\text{Rad}}^{(1)}}{d-1} \right)^p \right].$$

An interesting subcase is $p = 1$, for which the variance becomes zero. Thus, setting $D \in \{kd \mid k \in \mathbb{N}\}$ for $p = 1$ is equivalent to using the linear kernel with the original inputs.

Theorem 3.5.2 enables understanding the condition under which stacked TensorSRHT has a smaller variance than the unstructured Rademacher sketch in Eq. (3.1). Note that the term (A) in Eq. (3.28) is the variance of the approximate kernel with the Rademacher sketch with D features. On the other hand, the term (B) in Eq. (3.28) can be interpreted as the effect of the structured sketch. The term (B) always becomes non-negative when p is odd, and thus the overall variance of stacked TensorSRHT becomes smaller than the Rademacher sketch, as summarized in the following corollary. Thus, when p is odd, stacked TensorSRHT should be preferred over the Rademacher sketch.

Corollary 3.5.4. Let $p \in \mathbb{N}$ be odd. Then, for all input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the variance of the approximate kernel with stacked TensorSRHT in Eq. (3.28) is smaller or equal to the variance of the approximate kernel with the Rademacher sketch:

$$\frac{V_{\text{Rad}}^{(p)}}{D} - \frac{c(D, d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_{\text{Rad}}^{(1)}}{d-1} \right)^p \right] \leq \frac{V_{\text{Rad}}^{(p)}}{D}$$

Proof. Since, $V_{\text{Rad}}^{(1)} \geq 0$, we have $(\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} V_{\text{Rad}}^{(1)} \leq (\mathbf{x}^\top \mathbf{y})^2$. For odd p this leads to $\left((\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} V_{\text{Rad}}^{(1)} \right)^p \leq (\mathbf{x}^\top \mathbf{y})^{2p}$. The assertion immediately follows. \blacksquare

If p is even, on the other hand, the variance of stacked TensorSRHT can be larger than the Rademacher sketch for certain input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. For instance, if \mathbf{x} and \mathbf{y} are orthogonal, i.e., $\mathbf{x}^\top \mathbf{y} = 0$, then the variance of TensorSRHT in Eq. (3.28) is

$$\text{Eq. (3.28)} = \frac{V_{\text{Rad}}^{(p)}}{D} + \frac{c(D, d)}{D^2} \left(\frac{V_{\text{Rad}}^{(1)}}{d-1} \right)^p \geq \frac{V_{\text{Rad}}^{(p)}}{D}.$$

Therefore, for even p , we do not have a theoretical guarantee for the advantage of stacked TensorSRHT over the Rademacher sketch in terms of their variances. In practice, however, stacked TensorSRHT has often a smaller variance than the Rademacher sketch also for even p , as demonstrated in our experiments described later. Moreover, stacked TensorSRHT has a computational advantage over the Rademacher sketch, thanks to the fast Walsh-Hadamard transform.

Remark 3.5.5. *One can straightforwardly derive a probabilistic error bound for stacked TensorSRHT by using Theorem 3.5.2 and Chebyshev’s inequality. However, deriving an exponential tail bound for stacked TensorSRHT is more involved despite our derived closed form variance formula, because different features in the feature map $\Phi_{\mathcal{R}}(\mathbf{x})$ in Eq. (3.25) are dependent for stacked TensorSRHT and thus applying Bernstein’s inequality is not straightforward. One can find an exponential tail bound for TensorSRHT in Ahle et al. (2020, Lemma 33 in the longer version), while they analyze a slightly different version of TensorSRHT from ours and their bound is a uniform upper bound that holds for all input vectors simultaneously.*

Our variance formula in Eq. (3.28), which is a novel contribution to the literature, provides a precise characterization of how the variance of the approximate kernel depends on the input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, and shows when stacked TensorSRHT is more advantageous than the Rademacher sketch. Moreover, as the variance formula can be computed in practice, it can be used for designing an objective function for a certain optimization problem, as we do in Chapter 5 for designing a data-driven approach to feature construction.

Upsampled TensorSRHT

We now briefly discuss our upsampled TensorSRHT modification that is faster than stacked TensorSRHT if $D > d$ and more closely related to the original construction in Ahle et al. (2020) while still providing closed form variances. The disadvantage of upsampled TensorSRHT is that we have not derived a convex surrogate function for its variance, obstructing its use for the optimized Maclaurin method in Chapter 5.

For $D \leq d$, upsampled TensorSRHT and stacked TensorSRHT are equivalent. We will thus directly focus on the case $D > d$. Instead of generating the feature map in Eq. (3.25) a total number of $B = \lceil D/d \rceil$ times independently as for stacked TensorSRHT, we keep the construction in Eq. (3.25) and change only the sampling matrices $\{\mathbf{P}_{\pi_i}\}_{i=1}^p$.

For this purpose, we concatenate the index vector $(1, \dots, d)^\top \in \mathbb{R}^d$ a total number of B times with itself and call the resulting vector $\mathbf{p}_{\text{base}} \in \mathbb{R}^{\lceil D/d \rceil d}$. We then shuffle the elements of \mathbf{p}_{base} for each $i \in \{1, \dots, p\}$ independently at random and call the shuffled vectors $\{\mathbf{p}_i\}_{i=1}^p$. Each vector \mathbf{p}_i defines a sampling matrix $\mathbf{P}_i = (\mathbf{e}_{p_{i,1}}, \dots, \mathbf{e}_{p_{i,D}})^\top \in \{0, 1\}^{D \times d}$. We keep only the first D elements of every \mathbf{p}_i to create D random feature samples. Now we let $(\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,D}) := \mathbf{D}_i \mathbf{H}_d \mathbf{P}_i$.

Using $\{\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,D}\}_{i=1}^p$ as part of Eq. (3.25) defines our upsampled TensorSRHT sketch.

The computation of this sketch can be equivalently written as

$$\Phi_{\mathcal{R}}(\mathbf{x}) := (\mathbf{P}_1 \mathbf{H}_d \mathbf{D}_1 \mathbf{x}) \odot \dots \odot (\mathbf{P}_p \mathbf{H}_d \mathbf{D}_p \mathbf{x}) / \sqrt{D}, \quad (3.30)$$

which allows for a computation in $\mathcal{O}(p(d \log d + D))$ time.

The unbiasedness of upsampled TensorSRHT can be shown by following the same arguments as in Proposition 3.5.8. The following theorem provides a closed form expression for the variance, whose proof is given for the more general complex case in Appendix B.2.2.

Theorem 3.5.6 (Variance of Real Upsampled TensorSRHT). *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary, and $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{R}}(\mathbf{x})^\top \Phi_{\mathcal{R}}(\mathbf{y})$ be the approximate kernel with $\Phi(\mathbf{x}), \Phi(\mathbf{y}) \in \mathbb{R}^D$ given by the random feature map in Eq. (3.30). Then we have*

$$\mathbb{V} \left[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) \right] = \underbrace{\frac{V_{\text{Rad}}^{(p)}}{D}}_{(A)} - \underbrace{\left(1 - \frac{1}{D}\right) \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_{\text{Rad}}^{(1)}}{\lceil D/d \rceil d - 1} \right)^p \right]}_{(B)}, \quad (3.31)$$

where $V_{\text{Rad}}^{(p)} \geq 0$ and $V_{\text{Rad}}^{(1)} \geq 0$ are the variances of the real Rademacher sketch with a single feature in Eq. (3.7) with generic $p \in \mathbb{N}$ and $p = 1$, respectively, as before.

Remark 3.5.7. $c(D, d)$ in Eq. (3.28) has been replaced by $(1 - 1/D)$ in Eq. (3.31) because all weight vectors $\mathbf{s}_{i,\ell}$ in Eq. (3.25) have a non-zero covariance for upsampled TensorSRHT. At the same time, $\frac{V_{\text{Rad}}^{(1)}}{d-1}$ is replaced by $\frac{V_{\text{Rad}}^{(1)}}{\lceil D/d \rceil d - 1}$, which counterbalances the aforementioned upscaling from $c(D, d)$ to $(1 - 1/D)$. We observed that upsampled TensorSRHT and stacked TensorSRHT perform similarly well in preliminary experiments although we did not carry out a detailed comparison since we developed both sketches in separate works and for a different purpose.

3.5.2 Complex-Valued TensorSRHT

We present here a generalization of stacked and upsampled TensorSRHT by allowing for complex features. To this end, let $z \in \mathbb{C}$ be a random variable such that (i) $|z| = 1$ almost surely, (ii) $\mathbb{E}[z] = 0$ and (iii) z is symmetric, i.e., the distributions of z and $-z$ are the same. Define then $\mathbf{z}_1, \dots, \mathbf{z}_p \in \mathbb{C}^d$ as i.i.d. complex random vectors such that elements of each random vector \mathbf{z}_i are i.i.d. realizations of z :

$$\mathbf{z}_i = (z_{i,1}, \dots, z_{i,d})^\top \in \mathbb{C}^d, \quad z_{i,j} \stackrel{i.i.d.}{\sim} P_z \quad (j = 1, \dots, d), \quad (3.32)$$

where P_z denotes the probability distribution of z .

We propose the complex version of stacked TensorSRHT in detail, while the complex version of upsampled TensorSRHT is obtained using the same modifications. Let $\pi : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$ be a random permutation of indices $1, \dots, d$. For $i = 1, \dots, p$ and $\ell = 1, \dots, D$, we then define a random vector $\mathbf{s}_{i,\ell} \in \mathbb{C}^d$ as the Hadamard product of the random vector \mathbf{z}_i in (3.32) and the permuted column $\mathbf{h}_{\pi(\ell)}$ of the Hadamard matrix \mathbf{H}_d :

$$\mathbf{s}_{i,\ell} := \mathbf{z}_i \circ \mathbf{h}_{\pi(\ell)} = (z_{i,1}h_{\pi(\ell),1}, \dots, z_{i,d}h_{\pi(\ell),d})^\top \in \mathbb{C}^d, \quad (3.33)$$

With these weight vectors $\mathbf{s}_{i,\ell}$, we define a random feature map exactly in the same way as the feature map in Eq. (3.27) for the real stacked TensorSRHT in Section 3.5.1. We define the resulting feature map $\Phi_C : \mathbb{R}^d \rightarrow \mathbb{C}^D$ by

$$\Phi_C(\mathbf{x}) := \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{s}_{i,1}^\top \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{s}_{i,D}^\top \mathbf{x} \right) \right]^\top \in \mathbb{C}^D. \quad (3.34)$$

We call this feature construction *complex stacked TensorSRHT*. Substituting the diagonal real Rademacher matrices $\{\mathbf{D}_i\}_{i=1}^p$ in Eq. (3.30) by the matrices $\{\text{diag}(\mathbf{z}_i)\}_{i=1}^p$ instead results in the complex version of upsampled TensorSRHT.

Admissible examples of the distribution P_z in Eq. (3.32) include: (1) the uniform distribution on $\{1, -1\}$; (2) the uniform distribution on $\{1, -1, i, -i\}$; and (3) the uniform distribution on the unit circle in \mathbb{C}^d . The example (1) is where z is a real Rademacher random variable, and in this case the complex TensorSRHT coincides with the real TensorSRHT. Thus, the complex TensorSRHT is a strict generalization of the real TensorSRHT.

We first show that the stacked complex TensorSRHT provides an unbiased approximation of the polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y})^p$. Since the real TensorSRHT is a special case, its unbiasedness follows from this result.

Proposition 3.5.8. *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary, and $\hat{k}_C(\mathbf{x}, \mathbf{y}) = \Phi_C(\mathbf{x})^\top \overline{\Phi_C(\mathbf{y})}$ be the approximate kernel with $\Phi_C(\mathbf{x}), \Phi_C(\mathbf{y}) \in \mathbb{C}^D$ given by the random feature map in Eq. (3.34). Then we have $\mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = (\mathbf{x}^\top \mathbf{y})^p$.*

Proof. We first show $\mathbb{E}[\mathbf{s}_{i,\ell} \overline{\mathbf{s}_{i,\ell}}^\top] = \mathbf{I}_d$ for all $i = 1, \dots, p$ and $\ell = 1, \dots, D$. This follows from the fact that, for all $t, u = 1, \dots, d$, we have

$$\mathbb{E}[(\mathbf{s}_{i,\ell} \overline{\mathbf{s}_{i,\ell}}^\top)_{tu}] = \mathbb{E}[z_{i,t}h_{\pi(\ell),t} \overline{z_{i,u}h_{\pi(\ell),u}}] = \begin{cases} \mathbb{E}[|z_{i,t}|^2] \mathbb{E}[h_{\pi(\ell),t}^2] = 1 & (\text{if } t = u), \\ \mathbb{E}[z_{i,t}] \mathbb{E}[\overline{z_{i,u}}] \mathbb{E}[h_{\pi(\ell),t} h_{\pi(\ell),u}] = 0 & (\text{if } t \neq u), \end{cases}$$

Using this, we have

$$\mathbb{E} \left[\Phi_C(\mathbf{x})^\top \overline{\Phi_C(\mathbf{y})} \right] = \mathbb{E} \left[\frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p \mathbf{x}^\top \mathbf{s}_{i,\ell} \overline{\mathbf{s}_{i,\ell}}^\top \mathbf{y} \right] = \frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p \mathbf{x}^\top \mathbb{E}[\mathbf{s}_{i,\ell} \overline{\mathbf{s}_{i,\ell}}^\top] \mathbf{y} = (\mathbf{x}^\top \mathbf{y})^p$$

■

The proof for the unbiasedness of complex upsampled TensorSRHT is analogous.

We now study the variance of the approximate kernel given by complex stacked TensorSRHT in Eq. (3.34) as well as complex upsampled TensorSRHT. To this end, we use the same notation as Theorem 3.3.1 to write the real and imaginary parts of the random variable z as $z = a + ib$ with real-valued random variables $a, b \in \mathbb{R}$. The proof of the following theorem is provided in Appendix B.2.2.

Theorem 3.5.9 (Variance of Complex Stacked TensorSRHT). *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ be arbitrary, and $\hat{k}_c(\mathbf{x}, \mathbf{y}) = \Phi_c(\mathbf{x})^\top \overline{\Phi_c(\mathbf{y})}$ be the approximate kernel with $\Phi_c(\mathbf{x}), \Phi_c(\mathbf{y}) \in \mathbb{C}^D$ given by the complex random feature map in Eq. (3.34). For the random variable z defining Eq. (3.32), write $z = a + ib$ with $a, b \in \mathbb{R}$, and suppose that*

$$\mathbb{E}[ab] = 0, \quad \mathbb{E}[a^2] = q, \quad \mathbb{E}[b^2] = 1 - q \quad \text{where } 0 \leq q \leq 1.$$

Then we have

$$\mathbb{V}[\hat{k}_c(\mathbf{x}, \mathbf{y})] = \underbrace{\frac{V_q^{(p)}}{D}}_{(A)} - \underbrace{\frac{c(D, d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^p \right]}_{(B)}, \quad (3.35)$$

for complex stacked TensorSRHT, where $V_q^{(p)} \geq 0$ and $V_q^{(1)} \geq 0$ are Eq. (3.17) with the considered value of p and $p = 1$, respectively, and $c(D, d) \in \mathbb{N}$ is defined in (3.29).

For complex upsampled TensorSRHT, we have

$$\mathbb{V}[\hat{k}_c(\mathbf{x}, \mathbf{y})] = \underbrace{\frac{V_q^{(p)}}{D}}_{(A)} - \underbrace{\left(1 - \frac{1}{D}\right) \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{\lceil D/d \rceil d - 1} \right)^p \right]}_{(B)}. \quad (3.36)$$

Regarding Theorem 3.5.9, we make the following observations.

- The case $q = 1$ recovers Theorem 3.5.2 and Theorem 3.5.6 on the real TensorSRHT, where $z \in \{1, -1\}$ is a Rademacher random variable. The case $q = 1/2$ is the complex TensorSRHT with, for instance, P_z being the uniform distribution on $\{1, -1, i, -i\}$ or on the unit circle in \mathbb{C} . Other values of $q \in [0, 1]$ can also be considered, but we do not discuss them further.
- The first term (A) in Eq. (3.35) and Eq. (3.36) is the variance of the unstructured polynomial sketch in Eq. (3.10) with D features, since $V_q^{(p)}$ is its variance with a single feature ($D = 1$) in Eq. (3.17). The second term (B) in Eq. (3.35) and Eq. (3.36) is the effect of using the structured sketch. The quantity $V_q^{(1)}$ is the variance of the unstructured sketch in Eq. (3.10) with a single feature in Eq. (3.17) with $p = 1$.

- As for the real case, the variance (3.35) becomes zero when $p = 1$ and $D \in \{kd \mid k \in \mathbb{N}\}$.

As we discussed for the real TensorSRHT in Corollary 3.5.4, Theorem 3.5.9 enables understanding a condition under which the complex TensorSRHT is advantageous over the corresponding unstructured complex sketch in Eq. (3.10). As for the real case, the condition is that the degree p of the polynomial kernel is *odd*, as stated in the following.

Corollary 3.5.10. *Let $p \in \mathbb{N}$ be odd. Then, for all input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the variance of the approximate kernel with the complex stacked TensorSRHT in Eq. (3.35) is smaller or equal to the variance of the approximate kernel with the corresponding unstructured polynomial sketch:*

$$\frac{V_q^{(p)}}{D} - \frac{c(D, d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^p \right] \leq \frac{V_q^{(p)}}{D}$$

Proof. Since, $V_q^{(1)} \geq 0$, we have $(\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} V_q^{(1)} \leq (\mathbf{x}^\top \mathbf{y})^2$. For odd p this leads to $\left((\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} V_q^{(1)} \right)^p \leq (\mathbf{x}^\top \mathbf{y})^{2p}$. The assertion immediately follows. ■

As discussed for the real case, if p is even, the variance of the complex TensorSRHT can be larger than the corresponding unstructured sketch for certain input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ (e.g., when $\mathbf{x}^\top \mathbf{y} = 0$). Empirically, however, the complex TensorSRHT often has a smaller variance also for even p , as we demonstrate later. Moreover, the arguments of Corollary 3.5.10 can also be applied to upsampled TensorSRHT.

3.5.3 Comparing the Real and Complex TensorSRHT

Let us now compare the real and complex TensorSRHT. We focus on a comparison of stacked TensorSRHT here although the same arguments can be applied to upsampled TensorSRHT. To make the discussion clearer, suppose that the number of random features satisfies $D = Bd$ for some $B \in \mathbb{N}$, as in Remark 3.5.3. Then the variance formula in Eq. (3.35) simplifies to

$$\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] = \underbrace{\frac{V_q^{(p)}}{D}}_{(A)} - \frac{d-1}{D} \underbrace{\left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^p \right]}_{(B)}. \quad (3.37)$$

Recall that setting $q = 1$ and $q = 1/2$ recover the variances of real and complex TensorSRHT, respectively. Thus, let us compare these two cases. We make the following observations:

- As discussed in Section 3.3, it holds that $V_{1/2}^{(p)} \leq V_1^{(p)}$ and $V_{1/2}^{(1)} \leq V_1^{(1)}$ given that the input vectors $\mathbf{x} = (x_1, \dots, x_d)$, $\mathbf{y} = (y_1, \dots, y_d)$ satisfy the inequality in Eq. (3.20), i.e., $\sum_{i \neq j} x_i x_j y_i y_j \geq 0$, which is satisfied when \mathbf{x} and \mathbf{y} are non-negative vectors.
- Thus, if Eq. (3.20) is satisfied, the first term (A) becomes smaller for $q = 1/2$ (complex case) than $q = 1$ (real case). On the other hand, if p is odd, the second term (B) becomes smaller for $q = 1/2$ than $q = 1$; thus, the variance reduction (i.e., $-(B)$) is smaller for $q = 1/2$ than $q = 1$.

The above observations suggest that, even when Eq. (3.20) is satisfied, whether the complex TensorSRHT ($q = 1/2$) has a smaller variance than the real TensorSRHT ($q = 1$) depends on the balance between the two terms (A) and (B) and on the properties of the input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. We have not been able to provide a theoretical characterization of exact situations where the complex TensorSRHT has a smaller variance than the real TensorSRHT.

To complement the lack of a theoretical characterization, we performed experiments to compare the variances of real and complex TensorSRHT, whose results are shown in Fig. 3.3. We evaluated the variance formula in Eq. (3.37) for $q = 1$ (real) and $q = 1/2$ (complex), for 1000 pairs of input vectors \mathbf{x}, \mathbf{y} randomly sampled from a given dataset (EEG, CIFAR 10 ResNet34 features, MNIST and Gisette). For each pair \mathbf{x}, \mathbf{y} , we computed the ratio of Eq. (3.37) with $q = 1/2$ divided by Eq. (3.37) with $q = 1$, and Fig. 3.3 shows the empirical cumulative distribution function of this ratio for the 4 datasets. In these datasets, the input vectors are nonnegative.

Fig. 3.3 shows that, for 100%, 100%, 97.8%, 100% of the cases of the 4 datasets, respectively, the variance of the complex TensorSRHT is smaller than that of the real TensorSRHT. Moreover, the ratio of the variances tends to be even smaller for a larger value of p . These results suggest that the complex TensorSRHT is effective in reducing the variance of the real TensorSRHT, and the variance reduction is more significant for a larger value p of the polynomial degree. We leave a theoretical analysis for explaining this improvement of the complex TensorSRHT for a future work.

So far, we have compared the variances of real and complex Product-Sketches. In the following, we compare complex Rademacher and TensorSRHT sketches against their real analogs in terms of their *downstream task performance* for the task of Gaussian process classification.

3.5.4 Wall-Clock Time Comparison of Real and Complex Random Features in GP Classification

In this section, we empirically compare the complex Rademacher and TensorSRHT sketches that we introduced in this chapter against their real analogs, where we measure classification error against wall-clock time for the task of Gaussian process

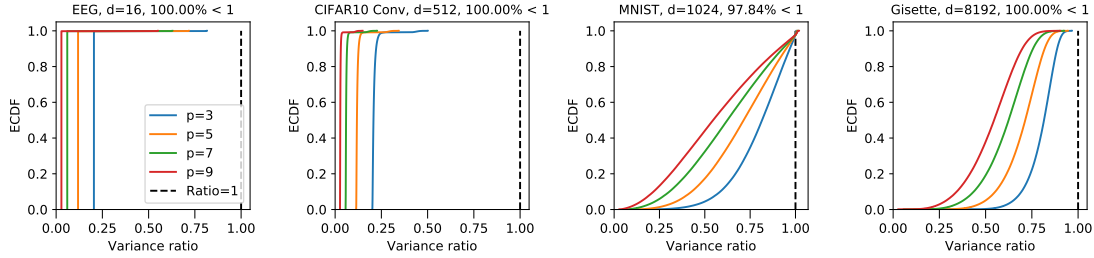


Figure 3.3: Empirical cumulative distribution of pairwise ratios $\text{Var}(\text{Compl. TensorSRHT}) / \text{Var}(\text{Real TensorSRHT})$ on a subsample (1000 samples) of four different datasets (EEG, CIFAR10 ResNet34 features, MNIST, Gisetite) with unit-normalized data where $D = d$. The datasets are not zero-centered and therefore entirely positive.

classification. This comparison is relevant because complex-valued random feature maps may incur a higher computational cost for the downstream task. We therefore investigate whether they lead to a better downstream task performance using the *same* computational budget as their real-valued counterparts. For completeness, we explain how to use complex-valued random features in Gaussian process inference and discuss the resulting computational costs in [Appendix B.3](#)⁴.

We consider Gaussian process classification using the polynomial kernel

$$k(\mathbf{x}, \mathbf{y}) = \sigma^2 \left(\left(1 - \frac{2}{a^2} \right) + \frac{2}{a^2} \mathbf{x}^\top \mathbf{y} \right)^p = \sigma^2 \left(1 - \frac{\|\mathbf{x} - \mathbf{y}\|^2}{a^2} \right)^p \quad (3.38)$$

with $p \in \mathbb{N}$, $a = 4$, $\sigma^2 > 0$, and $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$.

Setting. We use the Rademacher sketch and stacked TensorSRHT, and their respective complex versions. For each polynomial sketch, we compute the KL divergence ([B.23](#)) between the approximate and exact GP posteriors (see [Appendix B.3](#) for details), and record wall-clock time (in seconds) spent on constructing random features and on computing the approximate GP posterior.⁵ We use FashionMNIST ([Xiao et al., 2017](#)) for this experiment.

Results. [Fig. 3.4](#) describes the results. The approximate GPs using complex random features achieve equal or lower KL-divergences than those using real fea-

⁴Specifically, if one uses D complex features, then the inversion of the matrix in [Eq. \(B.20\)](#) requires 4 times as many floating point operations as the case of using D real features. Note that, if one instead uses $2D$ real features, then the inversion of the matrix in [Eq. \(B.20\)](#) requires 8 times as many operations as the case of using D real features. Thus, doubling the number of real features is 2 times more expensive than using complex features. See [Appendix B.3.3](#) for details.

⁵We recorded the time measurements on an NVIDIA P100 GPU and PyTorch version 1.10 with native complex linear algebra support.

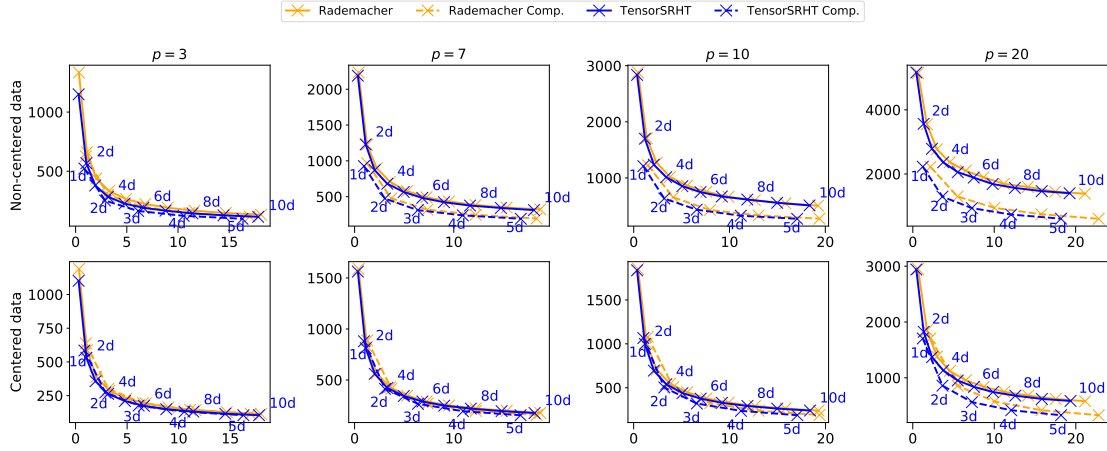


Figure 3.4: Results of the experiments in Section 3.5.4 on wall-clock time comparison of real and complex random features in GP classification on FashionMNIST. In each plot, the vertical axis shows the KL divergence (B.23) between the approximate and the exact GP posteriors for each polynomial sketch, and the horizontal axis is wall-clock time (in seconds) spent on constructing random features and on computing the approximate GP posterior. Each column corresponds to a different degree $p \in \{3, 7, 10, 20\}$ of the polynomial kernel in Eq. (3.38). The top row shows results on the non-centered (thus non-negative) data, and the bottom row to those on the zero-centered data. The number of random features is $D \in \{1d, \dots, 10d\}$ for real features, and $D \in \{1d, \dots, 5d\}$ for complex features, as annotated next to the respective measurements in each plot.

tures of the same computation time, for all the cases. In particular, the improvements of complex features are larger for higher polynomial degrees p and for the non-centered (and thus non-negative) data. This observations agrees with the discussion in Section 3.3 on when complex features yield lower variances than real features.

Chapter 4

Complex-to-Real Product-Sketches

In the previous chapter, we introduced complex Product-Sketches and analyzed their implications on the variances of the corresponding kernel estimate. Although complex polynomial sketches can yield much lower variances than real-valued analogs, they require the downstream model to handle complex data, which usually leads to an increased computational cost as we show for the example of Gaussian processes in Section B.3.

In this chapter, we show how to transform complex polynomial sketches into real ones that we call *Complex-to-Real (CtR)* polynomial sketches from now on. CtR-Sketches do not require the downstream model to handle complex data and can be used as a drop-in replacement inside any model that makes use of random feature approximations. We show that CtR-sketches maintain the variance reduction properties of complex Product-Sketches, thus leading to an even better downstream performance under a fixed computational budget.

4.1 Transforming Complex into Real Sketches

In the following, we show how to convert complex polynomial sketches into real-valued CtR-sketches. Let $\Phi_{\mathcal{C}} : \mathbb{R}^d \rightarrow \mathbb{C}^D$ be a complex polynomial sketch and $\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{C}}(\mathbf{x})^\top \overline{\Phi_{\mathcal{C}}(\mathbf{y})}$ the approximate kernel as defined in Chapter 3, where we make the use of complex projections explicit through the subscript \mathcal{C} . We further denote by $\text{Re}\{\cdot\}$ and $\text{Im}\{\cdot\}$ the real and imaginary parts of a complex vector.

$\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})$ is generally complex-valued and can hence be written as

$$\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y}) = \text{Re}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\} + i \cdot \text{Im}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}.$$

For an unbiased polynomial sketch, we have $\mathbb{E}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] = k(\mathbf{x}, \mathbf{y}) + 0 \cdot i$. From this it follows that $\mathbb{E}[\text{Re}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}] = k(\mathbf{x}, \mathbf{y})$ and $\mathbb{E}[\text{Im}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}] = 0$ through the linearity of the expectation.

We thus define $\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}) := \text{Re}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}$ to be our novel kernel estimate which, by expanding the real part of $\Phi_{\mathcal{C}}(\mathbf{x})^\top \overline{\Phi_{\mathcal{C}}(\mathbf{y})}$, can be written as

$$\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}) = \text{Re}\{\Phi_{\mathcal{C}}(\mathbf{x})\}^\top \text{Re}\{\Phi_{\mathcal{C}}(\mathbf{y})\} + \text{Im}\{\Phi_{\mathcal{C}}(\mathbf{x})\}^\top \text{Im}\{\Phi_{\mathcal{C}}(\mathbf{y})\}.$$

Note that $\text{Re}\{\Phi_{\mathcal{C}}(\mathbf{x})\}, \text{Im}\{\Phi_{\mathcal{C}}(\mathbf{x})\} \in \mathbb{R}^D$, which allows us to define a $2D$ -dimensional *real-valued* polynomial sketch

$$\begin{aligned} \Phi_{\text{CtR}}(\mathbf{x}) &:= (\text{Re}\{\Phi_{\mathcal{C}}(\mathbf{x})_1\}, \dots, \text{Re}\{\Phi_{\mathcal{C}}(\mathbf{x})_D\}, \\ &\quad \text{Im}\{\Phi_{\mathcal{C}}(\mathbf{x})_1\}, \dots, \text{Im}\{\Phi_{\mathcal{C}}(\mathbf{x})_D\})^\top \in \mathbb{R}^{2D}, \end{aligned} \quad (4.1)$$

for which we have

$$\Phi_{\text{CtR}}(\mathbf{x})^\top \Phi_{\text{CtR}}(\mathbf{y}) = \text{Re}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\} = \hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}). \quad (4.2)$$

We call Φ_{CtR} a *Complex-to-Real (CtR)* polynomial sketch and summarize its construction in [Algorithm 3](#), where we assume $\Phi_{\mathcal{C}}$ to be one of the Product-Sketches defined in [Chapter 3](#).

The dimension of the feature map Φ_{CtR} [\(4.1\)](#) is $2D$ for a given D random samples, while the dimension of the real and complex feature maps, $\Phi_{\mathcal{R}}$ and $\Phi_{\mathcal{C}}$, discussed in the previous chapter is only D . Feature maps with higher dimensions usually render the downstream task more expensive and the key question that we address in this chapter is thus: Does the CtR estimator in [Eq. \(4.2\)](#) yield lower variances than $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y}) = \Phi_{\mathcal{R}}(\mathbf{x})^\top \Phi_{\mathcal{R}}(\mathbf{y})$, where $\Phi_{\mathcal{R}}$ is the corresponding real-valued Product-Sketch defined in [Chapter 3](#), if Φ_{CtR} and $\Phi_{\mathcal{R}}$ have the *same* output dimension?

We will show in [Section 4.2.1](#) that this is indeed the case meaning that CtR-sketches inherit the variance reduction properties of complex Product-Sketches. This result implies that, under the conditions provided in this work, CtR-Product-Sketches should be preferred over traditional real-valued Product-Sketches contained in ([Kar and Karnick, 2012](#); [Hamid et al., 2014](#)). Unlike the complex sketches discussed in the previous chapter, CtR-sketches are real-valued and do not require the downstream model to handle complex data.

4.2 Variances of CtR Product-Sketches

We will now study the variances of CtR-Product-Sketches. Obtaining a closed form expression for their variances seems challenging at first sight, as we need to evaluate

$$\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] = \mathbb{V} \left[\frac{1}{D} \sum_{\ell=1}^D \text{Re} \left\{ \prod_{i=1}^p (\mathbf{w}_{i,\ell}^\top \mathbf{x}) \overline{(\mathbf{w}_{i,\ell}^\top \mathbf{y})} \right\} \right],$$

where the weights $\mathbf{W}_i = (\mathbf{w}_{i,1}, \dots, \mathbf{w}_{i,D})^\top \in \mathbb{C}^{D \times d}$ are defined through the respective sketch as shown in [Algorithm 3](#).

Algorithm 3: Complex-to-Real (CtR) Random Features**Result:** A feature map $\Phi_{\text{CtR}}(\mathbf{x})$ **Input:** Datapoint $\mathbf{x} \in \mathbb{R}^d$, projection dimension $D \in \mathbb{N}$, degree $p \in \mathbb{N}$;
Sample $\{\mathbf{W}_i\}_{i=1}^p$ with $\mathbf{W}_i \in \mathbb{C}^{D \times d}$ independently according to one of the following sketches:

- Gaussian: $(\mathbf{W}_i)_{\ell,k} \stackrel{i.i.d.}{\sim} \mathcal{CN}(0, 1)$
- Rademacher: $(\mathbf{W}_i)_{\ell,k} \stackrel{i.i.d.}{\sim} \text{Unif}(\{1, -1, i, -i\})$
- Stacked/Upsampled TensorSRHT: $\mathbf{W}_i = (\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,D})^\top$ with $\{\mathbf{s}_{i,\ell}\}_{\ell=1}^D$ defined in Eq. (3.34) in Section 3

Compute $\Phi_C(\mathbf{x}) := (\mathbf{W}_1 \mathbf{x} \odot \dots \odot \mathbf{W}_p \mathbf{x}) / \sqrt{D}$;Set $\Phi_{\text{CtR}}(\mathbf{x}) := (\text{Re}\{\Phi_C(\mathbf{x})_1\}, \dots, \text{Re}\{\Phi_C(\mathbf{x})_D\},$ $\text{Im}\{\Phi_C(\mathbf{x})_1\}, \dots, \text{Im}\{\Phi_C(\mathbf{x})_D\})^\top \in \mathbb{R}^{2D}$;

The challenge here is that the real part of the product of complex numbers corresponds to a large sum and we are not allowed to “swap” $\text{Re}(\cdot)$ with this product. However, we prove the following proposition that simplifies this problem and allows us to make use of the variances for complex Product-Sketches that we derived previously in Chapter 3.

Proposition 4.2.1. *Let $\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}) = \text{Re}\{\hat{k}_C(\mathbf{x}, \mathbf{y})\}$ for some complex-valued kernel estimate $\hat{k}_C(\mathbf{x}, \mathbf{y})$ with $\mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = k(\mathbf{x}, \mathbf{y})$. Then the following holds:*

$$\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] = \frac{1}{2} \text{Re} \left\{ \mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] + \mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] \right\} \quad (4.3)$$

$$\text{with } \mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] := \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] - (\mathbf{x}^\top \mathbf{y})^{2p}$$

being the pseudo-variance of $\hat{k}_C(\mathbf{x}, \mathbf{y})$ and $\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})]$ its variance.

Proof. For a complex random variable $z = a + ib$ with $a, b \in \mathbb{R}$, we have $|z|^2 = a^2 + b^2$ and $\text{Re}\{z^2\} = a^2 - b^2$. Combining both equations gives $a^2 = \frac{1}{2}(|z|^2 + \text{Re}\{z^2\})$. The scalar a is real-valued and its variance $\mathbb{V}[a] = \mathbb{E}[a^2] - \mathbb{E}[a]^2$ is thus

$$\mathbb{V}[a] = \frac{1}{2} \text{Re}\{\mathbb{E}[|z|^2] + \mathbb{E}[z^2] - 2\mathbb{E}[a]^2\}. \quad (4.4)$$

Now we set $z = \hat{k}_C(\mathbf{x}, \mathbf{y})$, from which it follows that $a = \text{Re}\{z\} = \hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})$. Eq. (4.4) can then be rewritten as

$$\begin{aligned} \mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] &= \frac{1}{2} \text{Re}\{\mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2] + \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] - 2\mathbb{E}[\text{Re}\{\hat{k}_C(\mathbf{x}, \mathbf{y})\}]^2\} \\ &= \frac{1}{2} \text{Re}\{\mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2] + \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] - 2\mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})]^2\} \\ &= \frac{1}{2} \text{Re}\{\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] + \mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})]\}, \end{aligned}$$

where $\mathbb{P}\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})] := \mathbb{E}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})^2] - \mathbb{E}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})] \in \mathbb{C}$ is called the *pseudo-variance* of $\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})$ (Park, 2018, Chapter 5). \blacksquare

Remark 4.2.2. *In fact, we show next that $\text{Im}\{\mathbb{P}\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]\} = 0$ for all the polynomial sketches discussed in this work. Hence, we can also write $\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] = \frac{1}{2}(\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})] + \mathbb{P}\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})])$ for them since $\mathbb{V}[z] \in \mathbb{R}$ for any $z \in \mathbb{C}$.*

In order to determine $\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})]$, we thus only need to know $\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]$ and $\mathbb{P}\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]$ for Gaussian, Rademacher and TensorSRHT sketches discussed in Algorithm 3. As we already derived $\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]$ for all these sketches in the previous chapter, the remaining task is to derive the pseudo-variances $\mathbb{P}\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]$, which can be easily done using very similar techniques that we already used for the derivation of the variances.

We derive the Pseudo-Variances of Gaussian and Rademacher sketches in Section C.1 and the ones for upsampled and stacked TensorSRHT in Section C.3. We summarize these novel results in the lower half of Table 4.1. For completeness, we add the variances of real and complex polynomial sketches from Chapter 3 $\Phi_{\mathcal{R}}$ to the upper half of Table 4.1. The variance of any CtR-Sketch discussed in this work can thus be obtained by substituting the corresponding variance and pseudo-variance expressions from Table 4.1 into Eq. (4.3).

Having worked out the variances of CtR-Product-Sketches allows us to compare these sketches to their real-valued analogs in the following and to draw conclusions when which sketch should be preferred.

4.2.1 Variance Comparison for Gaussian and Rademacher CtR-Sketches

We begin by studying the variance reduction properties of Rademacher and Gaussian CtR-sketches over their real-valued analogs that were originally proposed by Kar and Karnick (2012); Hamid et al. (2014).

Let $\Phi_{\mathcal{R}} : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$ (using $2D$ random features) be a real polynomial sketch as defined in Chapter 3 and $\Phi_{\text{CtR}} : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$ a CtR-sketch (using only D random features) as defined in Eq. (4.1). Let $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})$ and $\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})$ be the respective approximate kernels for some $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Then we can provide the following theorem for the case of Rademacher sketches.

Theorem 4.2.3 (CtR-Rademacher advantage). *Let $a := \sum_{i=1}^d \sum_{j \neq i}^d x_i x_j y_i y_j$ and $b(j) := \|\mathbf{x}\|^{2j} \|\mathbf{y}\|^{2j} - (\sum_{i=1}^d x_i^2 y_i^2)^j \geq 0$. Then $\mathbb{V}[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})]$ yields*

$$\frac{1}{2D} \sum_{k=2}^p \sum_{j=0}^{k-1} \binom{p}{k} \binom{k}{j} b(j) a^{p-j} \geq 0 \quad \text{if } a \geq 0.$$

Furthermore, CtR-Rademacher sketches achieve the lowest possible variance for $\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}) = \text{Re}\{\Phi_{\mathbb{C}}(\mathbf{x})^{\top} \overline{\Phi_{\mathbb{C}}(\mathbf{y})}\}$ with $\Phi_{\mathbb{C}}$ being defined through Eq. (3.10) in Chapter 3.

Sketch	Variance $\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]$ ($q = 1$) and $\mathbb{V}[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})]$ ($q = 2$)
Gaussian	$D^{-1}[(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + q(\mathbf{x}^\top \mathbf{y})^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p}]$
Rademacher	$D^{-1}[(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + q((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2))^p - (\mathbf{x}^\top \mathbf{y})^{2p}]$
TensorSRHT (stacked)	$\mathbb{V}_{\text{Rad.}}^{(p)} - (c(D, d)/D^2) \cdot [(\mathbf{x}^\top \mathbf{y})^{2p} - (\text{C}_{\text{Var.}})^p]$ with $\text{C}_{\text{Var.}} = (\mathbf{x}^\top \mathbf{y})^2 - (d-1)^{-1} \mathbb{V}_{\text{Rad.}}^{(1)}$
TensorSRHT (upsampled)	$\mathbb{V}_{\text{Rad.}}^{(p)} - (1 - 1/D) \cdot [(\mathbf{x}^\top \mathbf{y})^{2p} - (\text{C}_{\text{Var.}})^p]$ with $\text{C}_{\text{Var.}} = (\mathbf{x}^\top \mathbf{y})^2 - (\lceil D/d \rceil d - 1)^{-1} \mathbb{V}_{\text{Rad.}}^{(1)}$
Sketch	Pseudo-Variance $\mathbb{P}\mathbb{V}[\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})]$
Gaussian	$D^{-1}[(2(\mathbf{x}^\top \mathbf{y})^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p}]$
Rademacher	$D^{-1}[(2(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p}]$
TensorSRHT (stacked)	$\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)} - (c(D, d)/D^2) \cdot [(\mathbf{x}^\top \mathbf{y})^{2p} - (\text{C}_{\text{PVar.}})^p]$ with $\text{C}_{\text{PVar.}} = (\mathbf{x}^\top \mathbf{y})^2 - (d-1)^{-1} \mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)}$
TensorSRHT (upsampled)	$\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)} - (1 - 1/D) \cdot [(\mathbf{x}^\top \mathbf{y})^{2p} - (\text{C}_{\text{PVar.}})^p]$ with $\text{C}_{\text{PVar.}} = (\mathbf{x}^\top \mathbf{y})^2 - (\lceil D/d \rceil d - 1)^{-1} \mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)}$

Table 4.1: Variances of complex $\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})$ and real $\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})$ as well as pseudo-variances of complex $\hat{k}_{\mathbb{C}}(\mathbf{x}, \mathbf{y})$ are shown.

$\mathbb{V}_{\text{Rad.}}^{(p)}$, $\mathbb{V}_{\text{Rad.}}^{(1)}$ and $\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)}$, $\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)}$ are the Rademacher variances and pseudo-variances for a given degree p and $p = 1$, respectively.

$$c(D, d) := \lfloor D/d \rfloor d(d-1) + (D \bmod d)(D \bmod d - 1).$$

Proof. The variance reduction property is proven in Section C.2 and we show in Section C.1 that Rademacher sketches achieve the lowest possible pseudo-variance for the estimator defined through Eq. (3.10). Since they also achieve the lowest possible variance as shown in Theorem 3.3.1, this implies that they achieve the lowest CtR-variance. \blacksquare

Theorem 4.2.3 tells us that Φ_{CtR} should be preferred over $\Phi_{\mathcal{R}}$ when $a \geq 0$ for two given inputs $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and the variance gap increases as p increases. This is the same condition under which complex Rademacher sketches have lower variances than real ones (see Eq. (3.20) in Chapter 3), which holds, e.g., when \mathbf{x}, \mathbf{y} are non-negative. Even when this condition is not always met, CtR-Rademacher sketches perform similarly or better than real ones in our experiments in Section 4.3.1.

We can additionally provide the following theorem for Gaussian polynomial sketches proved in Section C.2.

Theorem 4.2.4 (CtR-Gaussian advantage). *For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\mathbb{V}[\hat{k}_{\mathcal{R}}(\mathbf{x}, \mathbf{y})] -$*

$\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})]$ yields

$$\frac{1}{2D} \sum_{k=0}^{p-1} \binom{p}{k} (2^k - 1) (\mathbf{x}^\top \mathbf{y})^{2k} (\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^{p-k} \geq 0.$$

Thus, regardless of the input data, Φ_{CtR} should be preferred over $\Phi_{\mathcal{R}}$ when using Gaussian polynomial sketches. The advantage again increases with p .

The results of Theorems 4.2.3 and 4.2.4 are remarkable because they show conditions under which the CtR-Gaussian and Rademacher sketches outperform the corresponding real analogs defined through Eq. (3.1) in Chapter 3, when both sketches use the *same* feature map dimension. CtR-Rademacher can thus be more efficient than Rademacher sketches with real-valued weights that already achieve a variance lower bound in Theorem 3.1.1. We now give an intuition for this variance reduction.

Intuition for the variance reduction. To obtain an intuition for the variance reduction of CtR Product-Sketches, we look at the variance structure of complex Product-Sketches first. We assume a single random feature $D = 1$ for simplicity here. The variance of the approximate kernel $\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})$ (3.11) can be written as:

$$\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] = \mathbb{V} \left[\prod_{i=1}^p (\mathbf{w}_i^\top \mathbf{x})(\mathbf{w}_i^\top \mathbf{y}) \right] = \prod_{i=1}^p \underbrace{\mathbb{E} [|(\mathbf{w}_i^\top \mathbf{x})(\mathbf{w}_i^\top \mathbf{y})|^2]}_{=: \mu_{2,\mathcal{C}}} - (\mathbf{x}^\top \mathbf{y})^{2p} \quad (4.5)$$

The value of $\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})]$ thus entirely depends on the second complex moment $\mu_{2,\mathcal{C}} \geq 0$. As shown in Section B.1.1, this second moment is smaller for complex Rademacher $\{\mathbf{w}_i\}_{i=1}^p$ than for real ones if $\sum_{i=1}^d \sum_{j=1}^d x_i x_j y_i y_j \geq 0$ holds. If we thus let $\mu_{2,\mathcal{C}} = c\mu_{2,\mathcal{R}}$ for some $c < 1$ in this case, we get $\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] = c^p \mu_{2,\mathcal{R}} - (\mathbf{x}^\top \mathbf{y})^{2p}$, where $\mu_{2,\mathcal{R}}$ denotes the respective second moment when $\{\mathbf{w}_i\}_{i=1}^p$ are real-valued. We thus obtain a variance advantage that increases with p .

This advantage is maintained for CtR-sketches because for any complex random variable, the relationship

$$\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] = \mathbb{V}[\text{Re}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}] + \mathbb{V}[\text{Im}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}]$$

and therefore $\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] \leq \mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})]$ hold.

Interestingly, when $p = 1$, CtR-sketches are equivalent to doubling the number of random features of the analogous sketch with real-valued weights. This can be seen when looking at Eq. (4.1) and recalling that the real and imaginary parts of $\{\mathbf{w}_i\}_{i=1}^p$ are sampled independently in our construction. As the resulting CtR feature map is $2D$ -dimensional, there is no statistical advantage w.r.t. the feature map dimension in this case. However, the product in Eq. (4.5) *increases* the variance advantage, giving CtR-sketches a variance improvement from $p \geq 2$.

We study the variance reduction properties of CtR-TensorSRHT next.

4.2.2 Variance of CtR-TensorSRHT

As we already showed in Corollary 3.5.4 for real TensorSRHT and Corollary 3.5.10 for complex TensorSRHT, we will now show that CtR-TensorSRHT can achieve a variance reduction over their non-structured counterpart (CtR-Rademacher sketches) for odd degrees p of the polynomial kernel. The analysis is almost the same as in Corollary 3.5.10.

Corollary 4.2.5. *Let $p \in \mathbb{N}$ be odd. Then, for all input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ with $\sum_{i=1}^d \sum_{j \neq i}^d x_i x_j y_i y_j \geq 0$, the variance of the approximate kernel with the CtR-TensorSRHT sketch in Table 4.1 is smaller or equal to the variance of the approximate kernel with the corresponding unstructured polynomial sketch.*

Proof. First recall from Eq. (4.3) that the variance of CtR-sketches is given by the average $\frac{1}{2} \text{Re}\{\mathbb{V}[\hat{k}_c(\mathbf{x}, \mathbf{y})] + \mathbb{P}\mathbb{V}[\hat{k}_c(\mathbf{x}, \mathbf{y})]\}$, where we can drop $\text{Re}\{\cdot\}$ because both the variances and pseudo-variances in Table 4.1 are already real-valued.

Next, we identify the following reduction term

$$R_{\text{Var./PVar.}} := (1 - 1/D) \left((\mathbf{x}^\top \mathbf{y})^{2p} - (C_{\text{Var./PVar.}})^p \right) \quad (4.6)$$

that is subtracted from the Rademacher variance $\mathbb{V}_{\text{Rad.}}^{(p)}$ and pseudo-variance $\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)}$ in Table 4.1, respectively, where we chose upsampled TensorSRHT here. Stacked TensorSRHT is analogous with $(1 - 1/D)$ being substituted by $c(D, d)/D^2$.

If p is odd, we have $(C_{\text{Var.}})^p \leq (\mathbf{x}^\top \mathbf{y})^{2p}$ because $\mathbb{V}_{\text{Rad.}}^{(1)} \geq 0$ and therefore $R_{\text{Var.}} \geq 0$ holds. In this case, the variance of complex/real TensorSRHT is upper-bounded by the complex/real Rademacher variance $\mathbb{V}_{\text{Rad.}}^{(p)} \geq 0$.

If we further have $\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)} = \sum_{i=1}^d \sum_{j \neq i}^d x_i x_j y_i y_j \geq 0$, the pseudo-variance of complex TensorSRHT is also upper-bounded by the Rademacher pseudo-variance. This is because $0 \leq (C_{\text{PVar.}})^p \leq (\mathbf{x}^\top \mathbf{y})^{2p}$ and $R_{\text{PVar.}} \geq 0$ hold. This completes the proof. \blacksquare

The condition $\sum_{i=1}^d \sum_{j \neq i}^d x_i x_j y_i y_j \geq 0$ is the same for which CtR-Rademacher sketches achieve lower variances than real Rademacher sketches, and we will see empirically in the following, that CtR-TensorSRHT sketches outperform real TensorSRHT sketches under this same condition. Therefore, the conditions for structured sketches to outperform non-structured ones and CtR-sketches outperforming their real analogs are well aligned. Moreover, CtR-TensorSRHT inherits the variance reduction of CtR-Rademacher sketches over their real analogs because the Rademacher variance and pseudo-variance are both included in the ones of complex TensorSRHT (see Table 4.1).

We now move on to comparing the variance of CtR-TensorSRHT against its real-valued analog.

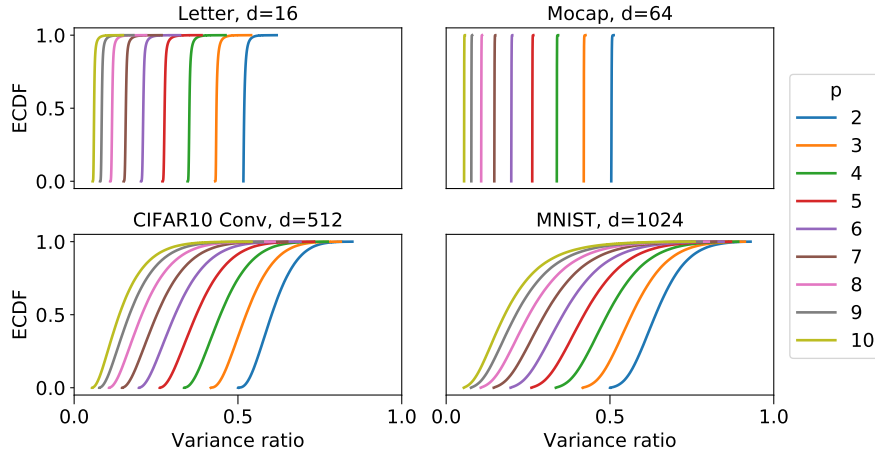


Figure 4.1: $\text{Var}(\text{CtR-TensorSRHT}) / \text{Var}(\text{TensorSRHT})$. The target kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^p$.

Comparing the variances of CtR-TensorSRHT against TensorSRHT.

Analyzing the variance reduction properties of CtR-TensorSRHT over its real analog is very involved as it has already been the case for complex TensorSRHT as we showed in Section 3.5.3. Therefore, we will follow the example of Section 3.5.3 and resort to an empirical comparison here.

Experimental setting. Since we obtained the variances of CtR-TensorSRHT and real TensorSRHT in closed form as shown in Table 4.1, we can carry out an exact empirical comparison of their variance ratio on pairs of data points contained in four example data sets (Letter and Mocap (Dua and Graff, 2017), CIFAR-10¹ (Krizhevsky et al., 2009) and MNIST (Lecun et al., 1998)). All data sets except for Mocap contain only non-negative inputs. For Mocap, we subtract the minimum value of each input dimension $\{x_i\}_{i=1}^d$ of $\mathbf{x} \in \mathbb{R}^d$ across the data set to achieve non-negativity. We also make sure that both feature maps, CtR-TensorSRHT and TensorSRHT use the same feature map dimension $D = 2d$ in this comparison, and we choose upsampled TensorSRHT here.

Results. Fig. 4.1 shows the results of the comparison through empirical cumulative distribution functions of the variance ratios evaluated on all input pairs inside 1000 random samples of each dataset. We observe that CtR-TensorSRHT outperforms real TensorSRHT (all variance ratios < 1) and the improvement increases with p .

Since the variance of TensorSRHT is closely related to Rademacher variances as mentioned before, we study the importance of the non-negativity condition of

¹We use the the convolutional outputs of a pretrained ResNet34 (He et al., 2016) on ImageNet (Russakovsky et al., 2015).

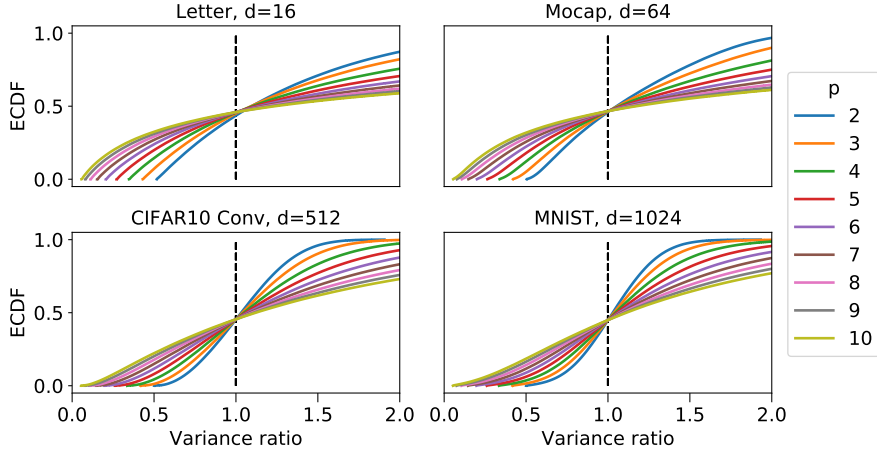


Figure 4.2: $\text{Var}(\text{CtR-TensorSRHT}) / \text{Var}(\text{TensorSRHT})$ on zero-centered data. The target kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^p$.

the data in [Theorem 4.2.3](#) for the case of CtR-TensorSRHT. For this purpose, we repeat the same experiment as before for zero-centered data that breaks the non-negativity condition. Results are shown in [Fig. 4.2](#). We see that CtR-TensorSRHT performs similarly as its real analog in this case (around half of the variance ratios are above and below 1). Therefore, a variance reduction happens in particular for non-negative data as is the case for CtR-Rademacher sketches. This makes sense because both the variance and the pseudo-variance of TensorSRHT are closely related to the ones of Rademacher sketches as shown in [Table 4.1](#).

So far, we have compared both complex and Complex-to-Real Product-Sketches against real Product-Sketches. The last open question is how Complex-to-Real sketches compare against complex sketches. To answer this question, we do *not* focus on a variance comparison between the two, but we consider studying the implications on the downstream task to be more important. This is because complex random feature maps require the downstream task to handle complex data incurring higher computational costs as shown previously in [Section 3.5.3](#) in [Chapter 3](#). We will therefore extend the Gaussian process classification example from [Section 3.5.3](#), and compare real, complex, and Complex-to-Real TensorSRHT sketches on this downstream task.

Comparing Real, Complex and CtR-Sketches. We compare the convergence of a GP classifier using real, complex and CtR-TensorSRHT towards one using the exact kernel on MNIST under a given time budget in [Fig. 4.3](#). We can see that a complex GP classifier using complex TensorSRHT needs roughly the same time for a projection dimension D as real and CtR-TensorSRHT for dimension $2D$. This is because complex TensorSRHT requires solving a complex GP model that takes more operations to solve for the same feature dimension as a real GP as explained [Appendix B.3](#).

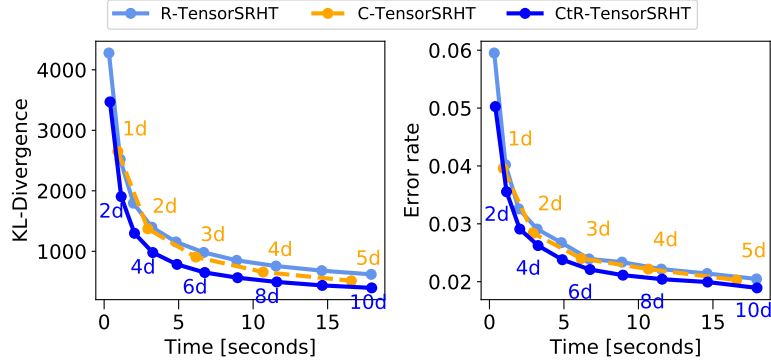


Figure 4.3: KL-divergence (left) and test error (right) for a given computation time budget computed on MNIST ($d = 1024$). Projection dimensions are annotated next to the respective measurements. We used the kernel $(\mathbf{x}^\top \mathbf{y} + 1)^6$.

This gives a clear advantage to the CtR-TensorSRHT sketch that achieves the lowest KL-divergence and test errors in the same amount of time. Therefore, we recommend the use of CtR-TensorSRHT over real and complex TensorSRHT. Here we used the kernel $(\mathbf{x}^\top \mathbf{y} + 1)^6$ and timings were measured on an NVIDIA P100 graphics card and PyTorch 1.10 (Paszke et al., 2019) with native complex linear algebra support.

4.3 Comparison against the State-of-the-Art

In this section, we compare our most efficient Product-Sketch, upsampled CtR-TensorSRHT, against TensorSketch (Pham and Pagh, 2013) that we introduced in Chapter 2. TensorSketch can be seen as the state-of-the-art *unbiased*² polynomial sketch. Unlike for CtR-TensorSRHT, the variance of TensorSketch is not available in the literature. We will therefore carry out an empirical variance comparison here. Furthermore, we add a feature construction time benchmark to this comparison, since TensorSketch leverages the fast Fourier transform and is therefore considered to be a particularly fast sketch.

Variance Comparison. We repeat the same experiment as in Fig. 4.1, but this time we evaluate the variance ratios of upsampled CtR-TensorSRHT against TensorSketch. As the variance of TensorSketch is not available in closed form, we

²For high-degree polynomial kernels on the unit-sphere it is still outperformed by biased Spherical Random Features (SRF) (Pennington et al., 2015). However, SRF require a preprocessing step for every polynomial kernel parameterization and only work for data lying on the unit-sphere, obstructing a straightforward application, e.g., in deep learning. We will carry out empirical evaluations against SRF in Chapter 5.

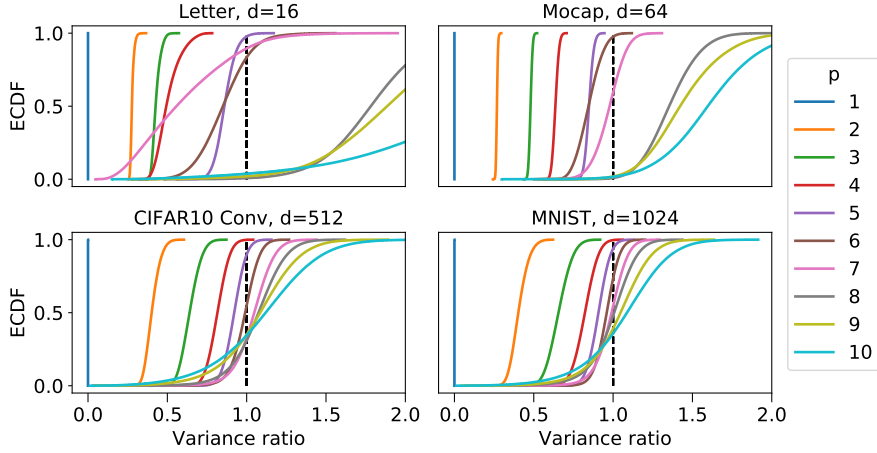


Figure 4.4: $\text{Var}(\text{CtR-TensorSRHT}) / \text{Var}(\text{TensorSketch})$. Empirical cumulative distribution function of variance ratios for feature maps with equal dimension $D = 2d$. The target kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^p$.

estimate it using 1000 Monte-Carlo samples of the kernel estimate. Fig. 4.4 shows a comparison of the variance ratios.

In general, upsampled CtR-TensorSRHT improves over TensorSketch for $p \leq 5$ (more than half of the ratios are less than one). This is an important advantage as such degrees are commonly used in practice (Pham and Pagh, 2013; Gao et al., 2016; Fukui et al., 2016). For $p = 1$, TensorSRHT always has a variance of zero as soon as $D = kd$ for $k \in \mathbb{N}$, which can be derived from the TensorSRHT variance equation in Table 4.1.

For large p , TensorSketch and upsampled CtR-TensorSRHT perform similarly well (around half the variance ratios are smaller/greater than 1) except for some cases with low-dimensional data where TensorSketch performs better. The ratios would be much worse if we used the real TensorSRHT instead as shown in Fig. 4.5, which underlines the relative improvement of the CtR-extension.

The variance comparison of CtR-TensorSRHT against TensorSketch is impressive given the fact that (upsampled) CtR-TensorSRHT has a faster running time, which we will compare next.

Feature construction time comparison. In the following, we carry out a feature construction time comparison of upsampled CtR-TensorSRHT against TensorSketch that has a time complexity $\mathcal{O}(p(D \log D + D))$ with D being the projection dimension.

Recall that our proposed upsampled TensorSRHT approach introduced in Section 3.5.1 in Chapter 3 has a time complexity of $\mathcal{O}(p(d \log d + D))$ and is thus faster in theory when $D > d$.

The left plot in Fig. 4.6 shows the results of this comparison, where we added

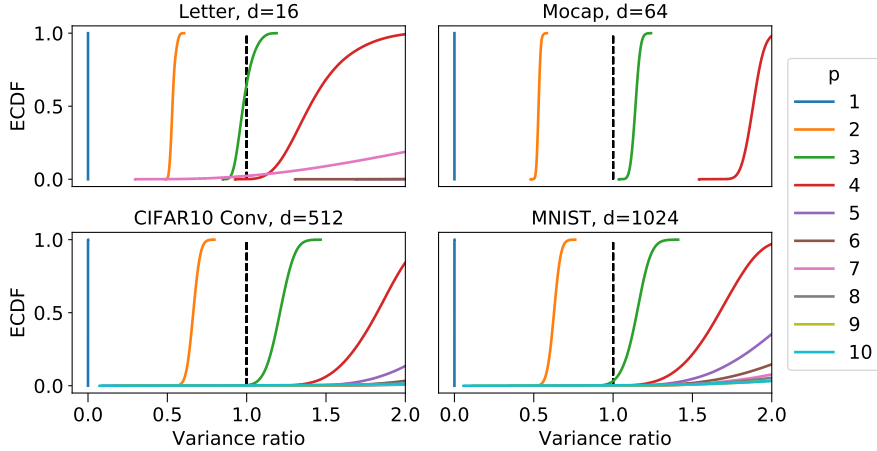


Figure 4.5: $\text{Var}(\text{TensorSRHT}) / \text{Var}(\text{TensorSketch})$. Empirical cumulative distribution function of variance ratios for feature maps with equal dimension $D = 2d$. The target kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^p$.

real Rademacher and upsampled TensorSRHT Product-Sketches to the plot to give a more complete picture. The construction times of real TensorSRHT and CtR-TensorSRHT have a smaller slope with respect to D than the other sketches leading to the lowest feature construction times, in particular when $D \gg d$. There is a small computational overhead for CtR-TensorSRHT compared to real TensorSRHT because CtR-TensorSRHT initially requires two Hadamard-projections (real and imaginary parts), but uses the same sampling matrix leading to the same scaling property with respect to D .

The right plot of Fig. Fig. 4.6 shows the kernel approximation error (relative Frobenius norm error) under a computational budget. We can see that our proposed CtR-TensorSRHT gives the best kernel approximations in the shortest amount of time when D is sufficiently large. As before, timings were measured on an NVIDIA P100 GPU and PyTorch 1.10 (Paszke et al., 2019) with native complex linear algebra support.

One may argue that the feature construction time is not the computational bottleneck in downstream tasks like closed form GP regression, where the bottleneck is the construction and inversion of a D by D matrix. This changes however, when using iterative solvers working with mini-batches. We carry out such a time benchmark for iterative solvers in the following.

Time benchmark for stochastic optimization. We compare CtR-TensorSRHT against TensorSketch for a GP classification experiment using stochastic variational inference and mini-batching. The goal is to measure the impact of the feature construction time on the training time of the learning algorithm. The feature construction time plays a crucial role here since every iteration of the algo-

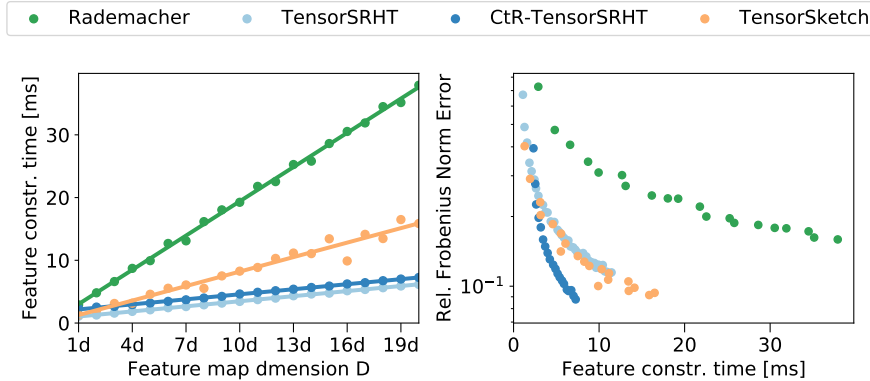


Figure 4.6: (Left) Feature construction time against feature map dimension D . (Right) Kernel approximation error (rel. Frobenius norm error) against feature construction time. We approximate the kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + 1)^6$ on 1000 random MNIST samples.

rithm carries out one forward and one backward pass through a two-layer neural network, where the first layer corresponds to the feature extraction layer, either TensorSRHT or TensorSketch in this case.

We choose $D = 10d$, where $d = 1024$ for the MNIST data set. We further set $p = 6$ as the polynomial degree implying the use of 6 random matrix multiplications. In such a scenario, the feature construction time has a big share on the total computation time of every iteration.

Fig. 4.7 shows the results of this stochastic variational inference experiment, where the goal is to obtain a factorized Gaussian posterior distribution over the weights of a GP classifier. We show the test error against the training time to reach this error. The feature maps approximate the kernel $(\gamma \mathbf{x}^\top \mathbf{y} + \nu)^6$, where $\gamma, \nu \geq 0$ are optimized through backpropagation. The mini-batch size is chosen to be 1000 and we use the Adam optimizer with learning rate 10^{-3} . We further choose 50 Monte-Carlo samples for the posterior weights at each iteration.

When using CtR-TensorSRHT, the experiments were run for 150 epochs while the ones for TensorSketch were run for 50 epochs. We can thus see that TensorSRHT achieves a much faster convergence (in wall-clock time), which is due to a faster computation of the feature map as shown in Fig. 4.6. CtR-TensorSRHT outperforms TensorSRHT in terms of test error and is only slightly slower.

4.3.1 Systematic Evaluation on Different Data Sets

In this section, we carry out a systematic comparison of the CtR polynomial sketches discussed in this work against their real-valued analogs as well as TensorSketch (Pham and Pagh, 2013) on the same four data sets that we used for

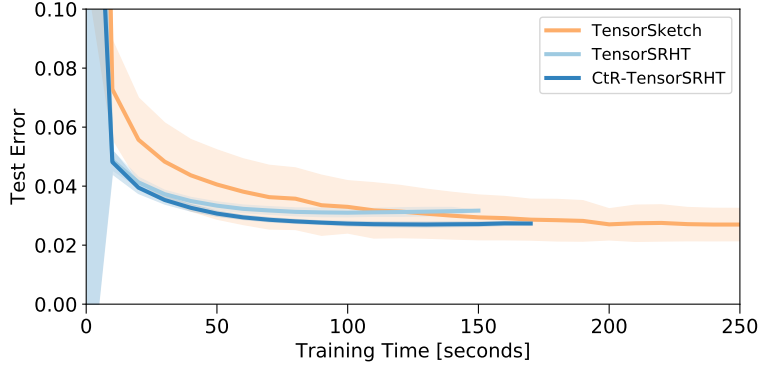


Figure 4.7: Stochastic variational inference on MNIST ($d = 1024$) for $D = 10d$ and $p = 6$. Results are averaged over 10 different runs. TensorSRHT finishes 150 training epochs before TensorSketch finishes 50 epochs and thus converges earlier.

the variance comparison experiments, namely Letter and Mocap (Dua and Graff, 2017), CIFAR-10 (Krizhevsky et al., 2009) and MNIST (Lecun et al., 1998), where we make sure that the data is positive-valued. This comparison includes the evaluation of kernel approximation errors as well as the downstream performance on the task of GP classification.

Target kernel and its approximation. The target polynomial kernel is always $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + 1)^{p^3}$ and we unit-normalize the data to make the kernel bounded and improve numerical stability. Since CtR polynomial sketches Φ_{CtR} are $2D$ -dimensional when sampling D random features, we compare them against real-valued polynomial sketches $\Phi_{\mathcal{R}}$ with $2D$ random features to obtain the *same* dimension of the feature map.

We measure kernel approximation quality through the relative Frobenius norm error, which is defined as $\|\hat{\mathbf{K}} - \mathbf{K}\|_F / \|\mathbf{K}\|_F$, where $\hat{\mathbf{K}}$ is the random feature approximation of the exact kernel matrix \mathbf{K} with $(\mathbf{K})_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ and $\mathbf{x}_i, \mathbf{x}_j \in X_{*,\text{sub}}$. $X_{*,\text{sub}}$ is a subset of the test data of size $N_{*,\text{sub}} = 1000$ of each dataset that is resampled for 100 different seeds used in these experiments. We use a 90/10 train/test split for Letter and Mocap that is recomputed for every seed, while it is provided for CIFAR-10 and MNIST.

Downstream task. We model GP classification as a multi-class GP regression problem with transformed labels (Milios et al., 2018) for which we obtain closed-form solutions. This allows us to measure the effects of the random feature approximations in isolation without needing to verify the convergence of iterative solvers. To measure the convergence towards a target GP when using random

³The kernel parameters are a standard choice used in, e.g., *scikit-learn* (Pedregosa et al., 2011).

features, we use the Kullback-Leibler (KL) divergence between the GP predictive distribution using the approximate kernel and the predictive distribution using the target kernel. These distributions are evaluated on the test data and the GPs are trained on a random subset of size 5000 of the training data as the reference GP cannot be scaled to a large number of datapoints. Classification errors instead, are obtained for the model trained on the whole training data, and evaluated on all test data.

Results. Fig. 4.8 and Fig. 4.9 show the results of this comparison for the polynomial degrees $p = 3$ and $p = 6$, respectively. We show kernel approximation errors (relative Frobenius norm errors) along with KL-Divergences to the ground-truth GP and test errors for each dataset. As shown in Fig. 4.4, CtR-TensorSRHT has a lower variance than TensorSketch for $p = 3$, and the advantage diminishes for larger p . At the same time, the advantage of CtR-sketches over real ones *increases* with p as shown in Fig. 4.1. These two effects are reflected in the corresponding kernel approximation errors in Fig. 4.8 and Fig. 4.9, respectively.

For $p = 3$, the errors are the lowest for CtR-TensorSRHT and also fluctuate the least. TensorSRHT thus performs better than unstructured Gaussian and Rademacher sketches as discussed in Section 3.5 in Chapter 3. The CtR extension makes TensorSRHT perform similarly on the downstream task of GP classification as TensorSketch (similar KL-divergence and classification error). The results are worse for TensorSRHT without the CtR-extension.

For $p = 6$ (Fig. 4.9), CtR-TensorSRHT and TensorSketch perform similarly well, while the relative advantage of CtR-TensorSRHT over TensorSRHT increases compared to the case of $p = 3$ in Fig. 4.8.

We add a CtR-modification of TensorSketch⁴ to both comparisons in Fig. 4.8 and Fig. 4.9. However, this modification does not yield any improvements in our experiments. CtR-extensions thus only seem to be beneficial for the Product-Sketches discussed in this work.

Lastly, we carry out the same comparison as before for zero-centered data to study the importance of the non-negativity condition of the data. Fig. 4.10 shows the results of this comparison. We can see that the relative advantage of CtR-sketches over their real analogs diminishes compared to the results shown in Fig. 4.9. However, CtR-sketches still achieve an advantage and never perform worse.

We will now conclude the empirical campaign carried out in this chapter with the application of our polynomial sketches to accelerate the computation of an evaluation metric used in *generative modelling*, namely the *Kernel Inception Distance*.

⁴We replace sign hash function $s : [d] \rightarrow \{1, -1\}$ with $s : [d] \rightarrow \{1, -1, i, -i\}$ in (Pham and Pagh, 2013, Def. 1) giving rise to a complex-valued TensorSketch.

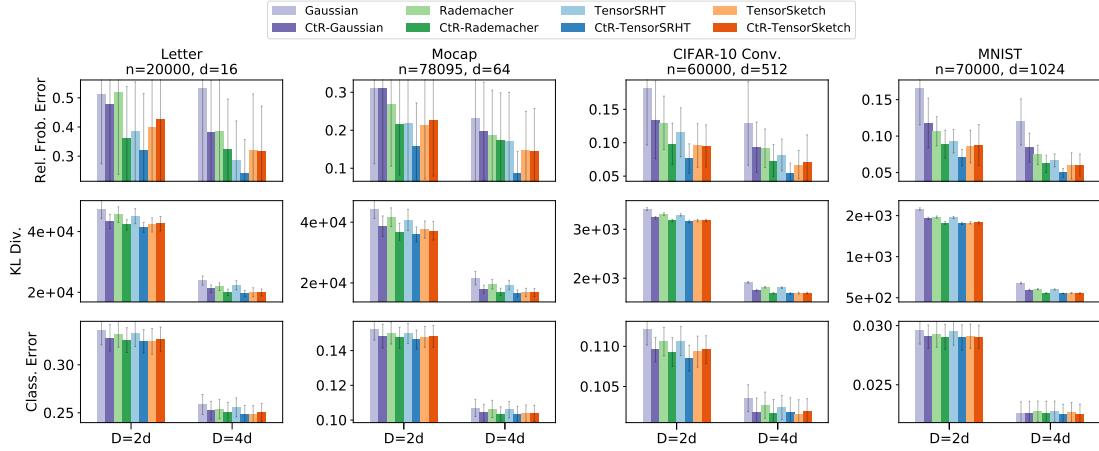


Figure 4.8: Gaussian Process classification experiments comparing different polynomial sketches along with their CtR extensions for $D \in \{2d, 4d\}$. The target polynomial kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^3$ and the data is unit-normalized.

4.3.2 Approximating the Kernel Inception Distance

The *Kernel Inception Distance (KID)* (Bińkowski et al., 2018) is used to measure the similarity between two image data sets (a lower distance refers to a higher similarity). It is especially used in generative modelling to evaluate the quality of a set of generated images by comparing them against a set of ground-truth images.

KID corresponds to the squared maximum mean discrepancy (MMD) (Gretton et al., 2012) using the polynomial kernel $(\frac{1}{d}\mathbf{x}^\top \mathbf{y} + 1)^p$, where this kernel is computed on top of the convolutional features obtained from each image through the Inception (Szegedy et al., 2016) neural network architecture.

As the computation of the MMD metric scales quadratically in the number of data points, approximation methods have to be used for large data sets such as CIFAR-10 that we use here. Therefore, the approximation of the KID is a great application for the polynomial sketches discussed in this work.

Experimental setting. We follow the experiment in the original work that introduced the KID (Bińkowski et al., 2018), which uses the average of blockwise MMD values (Block-MMD) as an approximation. This approximation will serve as a baseline to be compared against polynomial sketches. In this case, the exact squared MMD is computed on a subset of size n of each of the two data sets being compared and this process is repeated for 100 independent subset samples. The final KID estimate is then given by the average of the squared MMDs computed on the subsets. As the KID estimate is unbiased, we compute the standard deviation of these samples to measure the efficiency of the approximator. We use $n \in \{100, 1000, 5000, 10000, 15000, 20000\}$ for the experiments in this section.

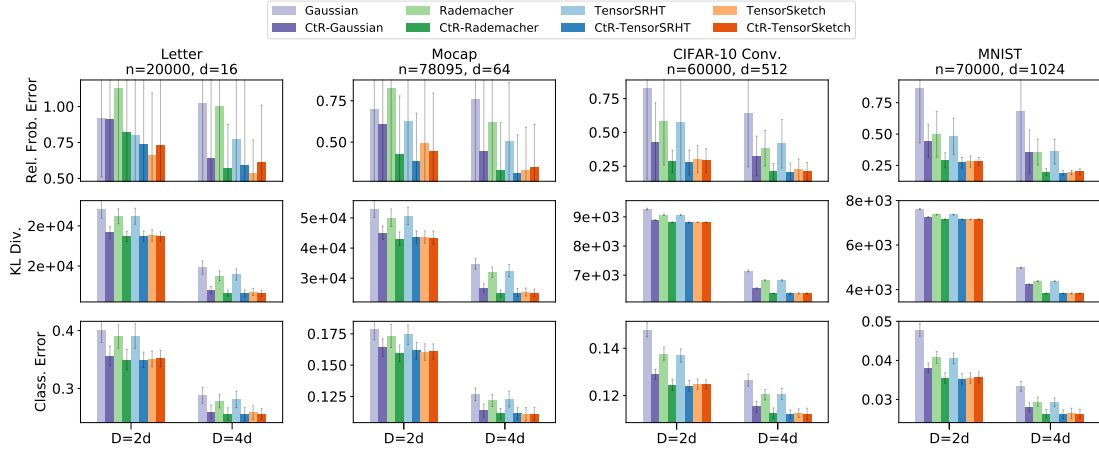


Figure 4.9: Gaussian Process classification experiments comparing different polynomial sketches along with their CtR extensions for $D \in \{2d, 4d\}$. The target polynomial kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^6$ and the data is unit-normalized.

We compare this baseline against squared MMD estimates obtained through polynomial sketches. We obtain unbiased squared MMD estimates using the method described in [Sutherland and Schneider \(2015, Section 3\)](#). We repeat this process 100 times as before and compute the standard deviation on the estimate. We use $D \in \{64, 128, 256, 512, 1024, d, 2d, 3d, 4d\}$ in our experiments, where $d = 2048$ is the dimension of the convolutional features obtained through the Inception neural network. The polynomial sketches being compared are CtR-TensorSRHT, TensorSRHT and TensorSketch.

Results. [Fig. 4.11](#) shows the standard deviations of the squared MMD (KID) estimates computed as explained before. We see that all polynomial sketches achieve around one order of magnitude lower standard deviations than the Block-MMD baseline in the same computation time indicating that they estimate the KID more efficiently. TensorSRHT and CtR-TensorSRHT perform slightly better than TensorSketch for $p = 3$ and CtR-TensorSRHT performs better for $p = 6$.

[Table 4.2](#) shows the same standard deviations as before but computed on 1000 instead of 100 resamples of the KID metric. The results are thus more accurate and confirm that CtR-TensorSRHT provides the best KID estimates for a given feature map dimension D . However, obtaining these accurate estimates is slower and therefore less samples are used in practice.

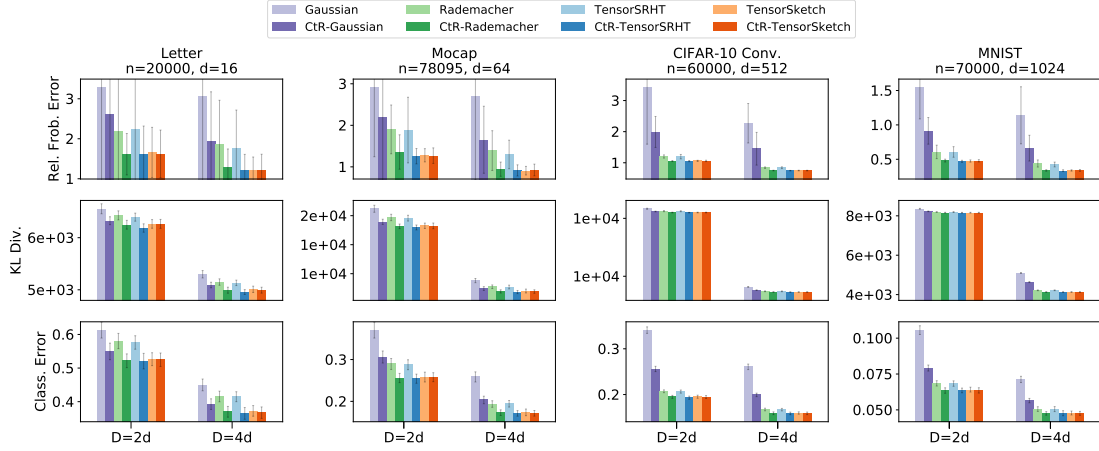
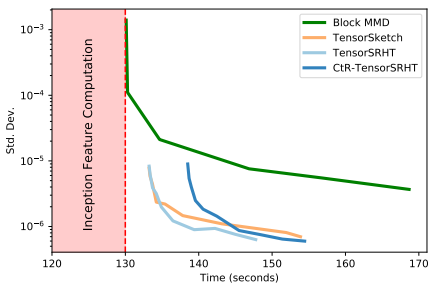


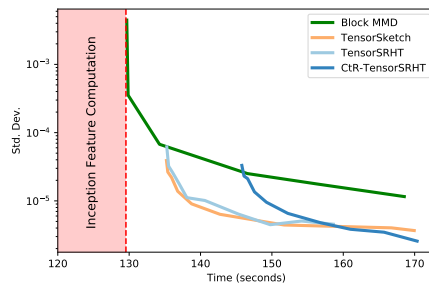
Figure 4.10: Gaussian Process classification experiments comparing different polynomial sketches along with their CtR-extensions for $D \in \{2d, 4d\}$. The target polynomial kernel is $(\mathbf{x}^\top \mathbf{y} + 1)^6$ and the data is unit-normalized as well as zero-centered.

Table 4.2: Standard deviation of the KID estimator between the CIFAR-10 train and test sets estimated on 1000 independent samples.

	$p = 3$				$p = 6$			
	$D = 1d$	$D = 2d$	$D = 3d$	$D = 4d$	$D = 1d$	$D = 2d$	$D = 3d$	$D = 4d$
TensorSketch	1.363E-06	9.639E-07	7.837E-07	6.841E-07	6.974E-06	4.909E-06	4.015E-06	3.436E-06
TensorSRHT	1.435E-06	1.047E-06	7.869E-07	6.946E-07	9.476E-06	6.897E-06	5.331E-06	4.871E-06
CtR-TensorSRHT	1.312E-06	8.747E-07	7.420E-07	5.941E-07	6.697E-06	4.856E-06	3.762E-06	3.260E-06



(a) $p = 3$



(b) $p = 6$

Figure 4.11: Standard deviation of the KID estimator between the CIFAR-10 train and test sets against computation time. Estimates are obtained using 100 independent resamples of the approximate metric. Note that the first 130 seconds are spent on computing the Inception features (the plot starts at 120 to become better readable).

Chapter 5

Approximating Dot Product Kernels

We discuss in this chapter how the polynomial sketches described so far can be used for approximating more general *dot product kernels*, i.e., kernels whose values depend only on the inner product of the input vectors.

There exist alternative random feature approximations for dot product kernels using Gegenbauer polynomials (Han et al., 2022) that have been concurrently developed while we worked on the journal submission associated to the contents of this chapter (Wacker et al., 2022a). We therefore only focus on random feature approximations via polynomial sketches and leave a comparison against random Gegenbauer expansions for future work. We further do not make use of the Complex-to-Real Product-Sketches developed in the previous chapter here. This is because we developed CtR sketches *after* developing the contents of this chapter. Moreover, we will see that the Maclaurin expansions in this chapter are applied to zero-centred data, obstructing the non-negativity condition giving advantages to CtR-sketches in Chapter 4.

In Sections 5.1 and 5.2, we first review a key result on the Maclaurin expansion of dot product kernels and the resulting random sketching approach by Kar and Karnick (2012), and show how the polynomial sketches described so far can be used. In Section 5.3, we then introduce a data-driven optimization approach to improving the random sketches based on the Maclaurin expansion. In Section 5.3.1, we describe how this approach can also be applied to approximate the Gaussian kernel. We then carry out an empirical campaign comparing our optimized Maclaurin method against Spherical Random Features (Pennington et al., 2015) and random Fourier features (Rahimi and Recht, 2007) in Section 5.4. Moreover, we identify a pathology of vanishing approximate kernels for Maclaurin-based approximations of the Gaussian kernel and we propose a cure in Section 5.5.1 with an additional empirical evaluation.

5.1 Maclaurin Expansion of Dot Product Kernels

Let $\mathcal{X} \subseteq \mathbb{R}^d$ be a subset, and let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a positive definite kernel on \mathcal{X} . The kernel k is called *dot product kernel*, if there exists a function $f : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$k(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}^\top \mathbf{y}) \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathcal{X}. \quad (5.1)$$

Examples of dot product kernels include polynomial kernels $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + \nu)^p$ with $\nu \geq 0$ and $p \in \mathbb{N}$, which have been our focus in this thesis, and exponential kernels $k(\mathbf{x}, \mathbf{y}) = \exp(\mathbf{x}^\top \mathbf{y}/l^2)$ with $l > 0$. Other examples of dot product kernels can be found in, e.g., [Smola et al. \(2000\)](#).

We focus on dot product kernels for which the function f in [Eq. \(5.1\)](#) is an analytic function whose Maclaurin expansion has non-negative coefficients: $f(x) = \sum_{n=0}^{\infty} a_n x^n$ and $a_n \geq 0$ for $n \in \{0\} \cup \mathbb{N}$. In other words, we consider dot product kernels that can be expanded as

$$k(\mathbf{x}, \mathbf{y}) = \sum_{n=0}^{\infty} a_n (\mathbf{x}^\top \mathbf{y})^n \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathcal{X}, \quad (5.2)$$

with $a_n \geq 0$ for all $n \in \{0\} \cup \mathbb{N}$.

Many dot product kernels can be expanded as [Eq. \(5.2\)](#). In fact, [Kar and Karnick \(2012, Theorem 1\)](#) show that, if \mathcal{X} is the unit ball of \mathbb{R}^d , the function k of the form of [Eq. \(5.1\)](#) is positive definite on \mathcal{X} if and only if it can be written as [Eq. \(5.2\)](#).

We show here a few concrete examples. The polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + \nu)^p$ with $p \in \mathbb{N}$ and $\nu \geq 0$ can be expanded as

$$(\mathbf{x}^\top \mathbf{y} + \nu)^p = \sum_{n=0}^p \binom{p}{n} \nu^{p-n} (\mathbf{x}^\top \mathbf{y})^n, \quad (5.3)$$

and thus $a_n = \binom{p}{n} \nu^{p-n} \geq 0$ for $n \in \{0, \dots, p\}$ and $a_n = 0$ for $n > p$ in [Eq. \(5.2\)](#). The exponential kernel $k(\mathbf{x}, \mathbf{y}) = \exp(\mathbf{x}^\top \mathbf{y}/l^2)$ can be expanded as

$$\exp\left(\frac{\mathbf{x}^\top \mathbf{y}}{l^2}\right) = \sum_{n=0}^{\infty} \frac{1}{n! l^{2n}} (\mathbf{x}^\top \mathbf{y})^n \quad (5.4)$$

and thus $a_n = 1/(n! l^{2n})$ for $n \in \mathbb{N}$ in [Eq. \(5.2\)](#).

Gaussian kernel as a weighted dot product kernel The Gaussian kernel defined as $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/(2l^2))$ with $l > 0$ can be written as a *weighted*

exponential kernel:

$$\begin{aligned} \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2l^2}\right) &= \exp\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{y}\|^2}{2l^2}\right) \exp\left(\frac{\mathbf{x}^\top \mathbf{y}}{l^2}\right) \\ &= \exp\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{y}\|^2}{2l^2}\right) \sum_{n=0}^{\infty} \frac{1}{n!l^{2n}} (\mathbf{x}^\top \mathbf{y})^n, \end{aligned} \quad (5.5)$$

where the second identity uses the Maclaurin expansion of the exponential kernel in Eq. (5.4). For approximating the Gaussian kernel, [Cotter et al. \(2011\)](#) proposed a finite dimensional feature map based on a truncation of this expansion.

5.2 Random Sketch based on the Maclaurin Expansion

We describe here the approach of [Kar and Karnick \(2012\)](#) on the unbiased approximation of dot product kernels based on the Maclaurin expansion in Eq. (5.2). We discuss this approach to provide a basis and motivation for our new approach for approximating dot product kernels.

First, we define a probability measure μ on $\{0\} \cup \mathbb{N}$. [Kar and Karnick \(2012\)](#) propose to define μ as

$$\mu(n) \propto c^{-(n+1)}, \quad n \in \{0\} \cup \mathbb{N}, \quad (5.6)$$

for a constant $c > 1$ (e.g., $c = 2$).

Using this probability measure and the Rademacher sketch, [Kar and Karnick \(2012\)](#) propose a doubly stochastic approximation of the dot product kernel in Eq. (5.2).

This approach first generates an i.i.d. sample of size $D \in \mathbb{N}$ from this probability measure μ

$$n_1, \dots, n_D \stackrel{i.i.d.}{\sim} \mu \quad (5.7)$$

and defines D_n for $n \in \{0\} \cup \mathbb{N}$ as the number of times n appears in n_1, \dots, n_D ; thus $\sum_{n=0}^{\infty} D_n = D$.

Then, for each $n \in \{0\} \cup \mathbb{N}$ with $D_n > 0$, construct a random feature map $\Phi_n : \mathcal{X} \rightarrow \mathbb{R}^{D_n}$ with D_n features of the form in Eq. (3.1) in Chapter 3 that provide an unbiased approximation of the polynomial kernel $k_n(\mathbf{x}, \mathbf{y}) := (\mathbf{x}^\top \mathbf{y})^n$ of degree n :

$$\mathbb{E}[\Phi_n(\mathbf{x})^\top \Phi_n(\mathbf{y})] = (\mathbf{x}^\top \mathbf{y})^n. \quad (5.8)$$

The original formulation of [Kar and Karnick \(2012\)](#) uses the Rademacher sketch as Φ_n , but one can use other sketches introduced in Chapters 2, 3 and 4, such as the Gaussian sketch and TensorSRHT.

Finally, by defining a random variable $n^* \sim \mu$, the dot product kernel in Eq. (5.2) is rewritten and approximated as

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \sum_{n=0}^{\infty} a_n (\mathbf{x}^\top \mathbf{y})^n = \sum_{n=0}^{\infty} \frac{a_n}{\mu(n)} \mu(n) (\mathbf{x}^\top \mathbf{y})^n = \mathbb{E}_{n^* \sim \mu} \left[\frac{a_{n^*}}{\mu(n^*)} (\mathbf{x}^\top \mathbf{y})^{n^*} \right] \\ &\approx \frac{1}{D} \sum_{n \in \{n_1, \dots, n_D\}} D_n \frac{a_n}{\mu(n)} (\mathbf{x}^\top \mathbf{y})^n = \frac{1}{D} \sum_{n: D_n > 0} D_n \frac{a_n}{\mu(n)} (\mathbf{x}^\top \mathbf{y})^n \\ &\approx \frac{1}{D} \sum_{n: D_n > 0} D_n \frac{a_n}{\mu(n)} \Phi_n(\mathbf{x})^\top \Phi_n(\mathbf{y}), \end{aligned} \quad (5.9)$$

where the first approximation is the Monte Carlo approximation of the expectation $\mathbb{E}_{n^* \sim \mu}$ using the i.i.d. sample in Eq. (5.7) and the second approximation is using the random feature map in Eq. (5.8). The approximation in Eq. (5.9) is unbiased, since the two approximations are statistically independent and both are unbiased.

The first approximation for Eq. (5.9) can be interpreted as first selecting polynomial degrees $n \in \{0\} \cup \mathbb{N}$ and assigning the number of features D_n to each selected degree, given a budget constraint $D = \sum_{n: D_n > 0} D_n$. While performing these assignments by random sampling as in Eq. (5.7) makes the approximation in Eq. (5.9) unbiased, the resulting variance of Eq. (5.9) can be large. In the next section, we introduce a data-driven optimization approach to this feature assignment problem, to achieve a good balance between the bias and variance.

5.3 Optimization for a Truncated Maclaurin Approximation

We develop here an optimization algorithm for selecting the polynomial degrees n and assigning the number of random features to each selected polynomial degree in the Maclaurin sketch in Eq. (5.9). The objective function is an estimate of the expected bias and variance of the resulting approximate kernel, and we define it using the variance formulas derived in Chapter 3.

We consider a biased approximation obtained by truncating the Maclaurin expansion in Eq. (5.2) up to the p -th degree polynomials, where p is to be determined by optimization. Let $D_{\text{total}} \in \mathbb{N}$ be the total number of random features, which is specified by a user. For each $n = 1, \dots, p$, let $D_n \in \{0\} \cup \mathbb{N}$ be the number of random features for approximating the n -th term $(\mathbf{x}^\top \mathbf{y})^n$ of the Maclaurin expansion in Eq. (5.2), such that $\sum_{n=1}^p D_n = D_{\text{total}}$. The numbers D_n are to be determined by optimization. Let $\Phi_n : \mathbb{R}^d \rightarrow \mathbb{C}^{D_n}$ be a (possibly complex) random feature map defined in Chapter 3 such that $\mathbb{E}[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}] = (\mathbf{x}^\top \mathbf{y})^n$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{X} \subset \mathbb{R}^d$. Note that Φ_n can be a real-valued feature map, but we use the notation for the complex case since it subsumes the real case.

We then define an approximation to the dot product kernel in Eq. (5.2) as

$$\hat{k}(\mathbf{x}, \mathbf{y}) := a_0 + \sum_{n=1}^p a_n \Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{X} \quad (5.10)$$

This approximation is biased, since it ignores the polynomial terms whose degrees are higher than p in the expansion of Eq. (5.2). One can reduce this bias by increasing p , but this may lead to a higher variance.

Therefore, there is a bias-variance trade-off in the choice of p . We describe below how to choose p and the number of features D_n of each random feature map $\Phi_n(\mathbf{x}), \Phi_n(\mathbf{y}) \in \mathbb{C}^{D_n}$ for $n = 1, \dots, p$.

Optimization Objective

For a given learning task, we are usually provided data points generated from an unknown probability distribution $P(\mathbf{x})$ on the input domain $\mathcal{X} \subseteq \mathbb{R}^d$. The approximate kernel $\hat{k}(\mathbf{x}, \mathbf{y})$ in Eq. (5.10) should be an accurate approximation of the target kernel $k(\mathbf{x}, \mathbf{y})$ for input vectors \mathbf{x}, \mathbf{y} drawn from this unknown data distribution $P(\mathbf{x})$. Therefore, we consider the following *integrated mean squared error* as our objective function:

$$\begin{aligned} & \int \int \mathbb{E} \left[\left(k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y}) \right)^2 \right] dP(\mathbf{x}) dP(\mathbf{y}) \quad (5.11) \\ &= \int \int \underbrace{\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})]}_{\text{variance}} dP(\mathbf{x}) dP(\mathbf{y}) + \int \int \underbrace{\left(k(\mathbf{x}, \mathbf{y}) - \mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})] \right)^2}_{\text{bias}^2} dP(\mathbf{x}) dP(\mathbf{y}) \quad (5.12) \end{aligned}$$

where the expectation $\mathbb{E}[\cdot]$ and variance $\mathbb{V}[\cdot]$ are taken with respect to the random feature maps in the approximate kernel in Eq. (5.10), and the identity follows from the standard bias-variance decomposition.

We study the variance and bias terms in Eq. (5.12). Let $\delta[D_n > 0]$ be an indicator such that $\delta[D_n > 0] = 1$ if $D_n > 0$ and $\delta[D_n > 0] = 0$ otherwise. Using this indicator, and since the p random feature maps Φ_1, \dots, Φ_p in Eq. (5.10) are statistically independent, the variance term in Eq. (5.12) can be written as

$$\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})] = \sum_{n=1}^p \delta[D_n > 0] a_n^2 \mathbb{V}[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}]. \quad (5.13)$$

Each individual term $\mathbb{V}[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}]$ in Eq. (5.13) is the variance of the approximate kernel $\hat{k}_n(\mathbf{x}, \mathbf{y}) := \Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}$ for approximating the polynomial kernel $k_n(\mathbf{x}, \mathbf{y}) := (\mathbf{x}^\top \mathbf{y})^n$ of degree $n = 1, \dots, p$. Therefore, one can explicitly compute $\mathbb{V}[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}]$ for any given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ using the variance formulas derived Chapter 3. For the convenience of the reader, we summarize the variance formulas for

Sketch	Variance
Real Gaussian	$D^{-1} \left[\left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + 2(\mathbf{x}^\top \mathbf{y})^2 \right)^n - (\mathbf{x}^\top \mathbf{y})^{2n} \right]$
Complex Gaussian	$D^{-1} \left[\left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + (\mathbf{x}^\top \mathbf{y})^2 \right)^n - (\mathbf{x}^\top \mathbf{y})^{2n} \right]$
Real Rademacher	$D^{-1} \left[\left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + 2 \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right) \right)^n - (\mathbf{x}^\top \mathbf{y})^{2n} \right]$
Complex Rademacher	$D^{-1} \left[\left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + (\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right)^n - (\mathbf{x}^\top \mathbf{y})^{2n} \right]$
Real stacked TensorSRHT	Real Rademacher Variance $-\frac{c(D,d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2n} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} \left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + (\mathbf{x}^\top \mathbf{y})^2 - 2 \sum_{k=1}^d x_k^2 y_k^2 \right) \right)^n \right]$
Comp. stacked TensorSRHT	Complex Rademacher Variance $-\frac{c(D,d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2n} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} \left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 - \sum_{k=1}^d x_k^2 y_k^2 \right) \right)^n \right]$
Conv. Sur. TensorSRHT	$\begin{cases} D^{-1} \left(V_q^{(n)} + (d-1) \text{Cov}_q^{(n)} \right) & \text{if } \text{Cov}_q^{(n)} > 0 \text{ or } D > d, \\ D^{-1} \left(V_q^{(n)} - \text{Cov}_q^{(n)} \right) + \text{Cov}_q^{(n)} & \text{otherwise.} \end{cases}$
(Real case: $q = 1$)	$V_q^{(n)} = \left(\ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 + ((2q-1)^2 + 1) \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{k=1}^d x_k^2 y_k^2 \right) \right)^n - (\mathbf{x}^\top \mathbf{y})^{2n}$
(Complex case: $q = 1/2$)	$\text{Cov}_q^{(n)} = \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^n - (\mathbf{x}^\top \mathbf{y})^{2n}$

Table 5.1: Closed-form expressions for the variance $\mathbb{V}[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})}]$ for different random feature maps $\Phi_n : \mathbb{R}^d \rightarrow \mathbb{C}^D$ to approximate polynomial kernel of order $n \in \mathbb{N}$. Here, $D \in \mathbb{N}$ is the number of random features and $c(D, d) := \lfloor D/d \rfloor d(d-1) + (D \bmod d)(D \bmod d - 1)$. See Chapter 3 for details and more generic results. We also show convex surrogate functions in Eq. (D.4) and Eq. (D.5) for the variance of stacked TensorSRHT derived in Appendix D.1.

specific cases in Table 5.1¹. Regarding the bias term in Eq. (5.12), the expectation of the approximate kernel (5.10) is given by

$$\mathbb{E} \left[\hat{k}(\mathbf{x}, \mathbf{y}) \right] = \sum_{n=0}^p \delta[D_n > 0] a_n (\mathbf{x}^\top \mathbf{y})^n, \quad (5.14)$$

since $\mathbb{E} \left[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})} \right] = (\mathbf{x}^\top \mathbf{y})^n$ for $n = 1, \dots, p$ with $D_n > 0$.

Note that the integrals in Eq. (5.12) with respect to P are not available in practice, as P is the unknown data distribution. We instead assume that an i.i.d. sample $\mathbf{x}_1, \dots, \mathbf{x}_m$ of size $m \in \mathbb{N}$ from P is available. This sample may be a subsample of a larger dataset from P . For example, in a supervised learning problem, $\mathbf{x}_1, \dots, \mathbf{x}_m$ may be a random subsample of training input points.

Using the i.i.d. sample $\mathbf{x}_1, \dots, \mathbf{x}_m$, the objective function in Eq. (5.12) can then

¹We omit the variance formulas of CtR-sketches (Chapter 4), since the contents of Chapter 4 were developed after we developed the contents of this chapter. CtR-sketches are thus not included in the empirical evaluations that follow. We also only use stacked TensorSRHT here because we are able to derive a convex surrogate function for it.

be unbiasedly approximated in a U-statistics form as

$$\begin{aligned} & \frac{1}{m(m-1)} \sum_{i \neq j} \mathbb{V}[\hat{k}(\mathbf{x}_i, \mathbf{x}_j)] + \frac{1}{m(m-1)} \sum_{i \neq j} \left(k(\mathbf{x}_i, \mathbf{x}_j) - \mathbb{E}[\hat{k}(\mathbf{x}_i, \mathbf{x}_j)] \right)^2 \\ &= \frac{1}{m(m-1)} \sum_{n=1}^p \delta[D_n > 0] a_n^2 \sum_{i \neq j} \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \Phi_n(\mathbf{x}_j) \right] \end{aligned} \quad (5.15)$$

$$+ \frac{1}{m(m-1)} \sum_{i \neq j} \left(k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{n=0}^p \delta[D_n > 0] a_n (\mathbf{x}_i^\top \mathbf{x}_j)^n \right)^2, \quad (5.16)$$

$$=: g(p, (D_n)_{n=1}^p) \quad (5.17)$$

where we used [Eq. \(5.13\)](#) and [Eq. \(5.14\)](#).

Finally, we formulate our optimization problem. To make the problem tractable, we search for the degree p of the approximate kernel in [Eq. \(5.10\)](#) from the range $\{p_{\min}^*, p_{\min}^* + 1, \dots, p_{\max}^*\}$, where $p_{\min}^*, p_{\max}^* \in \mathbb{N}$ with $p_{\min}^* < p_{\max}^*$ are lower and upper bounds of p selected by the user. We then define our optimization problem as follows:

$$\min_{p, (D_n)_{n=1}^p} g(p, (D_n)_{n=1}^p), \quad (5.18)$$

$$\begin{aligned} \text{subject to: } & p \in \{p_{\min}^*, p_{\min}^* + 1, \dots, p_{\max}^*\}, D_n \in \{0, \dots, D_{\text{total}}\}, \\ & \sum_{n=1}^p D_n = D_{\text{total}}, \text{ and } D_n \geq 1 \text{ if and only if } a_n > 0 \quad (n = 1, \dots, p), \end{aligned}$$

where $g(p, (D_n)_{n=1}^p)$ is defined in [Eq. \(5.17\)](#).

To present our approach to solving [Eq. \(5.18\)](#), we will first define a simplified optimization problem and describe an algorithm for solving it. We will then use this simplified problem and its solver to develop a solver for the full problem in [Eq. \(5.18\)](#).

Solving a Simplified Problem

We consider a simplified problem of [Eq. \(5.18\)](#) in which the polynomial degree $p \in \mathbb{N}$ is fixed and given, and the number of random features D_n is positive, $D_n \geq 1$, for every polynomial degree $n = 1, \dots, p$ with $a_n > 0$. Note that the bias term of the objective function $g(p, (D_n)_{n=1}^p)$, i.e. [Eq. \(5.16\)](#), only depends on $(D_n)_{n=1}^p$ through the indicator function $\delta[D_n > 0]$. Therefore, under the constraint that $D_n \geq 1$ for all $n = 1, \dots, p$ with $a_n > 0$, [Eq. \(5.16\)](#) becomes constant with respect to $(D_n)_{n=1}^p$.

Thus, the optimization problem [Eq. \(5.18\)](#) under the additional constraint of p being fixed and $D_n \geq 1$ for all $n = 1, \dots, p$ with $a_n > 0$ is equivalent to the

following optimization problem:

$$\begin{aligned} \min_{(D_n)_{n=1}^p} \quad & \frac{1}{m(m-1)} \sum_{n=1}^p a_n^2 \sum_{i \neq j} \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right] \\ \text{subject to} \quad & D_n \in \{0, \dots, D_{\text{total}}\}, \quad \sum_{n=1}^p D_n = D_{\text{total}}, \\ & D_n \geq 1 \text{ if and only if } a_n > 0 \quad (n = 1, \dots, p). \end{aligned} \quad (5.19)$$

This is a discrete optimization problem with one equality constraint, and it is an instance of the so-called *Resource Allocation Problem* (Floudas and Pardalos, 2009).

We discuss properties of the objective function in Eq. (5.19) and describe a solver. To this end, we first consider the case where $\Phi_n : \mathbb{R}^d \rightarrow \mathbb{C}^{D_n}$ is one of the unstructured polynomial sketches in Chapter 3; we will later explain its extension to structured sketches from Section 3.5. In this case, we have $\mathbb{V} \left[\Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})} \right] = C_{\mathbf{x}, \mathbf{y}}^{(n)} / D_n$ for a constant $C_{\mathbf{x}, \mathbf{y}}^{(n)}$ depending on $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and the polynomial degree $n \in \mathbb{N}$ but not on D_n , as summarized in Table 5.1. Therefore,

$$a_n^2 \sum_{i \neq j} \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right] = \frac{a_n^2}{D_n} \sum_{i \neq j} C_{\mathbf{x}_i, \mathbf{x}_j}^{(n)} \quad (5.20)$$

is convex and monotonically decreasing with respect to D_n . From this property, one can use the *Incremental Algorithm* (Floudas and Pardalos, 2009, p. 384) to directly solve the optimization problem (5.19).

Algorithm 4 describes the Incremental Algorithm for solving the simplified problem in Eq. (5.19). At every iteration, the algorithm finds $n \in \{1, \dots, p\}$ such that adding one more feature to the feature map Φ_n (i.e., $D_n = D_n + 1$) decreases the objective function most, and sets $D_n = D_n + 1$. Note again that a closed form expression for $\mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right]$ is available from Table 5.1.

Time and space complexities. The time and space complexities of Algorithm 4 are $\mathcal{O}(pD_{\text{total}})$ and $\mathcal{O}(p)$, respectively. Note that from Eq. (5.20), the objective function can be written as

$$f(D_1, \dots, D_p) := \sum_{n=1}^p a_n^2 \sum_{i \neq j} \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right] = \sum_{n=1}^p \frac{a_n^2}{D_n} \sum_{i \neq j} C_{\mathbf{x}_i, \mathbf{x}_j}^{(n)}$$

with a_n and $C_{\mathbf{x}_i, \mathbf{x}_j}^{(n)}$ not depending on the optimizing variable D_n . Therefore, one can precompute the term $\sum_{i \neq j} C_{\mathbf{x}_i, \mathbf{x}_j}^{(n)}$ for each $n = 1, \dots, p$ before starting the iterations in Algorithm 4, and during the iterations one can use the precomputed values of $\sum_{i \neq j} C_{\mathbf{x}_i, \mathbf{x}_j}^{(n)}$. Thus, while the complexity of precomputing $\sum_{i \neq j} C_{\mathbf{x}_i, \mathbf{x}_j}^{(n)}$ is $\mathcal{O}(m^2)$, where m is size of the dataset $\mathbf{x}_1, \dots, \mathbf{x}_m$ defining the objective function (5.19), the time and space complexities of Algorithm 4 do not depend on m .

Algorithm 4: Incremental Algorithm

Result: Optimal solution $D_1, \dots, D_p \geq 1$ to the optimization problem (5.19).

Input: Dot product kernel $k(\mathbf{x}, \mathbf{y}) = \sum_{n=0}^{\infty} a_n (\mathbf{x}^\top \mathbf{y})^n$ with $a_n \geq 0$, truncation order $p \in \mathbb{N}$, the total number of random features $D_{\text{total}} \in \mathbb{N}$; Initialize $D_1 = \dots = D_p = 1$ and $t = 0$;

Let $f(D_1, \dots, D_p) := \sum_{n=1}^p a_n^2 \sum_{i \neq j} \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right]$.

while $t < D_{\text{total}}$ **do**

Find $j^* = \arg \min_{j \in \{1, \dots, p\}} f(D_1, \dots, D_j + 1, \dots, D_p)$;

$D_{j^*} = D_{j^*} + 1$;

$t = t + 1$;

end

Structured case. We assumed here that Φ_n is one of the unstructured sketches studied in Chapter 3. This choice of Φ_n makes Eq. (5.20) convex and monotonically decreasing with respect to D_n , which enables the Incremental Algorithm to solve the optimization problem in Eq. (5.19).

However, if Φ_n is a structured sketch (i.e., either real or complex TensorSRHT) in Section 3.5, Eq. (5.20) is not convex with respect to D_n , and the Incremental Algorithm is not directly applicable. To overcome this problem, we propose to use *stacked TensorSRHT* along with the convex surrogate functions in Eq. (D.4) and Eq. (D.5) derived in Appendix D.1 to replace $\mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right]$ in the objective function (5.19), and then apply the Incremental Algorithm. We summarize the concrete form of the convex surrogate function in Table 5.1. For details, see Appendix D.1.

Solving the Full Problem

We now address the full problem in Eq. (5.18) using Algorithm 4 developed for the simplified problem in Eq. (5.19). Recall that, by fixing $p \in \{p_{\min}^*, \dots, p_{\max}^*\}$ and constraining $D_n \geq 1$ for all $n = 1, \dots, p$, the full problem in Eq. (5.18) becomes equivalent to the simplified problem in Eq. (5.19), which can be solved by Algorithm 4. Thus, we propose to solve the full problem in Eq. (5.18) by i) first performing Algorithm 4 for each $p \in \{p_{\min}, \dots, p_{\max}\}$, ii) then evaluate each solution $(D_n)_{n=1}^p$ by computing the objective function $g(p, (D_n)_{n=1}^p)$ in Eq. (5.17), and finally pick up p that gives the smallest objective function value.

Algorithm 5 summarizes the whole procedure for solving the full optimization problem in Eq. (5.18). Algorithm 5 returns the optimal truncation order $p^* \in \{p_{\min}, \dots, p_{\max}\}$ with the corresponding feature cardinalities $D^* = (D_1, \dots, D_{p^*})$. Given these values, one can construct a feature map as summarized in Algorithm 6. Note that the U-statistics in the empirical objective (5.17) can be precomputed for

Algorithm 5: Extended Incremental Algorithm

Result: Optimal polynomial degree $p^* \in \{p_{\min}, \dots, p_{\max}\}$ and feature cardinalities $D^* = (D_1, \dots, D_{p^*}) \in \mathbb{N}^{p^*}$ to the full optimization problem (5.18).

Input: Dot product kernel $k(\mathbf{x}, \mathbf{y}) = \sum_{n=0}^{\infty} a_n (\mathbf{x}^\top \mathbf{y})^n$ with $a_n \geq 0$, upper and lower bounds $p_{\min}, p_{\max} \in \mathbb{N}$, the total number of random features $D_{\text{total}} \in \mathbb{N}$;

Set $g^* = \infty$, $p^* = p_{\min}$ and $D^* = \{\}$;

forall $p \in \{p_{\min}, \dots, p_{\max}\}$ **do**

- | Solve Algorithm 4 to obtain D_1, \dots, D_p ;
- | Compute $g(p, (D_n)_{n=1}^p)$ in Eq. (5.17) ;
- | **if** $g(p, (D_n)_{n=1}^p) < g^*$ **then**
- | | $g^* = g(p, (D_n)_{n=1}^p)$;
- | | $D^* = (D_n)_{n=1}^p$;
- | | $p^* = p$;

end

Algorithm 6: Improved Random Maclaurin (RM) Features

Result: Feature map $\Phi(\mathbf{x}) \in \mathbb{C}^{D_{\text{total}}+1}$

Input: Dot product kernel $k(\mathbf{x}, \mathbf{y}) = \sum_{n=0}^{\infty} a_n (\mathbf{x}^\top \mathbf{y})^n$ with $a_n \geq 0$, polynomial degree $p^* \in \mathbb{N}$ and feature cardinalities D_1, \dots, D_{p^*} from Algorithm 5 ;

Initialize $\Phi(\mathbf{x}) := [\sqrt{a_0}]$

forall $n \in \{1, \dots, p^*\}$ **do**

- | Let $\Phi_n(\mathbf{x}) \in \mathbb{C}^{D_n}$ be an unbiased polynomial sketch of degree n with D_n features (see Chapter 3) ;
- | Append $\sqrt{a_n} \Phi_n(\mathbf{x})$ to $\Phi(\mathbf{x})$;

end

all $p_{\min}, \dots, p_{\max}$ before running any optimization algorithm. They do not need to be re-evaluated for every execution of Algorithm 4.

5.3.1 Approximating a Gaussian Kernel

Here we describe how to adapt Algorithm 4 and Algorithm 5 for approximating a Gaussian kernel of the form $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/(2l^2))$ with $l > 0$. By Eq. (5.5), this Gaussian kernel can be written as

$$k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{y}\|^2}{2l^2}\right) \sum_{n=0}^{\infty} a_n (\mathbf{x}^\top \mathbf{y})^n,$$

where $a_n := 1/(n!l^{2n})$ for $n \in \mathbb{N} \cup \{0\}$. Notice that $\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right)$ and $\left(-\frac{\|\mathbf{y}\|^2}{2l^2}\right)$ are scalar values and can be computed for any given input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$.

Thus, the objective function $g(p, (D_n)_{n=1}^p)$ in Eq. (5.17), which is an empirical approximation of the bias-variance decomposition of the mean squared error in Eq. (5.12) using an i.i.d. sample $\mathbf{x}_1, \dots, \mathbf{x}_m \stackrel{i.i.d.}{\sim} P$, can be adapted as

$$\begin{aligned} & g(p, (D_n)_{n=1}^p) \tag{5.21} \\ &= \frac{1}{m(m-1)} \sum_{n=1}^p \delta[D_n > 0] a_n^2 \sum_{i \neq j} \exp\left(-\frac{\|\mathbf{x}_i\|^2}{l^2}\right) \exp\left(-\frac{\|\mathbf{x}_j\|^2}{l^2}\right) \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right] \\ &+ \frac{1}{m(m-1)} \sum_{i \neq j} \left(k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{n=0}^p \delta[D_n > 0] a_n \exp\left(-\frac{\|\mathbf{x}_i\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{x}_j\|^2}{2l^2}\right) (\mathbf{x}_i^\top \mathbf{x}_j)^n \right)^2. \end{aligned}$$

Accordingly, the objective function of the simplified problem in Eq. (5.19) is adapted as

$$f(D_1, \dots, D_p) := \frac{1}{m(m-1)} \sum_{n=1}^p a_n^2 \sum_{i \neq j} \exp\left(-\frac{\|\mathbf{x}_i\|^2}{l^2}\right) \exp\left(-\frac{\|\mathbf{x}_j\|^2}{l^2}\right) \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right].$$

By these modifications, Algorithm 4 and Algorithm 5 can be used to obtain the optimal truncation order $p^* \in \{p_{\min}, \dots, p_{\max}\}$ and the corresponding feature cardinalities D_1, \dots, D_{p^*} . Lastly, Algorithm 6 can be adapted by multiplying the scalar value $\exp\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right)$ to the feature map $\Phi(\mathbf{x})$ obtained from Algorithm 6: the new feature map is defined as $\Phi'(\mathbf{x}) := \exp\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right) \Phi(\mathbf{x})$.

5.3.2 Numerical Illustration of the Objective Function

To gain an insight about the behavior of Algorithm 5, we provide a numerical illustration of the bias and variance terms in the objective function $g(p, (D_n)_{n=1}^p)$ in Eq. (5.17) (or its version adapted for the Gaussian kernel in Eq. (5.21)). To this end, we used the Fashion MNIST dataset (Xiao et al., 2017) and randomly sampled data points $\mathbf{x}_1, \dots, \mathbf{x}_m$ with $m = 500$ from the entire dataset of size 60,000. As a target kernel to approximate, we consider (i) a polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} / 8 + 7/8)^{20}$ of degree $p = 20$; and (ii) the Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2l^2))$, where the length scale $l > 0$ is given by the median heuristic (Garreau et al., 2017), i.e., the median of the pairwise Euclidean distances of $\mathbf{x}_1, \dots, \mathbf{x}_m$.

For the polynomial kernel (i), we computed (a) $\frac{a_n^2}{m(m-1)} \sum_{i \neq j} \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right]$ for each $n = 1, \dots, p$ ($= 20$), which is the variance component of the objective function in Eq. (5.17); and (b) $\frac{1}{m(m-1)} \sum_{i \neq j} \left(k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{\nu=0}^n a_\nu (\mathbf{x}_i^\top \mathbf{x}_j)^\nu \right)^2$ for each $n = 1, \dots, p$ ($= 20$), which is the bias component in Eq. (5.17) computed up to

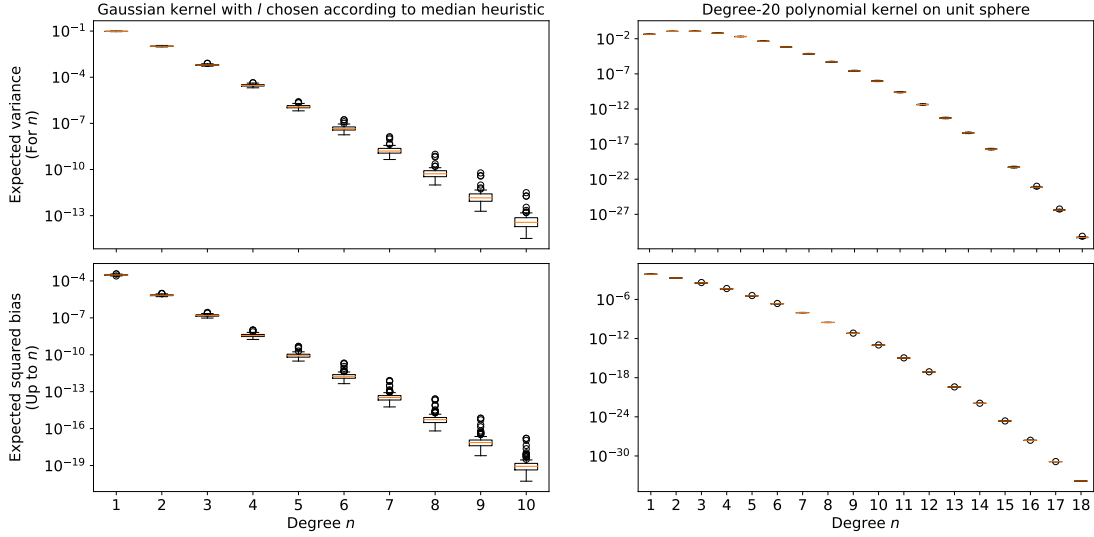


Figure 5.1: Numerical illustration of Section 5.3.2. The left two figures are box plots for the Gaussian kernel (i), and the right two figures are those for the polynomial kernel (ii). The top figures show the variance terms (a), and the bottom figures show the bias terms (b). See Section 5.3.2 for details.

n -th order. For the Gaussian kernel (ii), we computed corresponding quantities from the objective function in Eq. (5.21):

$$(a) \frac{a_n^2}{m(m-1)} \sum_{i \neq j} \exp\left(-\frac{\|\mathbf{x}_i\|^2}{l^2}\right) \exp\left(-\frac{\|\mathbf{x}_j\|^2}{l^2}\right) \mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right] \text{ and}$$

$$(b) \frac{1}{m(m-1)} \sum_{i \neq j} \left(k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{\nu=0}^n a_\nu \exp\left(-\frac{\|\mathbf{x}_i\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{x}_j\|^2}{2l^2}\right) (\mathbf{x}_i^\top \mathbf{x}_j)^\nu \right)^2$$

for $n = 1, \dots, 10$.

We used the real Gaussian sketch for the feature map Φ_n , for which Eq. (3.8) in Chapter 3 gives a closed form expression of the variance $\mathbb{V} \left[\Phi_n(\mathbf{x}_i)^\top \overline{\Phi_n(\mathbf{x}_j)} \right]$; see also Table 5.1. We set $D_n = 1$ for each n to be evaluated (i.e., $\Phi_n(\mathbf{x}) \in \mathbb{R}$).

To compute the means and standard deviations of the above quantities (a) and (b), we repeated this experiment 100 times by independently subsampling $\mathbf{x}_1, \dots, \mathbf{x}_m$ with $m = 500$ from the entire dataset each time. Fig. 5.1 describes the results. First, we can see that the standard deviations of the quantities (a) and (b) are relatively small, and thus a subsample $\mathbf{x}_1, \dots, \mathbf{x}_m$ of size $m = 500$ is sufficient for providing accurate approximations of the respective population quantities of (a) and (b) (where the empirical average with respect to $\mathbf{x}_1, \dots, \mathbf{x}_m$ is replaced by the corresponding expectation) in this setting.

Regarding the polynomial kernel (i), the variance terms (a) for polynomial degrees up to $n = 3$ have similar magnitudes, and they decay exponentially fast for polynomial degrees larger than $n = 3$ (notice that the vertical axis of the plot is in log scale). On the other hand, the bias term (b) decays exponentially fast as the polynomial degree n increases. These trends suggest that Algorithm 5

would assign more features to lower order degrees n , in particular to the degree 3 or less. One explanation of these trends is that the parametrization of the kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y}/8 + 7/8)^{20}$ gives larger coefficients to lower polynomial degrees in the the Maclaurin expansion (see Eq. (5.3)), and that the distribution of pairwise inner products of the data points $\mathbf{x}_1, \dots, \mathbf{x}_m$ is centered around zero in this experiment.

Regarding the Gaussian kernel (ii), both the variance term (a) and the bias term (b) decay exponentially fast as the polynomial degree n increases. This trend suggests that Algorithm 5 would assign more features to lower order polynomial degrees n .

To summarize, these observations suggest that, to minimize the mean squared error of the approximate kernel, it is more advantageous to assign more features to lower degree polynomial approximations. Algorithm 5 automatically achieves such feature assignments.

We will show later in this chapter for the Gaussian kernel that Algorithm 5 assigns more random features to *higher* polynomial degrees instead when the length scale $l > 0$ is small in comparison to the scaling of the data, i.e., when modelling high-frequency data and choosing l *without* the median heuristic. This is because the variances of approximate polynomial kernels of high degrees increase in this case. In the following experiments we stick to the median heuristic for now.

5.3.3 Gaussian Process Regression Toy Example

We performed a toy experiment on one-dimensional Gaussian process (GP) regression, whose results are described in Fig. 5.2. The purpose is to gain a qualitative understanding of the optimized Maclaurin approximation in Section 5.3 (Algorithm 5). For comparison, we also used Random Fourier Features (RFF) of Rahimi and Recht (2007) described in Chapter 2 in this experiment. We use the real Rademacher sketch in the optimized Maclaurin approach.

We define the ground-truth function as a sinc function, $f(x) = \sin(ax)/x$, with $a > 0$ for which we consider three settings: $a \in \{5, 2, 0.5\}$. We generated training data by adding independent Gaussian noise of variance $\sigma_{\text{noise}}^2 = 0.01$ to the ground-truth function $f(x)$. With this value of noise variance σ_{noise}^2 , we then fit a GP regressor (see Chapter 1) using the Gaussian kernel $k(x, y) = \sigma^2 \exp(-(x - y)^2/(2l^2))$ to the training data, where we determined the hyperparameters $l, \sigma^2 > 0$ by maximizing the log marginal likelihood (e.g., Rasmussen and Williams, 2006, Chapter 2). We used the resulting posterior GP as a ground-truth and call it “full GP”, treating it as a reference for assessing the quality of approximate GPs. As such, we used the same hyperparameters as the full GP in approximate GPs; this enables evaluating the effects of the approximation in the resulting GP predictive distributions.

We set the number of random features as $D = 10$. In this case, the optimized Maclaurin approach in Algorithm 5 selects the truncation degree $p^* = 9$ and simply

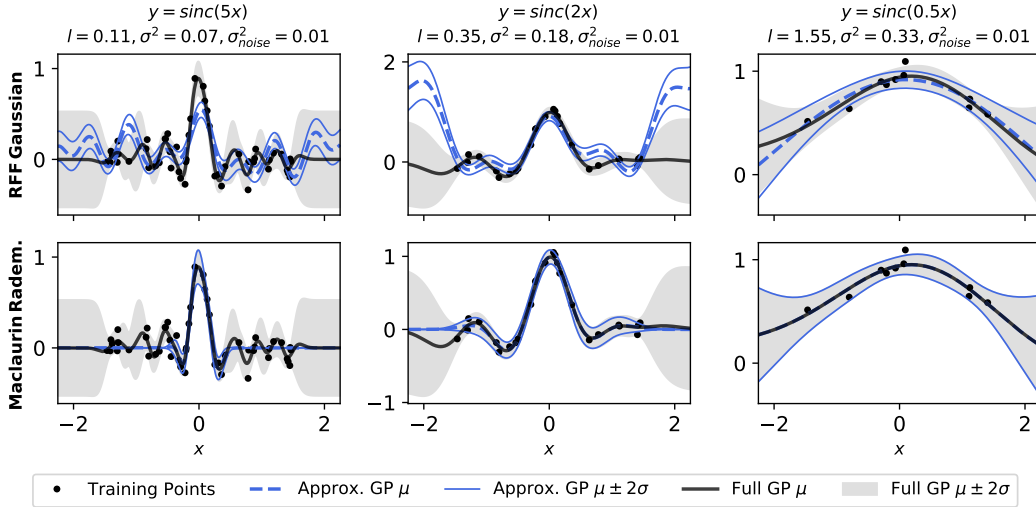


Figure 5.2: One-dimensional GP regression experiment in Section 5.3.3. The top row are the results of random Fourier features (Gaussian RFF) and the bottom row are those of the optimized Maclaurin approach. The left, middle, and right columns correspond to the ground-truth sinc functions with frequency 5, 2, and 0.5, respectively. The values of l and σ^2 are the kernel hyperparameters obtained by maximizing the log likelihood of training data in the full GP (i.e., without approximation). Dashed blue curves represent approximate GP posterior mean functions; blue curves represent the posterior means plus and minus 2 times approximate posterior standard deviations; black curves represent the posterior mean functions of the full GP; and the shaded areas are the full GP posterior means plus and minus 2 times the full GP posterior deviations.

allocates the feature cardinalities as $D_1 = \dots = D_9 = 1$. (Note that one feature is always allocated to the degree $n = 0$). This behavior is expected because the variance of the Rademacher sketch (see Table 5.1) is zero for all polynomial degrees n , as the input dimension is one ($d = 1$) in this experiment.²

We can make the following observations from Fig. 5.2. First, with the optimized Maclaurin approach, the approximate GP posterior mean function approximates the full GP posterior mean function around $x = 0$ more accurately than RFF. Moreover, the range of x on which the Maclaurin approach is accurate becomes wider for a lower frequency a of the ground-truth sinc function (for which the length scale l is larger). This tendency suggests that the Maclaurin approach may be more advantageous than RFF in approximating around $x = 0$ and when the length scale l is relatively large. Experiments in the next section, in particular those with high dimensional datasets, provide a further support for this observation.

One issue with the Maclaurin approximation is that, as can be seen from

²Thus, the error of the optimized Maclaurin approach stems solely from the finite truncation of the Maclaurin expansion in Eq. (5.10).

Fig. 5.2, the approximate GP posterior tends to collapse for an input location x far from 0. We will study this issue at the end of this chapter, and show that the collapsing GP posterior approximation is due to the vanishing of the approximate kernel $\hat{k}(\mathbf{x}, \mathbf{y})$ when $\|\mathbf{x}\|/l + \|\mathbf{y}\|/l$ becomes large. We will also propose a cure for this pathology that is relevant in particular when modelling data of high frequency, i.e., when l is short relative to the scaling of the data.

For now, we focus on a set of high-dimensional experiments, where we can make use of the median heuristic and approximations are well-behaved.

5.4 Empirical Evaluation of the Optimized Maclaurin Method

In this section, we perform systematic experiments to evaluate the optimized Maclaurin method proposed in this chapter. We use it in conjunction with the real and complex polynomial sketches introduced in Chapter 3, and we consider approximations of both polynomial kernels and Gaussian kernels. We do not make use of the Complex-to-Real sketches presented in Chapter 4 here, since we developed them independently from the methods discussed in this chapter, as mentioned at the beginning of this chapter.

We evaluate the performance of each approximation approach in terms of both i) the accuracy in kernel approximation and ii) the performance in downstream tasks. The downstream tasks we consider are Gaussian process regression and classification. For completeness, we explain how to use complex-valued random features in Gaussian process inference and discuss the resulting computational costs in Appendix B.3.

In Section 5.4.1, we first explain the setup of the experiments. In Section 5.4.2, we describe experiments on polynomial kernel approximation, comparing various approximation approaches. In Section 3.5.4, we report the results of wall-clock time comparison of real and complex random features, focusing on the downstream task performance of GP classification. In Section 5.4.3, we present detailed evaluations of the optimized Maclaurin approach for polynomial and Gaussian kernel approximations in GP classification and regression.

5.4.1 Experimental Setup

We explain here the common setup for the experiments in this section.

Datasets

Table 5.2 shows an overview of the datasets used in the experiments. All the datasets come from the UCI benchmark (Dua and Graff, 2017) except for Cod_rna (Uzilov et al., 2006), FashionMNIST (Xiao et al., 2017), and MNIST (Lecun et al.,

Classification	Num. data points N	Dimensionality d	Regression	Num. data points N	Dimensionality d
Cod_rna	331,152	8	Boston	506	16
Coverttype	581,012	64	Concrete	1,030	8
FashionMNIST	70,000	1,024	Energy	768	8
Magic	19,020	16	kin8nm	8,192	8
Miniboo	130,064	64	Naval	11,934	16
MNIST	70,000	1,024	Powerplant	9,568	4
Mocap	78,095	64	Protein	45,730	16

Table 5.2: Datasets used in the experiments. The left and right columns are datasets for classification and regression, respectively.

1998). We pad input vectors with zeros so that the input dimensionality d becomes a power of two to support Hadamard projections in TensorSRHT. The train/test split is 90/10 and is recomputed for every random seed for the UCI datasets; otherwise it is predefined.

For each dataset, we use its random subsets of size $m = \min(5000, N_{\text{train}})$ and $m_* = \min(5000, N_{\text{test}})$ to define training and test data in an experiment, respectively, where N_{train} and N_{test} are the sizes of the original training and test datasets. Denote by $X_{\text{sub}} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ and $X_{*,\text{sub}} = \{\mathbf{x}_{*,1}, \dots, \mathbf{x}_{*,m_*}\}$ those subsets for training and test, respectively. We repeat each experiment 10 times independently using 10 different random seeds, and hence with 10 different subset partitions.

Target Kernels to Approximate

We consider the approximation of (i) polynomial kernels and (ii) Gaussian kernels.

(i) Polynomial kernel approximation. We consider a polynomial kernel of the form

$$k(\mathbf{x}, \mathbf{y}) = \sigma^2 \left(\left(1 - \frac{2}{a^2} \right) + \frac{2}{a^2} \mathbf{x}^\top \mathbf{y} \right)^p = \sigma^2 \left(1 - \frac{\|\mathbf{x} - \mathbf{y}\|^2}{a^2} \right)^p \quad (5.22)$$

with $p \in \mathbb{N}$, $a \geq 2$, $\sigma^2 > 0$, and $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$. We choose this form of polynomial kernels because we use *Spherical Random Features (SRF)* of Pennington et al. (2015) as one of our baselines, and because SRF approximates the polynomial kernel in Eq. (5.22) defined on the unit sphere of \mathbb{R}^d . Following the experimental setup of Pennington et al. (2015), we set $a = 4$ and $p \in \{3, 7, 10, 20\}$ in Eq. (5.22). We set σ^2 as the variance of the labels of training subset X_{sub} .

To make SRF applicable, we unit-normalize the input vectors in each dataset so that they lie on the unit sphere in \mathbb{R}^d . In an experiment where we zero-centralize the input vectors, we unit-normalize after applying the zero-centering.

(ii) Gaussian kernel approximation. We consider the approximation of the Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2l^2))$, where we choose the length

scale $l > 0$ by the median heuristic (Garreau et al., 2017), i.e., as the median of pairwise Euclidean distances of input vectors in the training subset X_{sub} ; we set $\sigma^2 > 0$ as the variance of the labels of X_{sub} .

Error Metrics

We define several error metrics for studying the quality of each approximation approach.

Relative Frobenius norm error. To define this error metric, we need to define some notation. Let $\Phi : \mathbb{R}^d \rightarrow \mathbb{C}^D$ be the (either real or complex) feature map of a given approximation method. For test input vectors $\mathbf{X}_{*,\text{sub}} = \{\mathbf{x}_{*,1}, \dots, \mathbf{x}_{*,m_*}\}$, let $\hat{\mathbf{K}} \in \mathbb{C}^{m_* \times m_*}$ be the approximate kernel matrix such that $\hat{\mathbf{K}}_{i,j} = \Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$. Similarly, let $\mathbf{K} \in \mathbb{R}^{m_* \times m_*}$ be the exact kernel matrix such that $\mathbf{K}_{i,j} = k(\mathbf{x}_{*,i}, \mathbf{x}_{*,j})$ with k being the target kernel.

We then define the *relative Frobenius norm error* of $\hat{\mathbf{K}}$ against \mathbf{K} as:

$$\|\mathbf{K} - \hat{\mathbf{K}}\|_F / \|\mathbf{K}\|_F := \sqrt{\sum_{i=1}^m \sum_{j=1}^m |\mathbf{K}_{i,j} - \hat{\mathbf{K}}_{i,j}|^2} / \sqrt{\sum_{i=1}^m \sum_{j=1}^m \mathbf{K}_{i,j}^2}. \quad (5.23)$$

This error quantifies the quality of the feature map Φ in terms of the resulting approximation accuracy of the kernel matrix. As the target kernel matrix \mathbf{K} is real-valued, we discard the imaginary part of $\hat{\mathbf{K}}$ if it is complex-valued, unless otherwise specified.

We define other error metrics in terms of two downstream tasks: Gaussian process (GP) regression and classification (see Appendix B.3 for details of these GP tasks).

Kullback-Leibler (KL) divergence. We measure the *KL divergence* between two posterior predictive distributions at test input points: one is that of an approximate GP and the other is that of the exact GP without approximation; see Eq. (B.23) in Appendix B.3 for details. For GP classification, we measure the KL divergence between the corresponding latent GPs before transformation. Since there are as many GPs as the number of classes, we report the KL divergence averaged over those classes.

Prediction performance. For GP classification, we use the *test error rate* (i.e., the percentage of misclassified examples) for measuring the prediction performance. For GP regression, we report the *normalized mean squared error (norm. MSE)* between the posterior predictive outputs and true outputs, normalized by the variance of the test outputs. Here, we use the full training data of size N_{train} for

computing the approximate GP posterior and the full test data of size N_{test} for evaluating the prediction performance.³

Mean negative log likelihood (MNLL). We compute the *mean negative log likelihood (MNLL)* of the test data for the approximate GP predictive distribution. MNLL can capture the quality of prediction uncertainties of the approximate GP model (e.g. [Rasmussen and Williams, 2006](#), p. 23). We use the full training and test data for computing the MNLL, as for the prediction performance.

Other Settings

Optimized Maclaurin approach. For the optimized Maclaurin approach in [Algorithm 5](#), we set $p_{\min} = 2$ and $p_{\max} = 10$. We use the training subset $X_{\text{sub}} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ to precompute the U-statistics in [Eq. \(5.15\)](#) and [Eq. \(5.16\)](#).

Regularization parameters. We select the regularization parameter in GP classification and regression by a training-validation procedure. That is, we use the 90 % of training data for training and the remaining 10 % for validation, and select the regularization parameter that maximizes the MNLL on the validation set. For GP classification, we choose the regularization parameter from the range $\alpha \in \{10^{-5}, \dots, 10^{-0}\}$. For GP regression, we choose the noise variance from the range $\sigma_{\text{noise}}^2 \in \{2^{-15}, \dots, 2^{15}\}$. See [Appendix B.3](#) for the definition of these parameters.

Importantly, we perform this selection procedure using a baseline approach,⁴ and after selecting the regularization parameter, we set the *same* regularization parameter for all the approaches (including our optimized Maclaurin approach) for computing error metrics. In this way, we make sure that the selected regularization parameter is not in favour of our approaches (and in this sense we give an advantage to the baseline).

5.4.2 Polynomial Kernel Approximation

We first study the approximation of the polynomial kernels in [Eq. \(5.22\)](#), comparing different polynomial sketches and the optimized Maclaurin method in terms of the relative Frobenius norm error in [Eq. \(5.23\)](#) on FashionMNIST. [Fig. 5.3](#) describes the results. We consider the following polynomial sketches in this experiment:

³We did not use the full training and test datasets for evaluating the KL divergence, since it requires computing the exact GP posterior on the full training data of size N_{train} , which costs $\mathcal{O}(N_{\text{train}}^3)$ and is not feasible for datasets with large N_{train} .

⁴More specifically, we use the Spherical Random Features (SRF) ([Pennington et al., 2015](#)) when the target kernel is a polynomial kernel, and Random Fourier Features ([Rahimi and Recht, 2007](#)) when the target kernel is Gaussian, for selecting the regularization parameter.

(i) Rademacher Product-Sketch (Chapter 3). We use the real Rademacher sketch, i.e., the unstructured polynomial sketch in Eq. (3.1) with Rademacher weights (“Radem.” in Fig. 5.3).

(ii) Stacked TensorSRHT (Section 3.5). We consider the real stacked TensorSRHT in Eq. (3.27) with Rademacher weights (“TensorSRHT” in Fig. 5.3), and the complex TensorSRHT in Eq. (3.34) with complex Rademacher weights; see also Algorithm 2. We consider two versions of the complex TensorSRHT: one that keeps the imaginary part in the approximate kernel matrix (“TensorSRHT Comp. (Keep imag.)” in Fig. 5.3), and one that discards the imaginary part (“TensorSRHT Comp. (Disc. imag.)” in Fig. 5.3).

(iii) Random Maclaurin (Section 5.2). We use the Random Maclaurin approach by Kar and Karnick (2012) explained in Section 5.2. To improve its performance, we truncate the support of the importance sampling measure $\mu(n) = 2^{-(n+1)}$ in Eq. (5.6) to degrees $n \in \{1, \dots, p\}$.⁵ Note that the term $n = 0$ in Eq. (5.2) associated with coefficient a_0 does not need to be approximated, as we append $\sqrt{a_0}$ to the feature map. We consider two versions of the Random Maclaurin approach: one using the real Rademacher sketch (“Rand. Macl. Radem.” in Fig. 5.3) and one using the real TensorSRHT (“Rand. Macl. TensorSRHT” in Fig. 5.3).

(iv) Optimized Maclaurin (Section 5.3). We consider our optimized Maclaurin approach introduced in this chapter using the real TensorSRHT (“Opt. Macl. TensorSRHT” in Fig. 5.3) and using the complex TensorSRHT. We consider two versions of the latter: one keeping the imaginary part of the approximate kernel matrix (“Opt. Macl. TensorSRHT Comp. (Keep imag.)” in Fig. 5.3) and one discarding the imaginary part (“Opt. Macl. TensorSRHT Comp. (Disc. imag.)” in Fig. 5.3).

(v) TensorSketch For completeness, we also include in this experiment *TensorSketch* of Pham and Pagh (2013), a state-of-the-art polynomial sketch (“TensorSketch” in Fig. 5.3).

Setting. We perform the experiments using FashionMNIST (“Non-centered data” in Fig. 5.3) and its centered version for which we subtract the mean of the input vectors from each input vector (“Centered data” in Fig. 5.3). For each approach, the number of random features is $D \in \{d, 2d\}$, where $d = 1,024$ for FashionMNIST.

⁵Without this restriction of the support, the randomized Maclaurin approach may sample polynomial degrees n such that $n > p$ from $\mu(n)$, for which the associated coefficient in the Maclaurin expansion in Eq. (5.3) is zero. Therefore, the resulting feature maps may contain zeros, which are redundant and makes the kernel approximation inefficient.

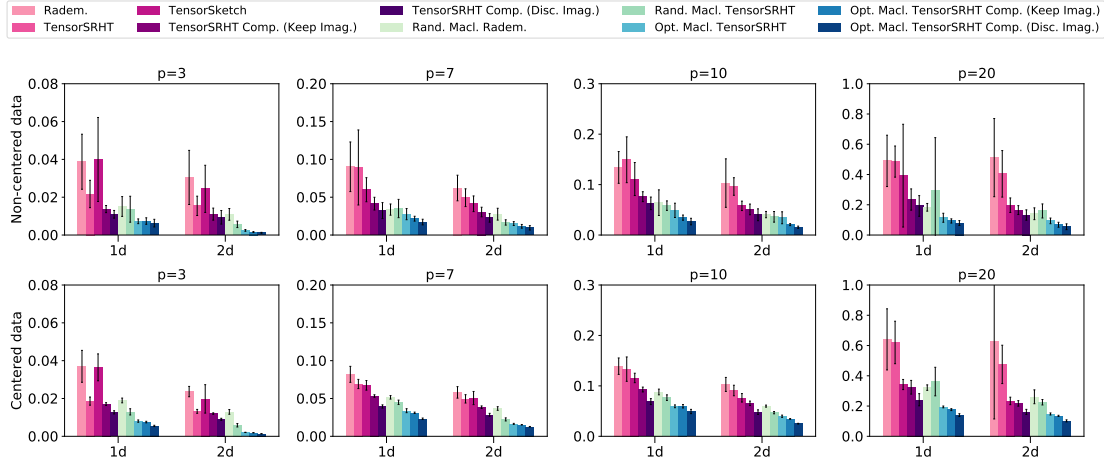


Figure 5.3: Results of the experiments in Section 5.4.2 using FashionMNIST. Each plot shows the relative Frobenius norm errors in Eq. (5.23) of different sketches for approximating the polynomial kernel in Eq. (5.22) with $p \in \{3, 7, 10, 20\}$ and $D \in \{1d, 2d\}$. The top and bottom rows show results with the data without and with zero-centring, respectively.

From the results in Fig. 5.3, we can make the following observations.

Effectiveness of the optimization approach. The optimized Maclaurin approach with the real TensorSRHT (“Opt. Mac. TensorSRHT”) achieves smaller errors than the corresponding random Maclaurin approach (“Rand. Macl. TensorSRHT”) for almost all cases, and with a large margin for $p = 3$ and $p = 20$. This improvement demonstrates the effectiveness of the proposed optimization approach that allocates more features to polynomial degrees with larger variance reduction.

Variance reduction by complex features. Complex TensorSRHT (“TensorSRHT Comp.”) achieves significantly smaller errors than the real TensorSRHT (“TensorSRHT”), in particular for higher polynomial degrees p . Observe that even the complex TensorSRHT with $D = d$ achieves smaller errors than the real TensorSRHT with $D = 2d$ for the non-centered data. These improvements show the effectiveness of complex features in variance reduction, corroborating the preliminary results shown in Figures 3.2 and 3.3.

The optimized Maclaurin approach using complex features (“Opt. Macl. TensorSRHT Comp”) also achieves smaller errors than the optimized Maclaurin approach using real features (“Opt. Macl. TensorSRHT”), but the improvements are relatively smaller than those for TensorSRHT. These milder improvements would be because the optimized Maclaurin approach tends to assign more features to lower degree polynomials, but the improvements for these lower degree polynomials by complex features are relatively smaller than for higher degree polynomials.

Effectiveness of complex features on non-negative data. The improvements by complex features are more significant for the non-centered data than those for the centered-data. The non-centered data here consist of *non-negative* input vectors, as FashionMNIST consists of such vectors. This observation agrees with the discussion in Section 3.3 in Chapter 3 suggesting that complex features yield an approximate kernel whose variance is smaller than that of real features, if the input vectors are non-negative.

TensorSRHT v.s. TensorSketch. While the real TensorSRHT produces larger errors than TensorSketch for all the cases except $p = 3$, the complex TensorSRHT outperforms TensorSketch for all the cases. This comparison shows the use of complex features can make TensorSRHT competitive to the state-of-the-art (and one can further improve its performance by using it in the optimized Maclaurin approach).

5.4.3 Systematic Evaluation of the Optimized Maclaurin Approach

Lastly, we systematically evaluate the performance of the optimized Maclaurin approach against Spherical Random Features (Pennington et al., 2015) and random Fourier features (Rahimi and Recht, 2007). We run experiments on approximate GP classification and regression on a variety of datasets, using a high-degree polynomial kernel and the Gaussian kernel.

Optimized Maclaurin approach. We consider the optimized Maclaurin approach in Section 5.3 using the real Rademacher sketch (“Opt. Macl. Radem.”), one using the real TensorSRHT (“Opt. Macl. TensorSRHT”) and its complex extension (“Opt. Macl. TensorSRHT Comp.”).

Baselines. We use here approximation approaches based on Random Fourier Features (RFF) (Rahimi and Recht, 2007) and their extensions such as *Spherical Random Features* (SRF) (Pennington et al., 2015) and *Structured Orthogonal Random Features* (SORF) (Yu et al., 2016) as baselines. The latter two approaches constitute the state-of-the-art and we refer the reader to Chapter 2 for an explanation.

These approaches generate a set of frequency samples $\omega_1, \dots, \omega_{D/2} \in \mathbb{R}^d$ (suppose D is even for simplicity) from a certain spectral density, and construct a

feature map⁶ of dimension D as, for any $\mathbf{x} \in \mathbb{R}^d$,

$$\Phi_{\mathcal{R}}(\mathbf{x}) = \sqrt{\frac{2}{D}} [\cos(\mathbf{w}_1^\top \mathbf{x}), \dots, \cos(\mathbf{w}_{D/2}^\top \mathbf{x}), \sin(\mathbf{w}_1^\top \mathbf{x}), \dots, \sin(\mathbf{w}_{D/2}^\top \mathbf{x})]^\top \in \mathbb{R}^D. \quad (5.24)$$

Each approach has its own way of generating the frequency samples $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_{D/2}$: the original RFF generates them in an i.i.d. manner from the spectral density of a kernel, SORF uses structured orthogonal matrices (thus we may call it ‘‘RFF Orth.’’), and SRF uses a certain optimized spectral density.

For a thorough comparison, we also consider a complex version of these RFF-based approaches. By generating frequency samples $\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_D \in \mathbb{R}^d$ in the specific way of each approach, one can define a corresponding complex feature map as, for any $\mathbf{x} \in \mathbb{R}^d$,

$$\Phi_{\mathcal{C}}(\mathbf{x}) := \sqrt{\frac{1}{D}} [\exp(i\boldsymbol{\omega}_1^\top \mathbf{x}), \dots, \exp(i\boldsymbol{\omega}_D^\top \mathbf{x})]^\top \in \mathbb{C}^D. \quad (5.25)$$

One can see⁷ that Eq. (5.25) is a complex version of Eq. (5.24) by defining an approximate kernel with $\Phi_{\mathcal{C}}(\mathbf{x})$ and taking its real part, which recovers Eq. (5.24) of dimension $2D$.

Approximate GP Inference with a High-degree Polynomial Kernel

We first consider approximate GP classification and regression with a high-degree polynomial kernel.

Setting. We set the polynomial degree as $p = 20$, to make it challenging to approximate the polynomial kernel. We apply zero-centering to each dataset (i.e., we subtract the mean of input vectors from each input vector), as it improves the MNLL values on most datasets (see Appendix D.2 for supplementary experiments). We evaluate all the four error metrics in Section 5.4.1, including the relative Frobenius norm error in Eq. (5.23). For each approach, the number of random features is $D \in \{d, 3d, 5d\}$ with d being the dimensionality of input vectors.

⁶There is another popular version of the feature map in Eq. (5.24) defined as $\Phi_{\mathcal{R}}(\mathbf{x}) = \sqrt{\frac{2}{D}} [\cos(\mathbf{w}_1^\top \mathbf{x} + b_1), \dots, \cos(\mathbf{w}_D^\top \mathbf{x} + b_D)]^\top \in \mathbb{R}^D$ with b_1, \dots, b_D uniformly sampled on $[0, 2\pi]$. Following Sutherland and Schneider (2015) who suggested the superiority of Eq. (5.24), we use Eq. (5.24) here in all the methods using RFF, including SRF and SORF.

⁷Define an approximate kernel with Eq. (5.25) as $\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi_{\mathcal{C}}(\mathbf{x})^\top \overline{\Phi_{\mathcal{C}}(\mathbf{y})} = \frac{1}{D} \sum_{i=1}^D \exp(i\boldsymbol{\omega}_i^\top (\mathbf{x} - \mathbf{y})) = \frac{1}{D} \sum_{i=1}^D \exp(i\boldsymbol{\omega}_i^\top \mathbf{x}) \exp(-i\boldsymbol{\omega}_i^\top \mathbf{y})$. By taking its real part, we have $\text{Re}\{\hat{k}(\mathbf{x}, \mathbf{y})\} = \frac{1}{D} \sum_{i=1}^D \cos(\boldsymbol{\omega}_i^\top (\mathbf{x} - \mathbf{y})) = \frac{1}{D} \sum_{i=1}^D (\cos(\boldsymbol{\omega}_i^\top \mathbf{x}) \cos(\boldsymbol{\omega}_i^\top \mathbf{y}) + \sin(\boldsymbol{\omega}_i^\top \mathbf{x}) \sin(\boldsymbol{\omega}_i^\top \mathbf{y})) =: \Phi_{\mathcal{R}}(\mathbf{x})^\top \Phi_{\mathcal{R}}(\mathbf{y})$, where $\Phi_{\mathcal{R}}(\mathbf{x}) := \sqrt{\frac{1}{D}} [\cos(\boldsymbol{\omega}_1^\top \mathbf{x}), \dots, \cos(\boldsymbol{\omega}_D^\top \mathbf{x}), \sin(\boldsymbol{\omega}_1^\top \mathbf{x}), \dots, \sin(\boldsymbol{\omega}_D^\top \mathbf{x})]^\top \in \mathbb{R}^{2D}$ is the $2D$ -dim. version of Eq. (5.24).

Baselines. As a baseline, we use SRF (Pennington et al., 2015), a state-of-the-art approach to approximating polynomial kernels defined on the *unit sphere* in \mathbb{R}^d . Pennington et al. (2015) show that SRF works particularly well for approximating high degree polynomial kernels, and significantly outperforms the Random Maclaurin approach (Kar and Karnick, 2012) and TensorSketch (Pham and Pagh, 2013) for such kernels.

We also consider two other extensions of SRF for baselines. SRF generates the frequency samples $\omega_1, \dots, \omega_{D/2}$ in Eq. (5.24) from an optimized spectral density, by first drawing samples from the unit sphere in \mathbb{R}^d . Therefore, by replacing these samples on the unit sphere by structured orthogonal projections of SORF (Yu et al., 2016), one can construct a structured version of SRF. We use this structured SRF as another baseline (“SRF Orth.” in Fig. 5.4). Moreover, we consider a complex extension of the structured SRF in the form of Eq. (5.25) (“SRF Orth. Comp.” in Fig. 5.4). While these extensions are themselves novel, we include them in the experiments, as they improve over the vanilla SRF and make the experiments more competitive.

Fig. 5.4 describes the results of approximate GP classification on four datasets from Table 5.2. We present the results on the other four datasets as well as the results of GP regression in Appendix D.2 to save the space. We can make the following observations from these results.

Relative Frobenius norm error. For most cases, the optimized Maclaurin approaches with TensorSRHT achieve lower relative Frobenius norm errors than the SRF approaches. Specifically, “Opt. Macl. TensorSRHT” outperforms “SRF” and “SRF Orth.”, and “Opt. Macl. TensorSRHT Comp.” outperforms “SRF Orth. Comp.”

KL divergence. While the optimized Maclaurin approaches achieve lower KL divergences than the SRF approaches for most cases, the margins are smaller than those for the relative Frobenius norm errors. One possible reason is that the Maclaurin approaches in general (either random or optimized) can be inaccurate in approximating the GP posterior variances at test inputs far from $\mathbf{x} = \mathbf{0}$, as discussed in Section 5.3.3.

Classification errors and mean negative log likelihood (MNLL). The optimized Maclaurin approaches with TensorSRHT achieve equal or lower classification errors and MNLL than the SRF approaches. These results suggest that the optimized Maclaurin approaches are promising not only in kernel approximation accuracy but also in downstream task performance. Recall that we selected the regularization parameter in GP classification by maximizing the MNLL of *SRF* (on the validation set), and used the same regularization parameter in the other approaches (See Section 5.4.1). Therefore, the results of Fig. 5.4 are in favor of

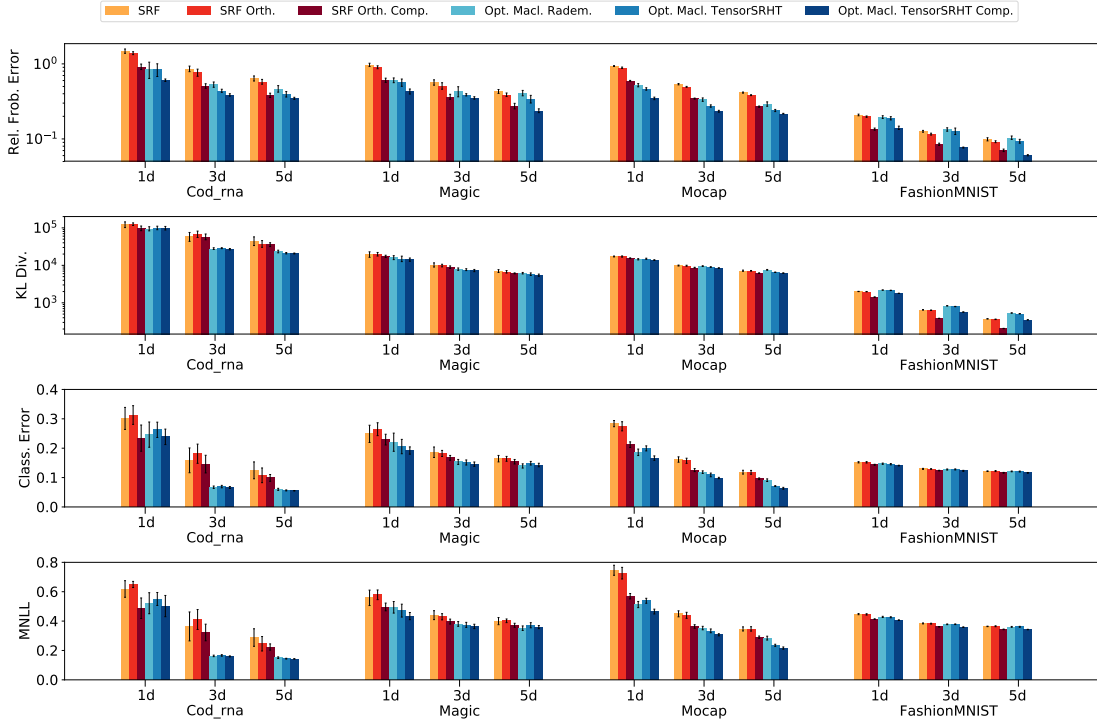


Figure 5.4: Results of the experiments in Section 5.4.3 on approximate GP classification with a high-degree polynomial kernel. Lower values are better for all the metrics. For each dataset, we show the number of random features $D \in \{1d, 3d, 5d\}$ used in each method on the horizontal axis, with d being the input dimensionality of the dataset. We put the legend labels and the bars in the same order.

the SRF approaches, and the optimized Maclaurin approaches may perform even better if we choose the regularization parameter for them separately.

Approximate GP Inference with a Gaussian kernel

We next consider GP classification using a Gaussian kernel. As in Section 5.4.3, we apply zero-centring to the input vectors of each dataset.

Baselines. We use RFF, SORF (“RFF Orth.”) and a complex extension of SORF (“RFF Orth. Comp.”) as baselines (see the beginning of Section 5.4.3 for details). SORF is a state-of-the-art approach to approximating a Gaussian kernel (e.g. Choromanski et al., 2018). As in Section 5.4.3, we consider its complex extension to make the experiments more competitive.

Results. Fig. 5.5 summarizes the results on four datasets from Table 5.2. We show the results on the rest of datasets as well as the results of GP regression in Appendix D.2. We can make similar observations for Fig. 5.5 as for the poly-

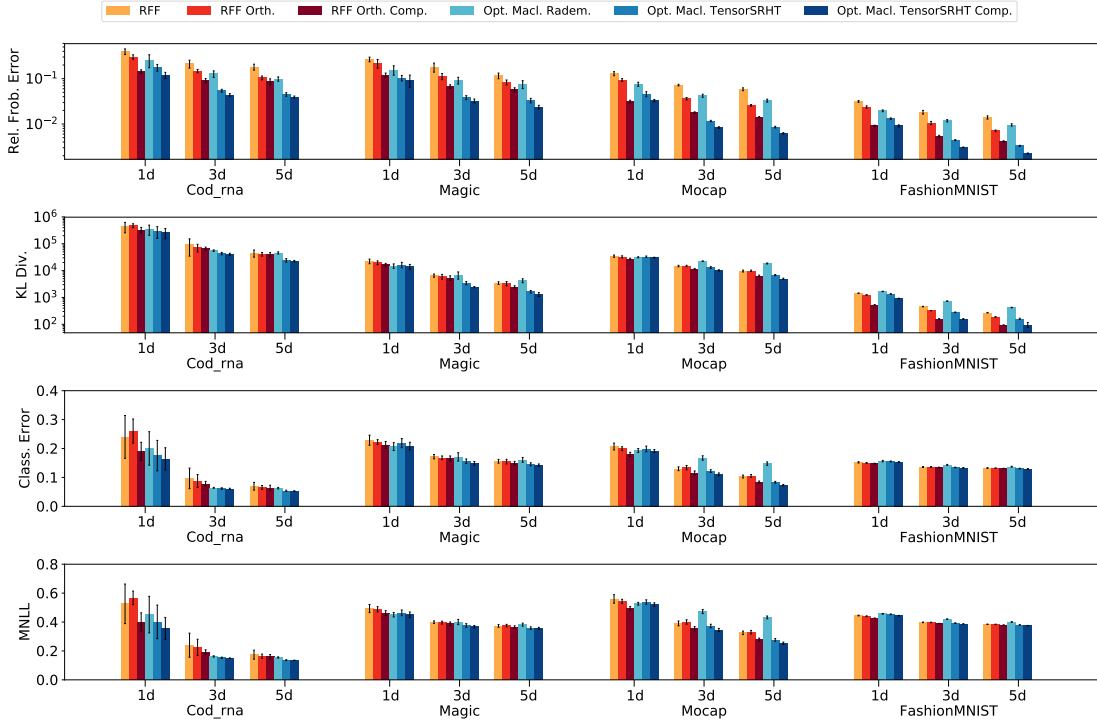


Figure 5.5: Results of the experiments in Section 5.4.3 on approximate GP classification with a Gaussian kernel. Lower values are better for all the metrics. For each dataset, we show the number of random features $D \in \{1d, 3d, 5d\}$ used in each method on the horizontal axis, with d being the input dimensionality of the dataset. We put the legend labels and the bars in the same order.

nomial kernel experiments in Section 5.4.3 (and thus we omit explaining them). The results suggest the effectiveness of the optimized Maclaurin approach with TensorSRHT in approximating the Gaussian kernel.

Influence of the Data Distribution on the Kernel Approximation

Lastly, we investigate a characterization of datasets for which the optimized Maclaurin approach performs well. We focus on polynomial kernel approximation, and make a comparison with SRF as in Section 5.4.3.

Fig. 5.6 describes a histogram of pairwise distances $\{\|\mathbf{x}_{*,i} - \mathbf{x}_{*,j}\|\}_{i \neq j}$ of the input vectors in a test subset $X_{*,\text{sub}} = \{\mathbf{x}_{*,1}, \dots, \mathbf{x}_{*,m_*}\}$, obtained after zero-centering and unit-normalization, of each of four representative datasets (kin8nm, Cod_rna, Naval, and Protein). For these datasets, the optimized Maclaurin approach and SRF show stark contrasts in their performances; see Section 5.4.3 and Appendix D.2. Note that the polynomial kernel in Eq. (5.22) is a shift-invariant kernel on the unit sphere of \mathbb{R}^d , and thus its value depends only on the distance $\tau := \|\mathbf{x} - \mathbf{y}\|$ between the input vectors \mathbf{x}, \mathbf{y} as long as $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$. This

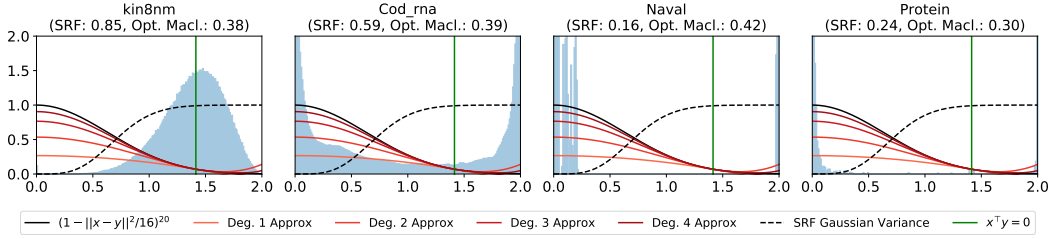


Figure 5.6: Histograms of pairwise Euclidean distances $\{\|\mathbf{x}_{*,i} - \mathbf{x}_{*,j}\|\}_{i \neq j}$ for test subsets of four datasets (Section 5.4.3). On the top of each figure, we show the relative Frobenius norm errors (5.23) of the optimized Maclaurin approach with real TensorSRHT and of SRF with structured orthogonal projections. The black curve represents the polynomial kernel in Eq. (5.22) with $p = 20$ as a function of $\tau := \|\mathbf{x} - \mathbf{y}\|$ (the horizontal axis); the orange curves describe its degree $n \in \{1, 2, 3, 4\}$ approximations (i.e., the truncation of the Maclaurin expansion (5.3) of the polynomial kernel up to the n -th degree terms.). The dashed curve represents the variance of the SRF approximation as a function of $\tau = \|\mathbf{x} - \mathbf{y}\|$. The green vertical line shows the value of $\tau = \|\mathbf{x} - \mathbf{y}\| = \sqrt{2}$ for which the input vectors \mathbf{x}, \mathbf{y} are orthogonal, $\mathbf{x}^\top \mathbf{y} = 0$.

motivates us studying here the distribution of pairwise distances and its effects on approximating the polynomial kernel in Eq. (5.22).

In Fig. 5.6, the optimized Maclaurin approach yields lower relative Frobenius norm errors (5.23) than SRF for the left two plots, while the optimized Maclaurin approach is less accurate than SRF for the right two plots. For the datasets of the right two plots (Naval and Protein), the pairwise distances $\{\|\mathbf{x}_{*,i} - \mathbf{x}_{*,j}\|\}_{i \neq j}$ concentrate around $\tau = 0$ (and there is a smaller mass around $\tau = 2$). In comparison, for the datasets of the left two plots (kin8nm and Cod_rna), the pairwise distances are relatively more evenly distributed across the possible range $\tau \in [0, 2]$.

The above observation suggests that the optimized Maclaurin approach is more suitable for datasets in which the pairwise distances $\{\|\mathbf{x}_{*,i} - \mathbf{x}_{*,j}\|\}_{i \neq j}$ are not concentrating around 0, i.e., datasets in which there is a diversity in the input vectors $\{\mathbf{x}_{*,1}, \dots, \mathbf{x}_{*,m_*}\}$. In fact, for approximating the polynomial kernel (black curve in Fig. 5.6), the finite-degree Maclaurin approximations (orange curves) tend to be less accurate for input vectors \mathbf{x}, \mathbf{y} close to each other, $\tau = \|\mathbf{x} - \mathbf{y}\| \approx 0$, and become relatively more accurate as input vectors \mathbf{x}, \mathbf{y} approach orthogonality, i.e. $\mathbf{x}^\top \mathbf{y} = 0$ (or $\tau = \|\mathbf{x} - \mathbf{y}\| = \sqrt{2}$; the vertical green line); see also the Maclaurin expansion (5.3) of the polynomial kernel. On the other hand, the variance of SRF is the lowest around $\tau = \|\mathbf{x} - \mathbf{y}\| = 0$ and increases as τ tends to 2. Therefore, the SRF performs well if the pairwise distances $\{\|\mathbf{x}_{*,i} - \mathbf{x}_{*,j}\|\}_{i \neq j}$ concentrate around 0, and may become inaccurate if they do not.

5.5 Pathology when Approximating the Gaussian Kernel

In this section, we uncover a pathology related to Maclaurin-based approximations of the Gaussian kernel that arises in particular when modelling high-frequency data characterized by a short length scale. This pathology causes the approximate GPR predictor in Section 5.3.3 to collapse away from $x = 0$ (see especially the left plot of Fig. 5.2). We will later on propose a cure for this pathology avoiding the collapsing GP predictions.

For ease of notation, we now absorb the length scale parameter $l > 0$ of the Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/(2l^2))$ in the input data, i.e., we define $\tilde{\mathbf{x}} := \mathbf{x}/l$ and $\tilde{\mathbf{y}} := \mathbf{y}/l$. In the following theorem, we show that Maclaurin-based kernel approximations vanish when the norms of the inputs become large. Since this effect is only related to the truncation of the Maclaurin series, it happens also without random feature approximations. The vanishing kernels in turn cause GP predictions to collapse, as we discuss later.

Theorem 5.5.1 (Vanishing Maclaurin approximation of Gaussian kernels). *The magnitude of the finite Maclaurin approximation $k_p(\mathbf{x}, \mathbf{y}) := \sigma^2 \exp(-(\|\tilde{\mathbf{x}}\|^2 + \|\tilde{\mathbf{y}}\|^2)/2) \sum_{n=0}^p 1/n! (\tilde{\mathbf{x}}^\top \tilde{\mathbf{y}})^n$ of the Gaussian kernel approaches zero as $\|\tilde{\mathbf{x}}\|^2 + \|\tilde{\mathbf{y}}\|^2$ increases. The error $|k(\mathbf{x}, \mathbf{y}) - k_p(\mathbf{x}, \mathbf{y})|$ between the exact kernel k and its approximation k_p is the largest for parallel \mathbf{x}, \mathbf{y} and zero when they are orthogonal.*

Proof. We start by deriving an upper bound for $|k_p(\mathbf{x}, \mathbf{y})|$. We leave out the length scale here for ease of notation:

$$\begin{aligned} |k_p(\mathbf{x}, \mathbf{y})| &= \sigma^2 \exp(-(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2) \left| \sum_{n=0}^p \frac{1}{n!} (\mathbf{x}^\top \mathbf{y})^n \right| \\ &\leq \sigma^2 \exp(-(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2) \sum_{n=0}^p \frac{1}{n!} (\|\mathbf{x}\| \|\mathbf{y}\|)^n \end{aligned} \quad (5.26)$$

Next, we notice that $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\mathbf{x}^\top \mathbf{y} \geq 0$ for any \mathbf{x}, \mathbf{y} . Thus, we can choose them to be parallel. So $\|\mathbf{x}\| \|\mathbf{y}\| \leq (\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2$. From this inequality, it follows

$$\underbrace{\sum_{n=0}^p \frac{1}{n!} (\|\mathbf{x}\| \|\mathbf{y}\|)^n}_{(A)} \leq \underbrace{\sum_{n=0}^p \frac{1}{n!} ((\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2)^n}_{(B)} \leq \underbrace{\exp((\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2)}_{(C)}.$$

Now, the gap $(C) - (B) = \sum_{n=p+1}^{\infty} \frac{1}{n!} ((\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2)^n$ increases as $\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ increases, which must decrease the ratio $(A)/(C)$ and thus the upper bound of $|k_p(\mathbf{x}, \mathbf{y})|$ (5.26) goes to zero as $\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ increases.

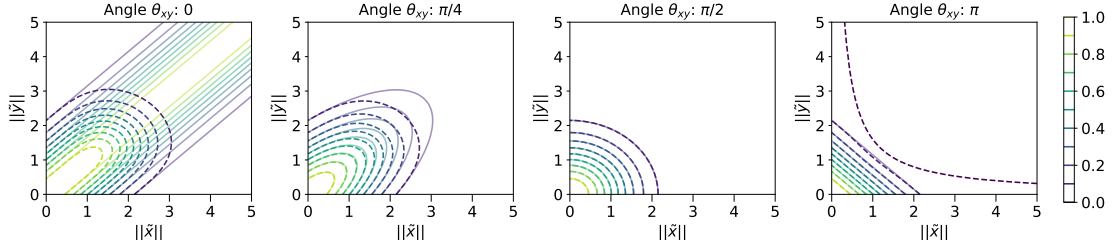


Figure 5.7: $k_p(\mathbf{x}, \mathbf{y})$ (dashed) vs. $k(\mathbf{x}, \mathbf{y})$ (transparent and solid) for fixed $p = 3$ and $\sigma^2 = l = 1$, where we use the relationship $\tilde{\mathbf{x}}^\top \tilde{\mathbf{y}} = \|\tilde{\mathbf{x}}\| \|\tilde{\mathbf{y}}\| \cos(\theta_{xy})$.

Next, we look at the approximation error related to the truncation of the Maclaurin series:

$$|k(\mathbf{x}, \mathbf{y}) - k_p(\mathbf{x}, \mathbf{y})| = \sigma^2 \exp(-(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)/2) \left| \sum_{n=p+1}^{\infty} \frac{1}{n!} (\mathbf{x}^\top \mathbf{y})^n \right|$$

If the angle between \mathbf{x} and \mathbf{y} is zero such that $\mathbf{x}^\top \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\|$, all addends in the infinite sum are maximized. The error thus becomes the largest. The error is zero when they are orthogonal. ■

We visualize the implications of [Theorem 5.5.1](#) in [Fig. 5.7](#), where we compare the approximation $k_p(\mathbf{x}, \mathbf{y})$ with the exact Gaussian kernel $k(\mathbf{x}, \mathbf{y})$ for $p = 3$ over a range of values $\|\tilde{\mathbf{x}}\|, \|\tilde{\mathbf{y}}\|$ as well as the angle θ_{xy} between \mathbf{x} and \mathbf{y} . One can see that $k_p(\mathbf{x}, \mathbf{y})$ approaches zero with increasing $\|\tilde{\mathbf{x}}\|^2 + \|\tilde{\mathbf{y}}\|^2$ regardless of θ_{xy} . This deteriorates the approximation quality, in particular as θ_{xy} goes to zero. This development accelerates when choosing a shorter length scale l .

We now discuss the implications of [Theorem 5.5.1](#) for the case of Gaussian process regression. Recall from [Chapter 1](#) that the GPR predictor is given through

$$\mu_* := \mathbf{k}_{\mathbf{f}*}^\top (\mathbf{K}_{\mathbf{ff}} + \sigma_{\text{noise}}^2 \mathbf{I})^{-1} \mathbf{y} \quad \text{and} \quad \sigma_*^2 := k_{**} - \mathbf{k}_{\mathbf{f}*}^\top (\mathbf{K}_{\mathbf{ff}} + \sigma_{\text{noise}}^2 \mathbf{I})^{-1} \mathbf{k}_{\mathbf{f}*}, \quad (5.27)$$

where $\mathbf{K}_{\mathbf{ff}}$ is the kernel matrix evaluated on the training data $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{k}_{\mathbf{f}*} = (k(\mathbf{x}_1, \mathbf{x}_*), \dots, k(\mathbf{x}_N, \mathbf{x}_*))^\top$, $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$ and $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the identity matrix.

When using k_p , the GPR predictive means μ_* and variances σ_*^2 in [Eq. \(5.27\)](#) collapse to zero for test points $\tilde{\mathbf{x}}^*$ with large $\|\tilde{\mathbf{x}}^*\|$ since $k_p(\mathbf{x}^*, \mathbf{x})$ goes to zero for any $\mathbf{x} \in \mathbb{R}^d$ in this case. This effect is shown in the left plot of [Fig. 5.2](#). We will discuss this example in greater detail in [Section 5.5.3](#).

A similar pathology was identified for the *Relevance Vector Machine* ([Tipping, 1999](#)) that was solved by adding new basis functions centered on the test points ([Rasmussen and Quiñonero Candela, 2005](#)). In the following, we present a similar strategy by centering our entire approximation around individual test points.

5.5.1 Curing the Pathology

We will now exploit a property of the Gaussian kernel that allows us to cure the aforementioned pathology.

The Gaussian kernel is shift-invariant, i.e., $k(\mathbf{x} + \boldsymbol{\delta}, \mathbf{y} + \boldsymbol{\delta}) = k(\mathbf{x}, \mathbf{y})$ for any $\boldsymbol{\delta} \in \mathbb{R}^d$, because $\|(\mathbf{x} + \boldsymbol{\delta}) - (\mathbf{y} + \boldsymbol{\delta})\| = \|\mathbf{x} - \mathbf{y}\|$. Thus, when making a prediction at a test input \mathbf{x}_* , one can subtract \mathbf{x}_* from all inputs used in Eq. (5.27) without changing the result of the prediction. More specifically, the values μ_* and σ_*^2 in Eq. (5.27) do not change if we substitute $k(\mathbf{x}, \mathbf{y})$ by $k(\mathbf{x} - \mathbf{x}_*, \mathbf{y} - \mathbf{x}_*)$ for the computation of $\mathbf{k}_{\mathbf{f}*}$, $\mathbf{K}_{\mathbf{ff}}$ and k_{**} .

Now we consider an approximate Gaussian kernel with random features $\hat{k}_p(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$ with truncation degree p of the Maclaurin series as defined in Section 5.3.1. Thus, for a set of polynomial sketches $\{\Phi_n(\mathbf{x})\}_{n=1}^p$ of dimensions $\{D_n\}_{n=1}^p$, the approximate kernel yields

$$\hat{k}_p(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp\left(-\frac{\|\mathbf{x}\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{y}\|^2}{2l^2}\right) \left(1 + \sum_{n=1}^p \frac{1}{n!l^{2n}} \Phi_n(\mathbf{x})^\top \Phi_n(\mathbf{y})\right) \quad (5.28)$$

and we have $\mathbb{E}[\hat{k}_p(\mathbf{x}, \mathbf{y})] = k_p(\mathbf{x}, \mathbf{y})$. The dimension of the feature map $\Phi(\mathbf{x})$ here is thus $D = 1 + D_1 + \dots + D_p$.

The approximate kernel \hat{k}_p (5.28) is greatly affected by the subtraction of the test point \mathbf{x}_* . To see this, we define $\hat{k}_p^*(\mathbf{x}, \mathbf{y}) := \hat{k}_p(\mathbf{x} - \mathbf{x}_*, \mathbf{y} - \mathbf{x}_*)$ with:

$$\hat{k}_p^*(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}^*\|^2 + \|\mathbf{y} - \mathbf{x}^*\|^2}{2l^2}\right) \sum_{n=1}^p \frac{1}{n!} \Phi_n\left(\frac{\mathbf{x} - \mathbf{x}^*}{l}\right)^\top \Phi_n\left(\frac{\mathbf{y} - \mathbf{x}^*}{l}\right) \quad (5.29)$$

where $\mathbb{E}[\hat{k}_p^*(\mathbf{x}, \mathbf{y})] = k_p(\mathbf{x} - \mathbf{x}^*, \mathbf{y} - \mathbf{x}^*) =: k_p^*(\mathbf{x}, \mathbf{y})$. As all $\{\Phi_n\}_{n=1}^p$ are sampled independently, the variance of $\hat{k}_p^*(\mathbf{x}, \mathbf{y})$ (5.29) can be written as

$$\mathbb{V}[\hat{k}_p^*(\mathbf{x}, \mathbf{y})] \propto \sum_{n=1}^p \left(\frac{1}{n!}\right)^2 \mathbb{V}\left[\Phi_n\left(\frac{\mathbf{x} - \mathbf{x}^*}{l}\right)^\top \Phi_n\left(\frac{\mathbf{y} - \mathbf{x}^*}{l}\right)\right], \quad (5.30)$$

When setting $\mathbf{x} = \mathbf{x}^*$ or $\mathbf{y} = \mathbf{x}^*$, the variance terms in Eq. (5.30) become *zero* for any of the polynomial sketches discussed in Chapter 3 as can be seen from Table 5.1. $\hat{k}_p^*(\mathbf{x}^*, \mathbf{y})$, $\hat{k}_p^*(\mathbf{x}, \mathbf{x}^*)$ and $\hat{k}_p^*(\mathbf{x}^*, \mathbf{x}^*)$ thus become *deterministic*.

We further have $\hat{k}_p^*(\mathbf{x}^*, \mathbf{y}) = \sigma^2 \exp(-\|\mathbf{x}^* - \mathbf{y}\|^2 / (2l)^2)$, $\hat{k}_p^*(\mathbf{x}, \mathbf{x}^*) = \sigma^2 \exp(-\|\mathbf{x} - \mathbf{x}^*\|^2 / (2l)^2)$ and $\hat{k}_p^*(\mathbf{x}^*, \mathbf{x}^*) = \sigma^2$, which are all equal to the exact kernel k evaluated at these points. Therefore, $\mathbf{k}_{\mathbf{f}*}$ and k_{**} in Eq. (5.27) become exact for our Maclaurin approximation. $\mathbf{K}_{\mathbf{ff}}$ unfortunately remains affected by the vanishing approximate kernels described by Theorem 5.5.1 and by non-zero random feature variances. A crucial advantage of using \hat{k}_p^* instead of \hat{k}_p is that the GP predictive distribution Eq. (5.27) does not collapse to zero anymore as $\|\tilde{\mathbf{x}}^*\|$ grows. We illustrate this effect on the same synthetic example that we used in Fig. 5.2.

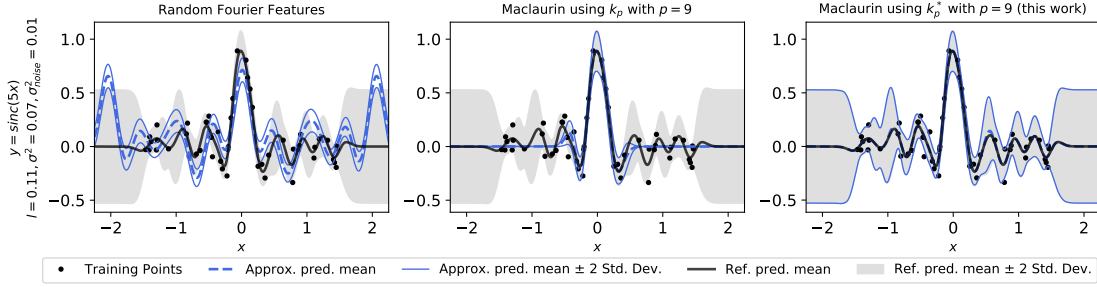


Figure 5.8: Approximating the predictive distribution of a reference GPR using different approximation methods and $D = 10$.

As before, we show the reference GPR along with three different approximation schemes in Fig. 5.8. The baseline Random Fourier features (RFF) [Rahimi and Recht \(2007\)](#) (left) struggles to recover the reference GPR for $D = 10$ random features. The dimension D of the feature map here is equal to $p + 1$ for the optimized Maclaurin approximation, since the example is one-dimensional and a single random feature perfectly approximates a degree- p polynomial kernel for the Rademacher sketch in this case. So $p = 9$ with $D_1 = \dots = D_p = 1$ is chosen by the optimized Maclaurin method.

The Maclaurin approach using \hat{k}_p^* (right) gives the best approximations while the predictive distribution of the one using \hat{k}_p (middle) collapses to zero very quickly at points $x^* \in \mathbb{R}$ away from zero. This is because $\|\tilde{x}^*\|$ becomes large very quickly and \hat{k}_p vanishes when being evaluated at these points. For values x^* far away from the training data, the GP predictor (5.27) using k_p^* even recovers the GPR prior distribution as desired, which can be explained as follows. When x^* is far from the training data, $\mathbf{k}_{\mathbf{f}^*}$ becomes zero. Then $\mu_* = 0$ and $\sigma_*^2 = k_{**} = \sigma^2$ in Eq. (5.27). Since k_{**} and $\mathbf{k}_{\mathbf{f}^*}$ in Eq. (5.27) are accurate when using k_p^* as explained earlier, the convergence to the prior is kept for k_p^* .

5.5.2 Reducing Computational Costs through Clustering

A caveat of using the approximate kernel $\hat{k}_p^*(\mathbf{x}, \mathbf{y})$ (5.29) described in Section 5.5.1 is that we need to recompute the GPR predictor (5.27) for every test point \mathbf{x}^* separately, which becomes expensive when many test points need to be predicted, even when using random features.

Recall from Chapter 1 that the GPR predictor using random features can be written as

$$\mu_* := \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{X})^\top \mathbf{y} / \sigma_{\text{noise}}^2 \quad \text{and} \quad \sigma_*^2 := \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{x}_*), \quad (5.31)$$

with $\mathbf{A} := \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) / \sigma_{\text{noise}}^2 + \mathbf{I}$ and $\Phi(\mathbf{X}) = (\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_N))^\top \in \mathbb{R}^{N \times D}$.

We denote the number of test inputs by N_* . A direct computation of Eq. (5.31) now costs a total of $\mathcal{O}(N_* N D^2)$ for all test points.

Algorithm 7: Precomputing Localized Maclaurin Features for the Gaussian Kernel

Input: Hyperparameters $l, \sigma^2, \sigma_{\text{noise}}^2 > 0$; training data $\{\mathbf{x}_i\}_{i=1}^N$; test data $\{\mathbf{x}_i^*\}_{i=1}^{N^*}$; number of features $D \geq 1$;
// 1) Training
if $d > 1$ **then** // We use random features
| Center the training data by subtracting the training mean ;
| Find $(p^*, D_1, \dots, D_{p^*})$, the optimal feature allocation for the random Maclaurin method (see Algorithm 5) ;
end
 $C := \{(\sum_{i=1}^N \mathbf{x}_i)/N\}$ // Initialize centroids with training mean ;
while *True* **do** // Farthest point clustering
| Let $\delta_i = \min\{\|\mathbf{x}_i - \mathbf{c}\| \mid \mathbf{c} \in C\}$ for $i = 1, \dots, N$;
| // Check if training Max.-Min.-Distance is below threshold
| **if** $\max\{\delta_i\}_{i=1}^N < \theta$ **then**
| | break;
| **end**
| Add \mathbf{x}_i with $i = \arg \max_i \{\delta_i\}_{i=1}^N$ to C ;
| Precompute $\mathbf{A}_c^{-1} := (\Phi(\mathbf{X} - \mathbf{c})^\top \Phi(\mathbf{X} - \mathbf{c})/\sigma_{\text{noise}}^2 + \mathbf{I})^{-1}$ with Φ defined as in Section 5.3.1 ;
end
forall $\{\mathbf{x}_i^*\}_{i=1}^{N^*}$ **do** // 2) Inference
| Assign \mathbf{x}_i^* to closest centroid $\mathbf{c} \in C$;
| Compute μ_* and σ_*^2 in Eq. (5.31) by setting $\mathbf{A}^{-1} = \mathbf{A}_c^{-1}$, $\mathbf{x}^* = \mathbf{x}_i^* - \mathbf{c}$ and $\mathbf{X} = \mathbf{X} - \mathbf{c}$;
end

The problem is “embarrassingly parallel” and the computation of Eq. (5.31) for every \mathbf{x}^* could be easily distributed on a cluster of compute nodes or parallelized using a single GPU, e.g., using JAX (Bradbury et al., 2018). However, we propose a different approach here that requires no separate training at test time while staying as close as possible to the training time of $\mathcal{O}(ND^2)$ as is generally desired for random feature approximations.

Our approach is to cluster the *training* data into N_C clusters and to use the centroids of these clusters as *pseudo test inputs*. We choose a farthest point clustering for simplicity as it does not require convergence verification and determines the number of clusters using a threshold $\theta > 0$, but any clustering algorithm can be used instead (even a random selection of training points). As the centroids are known during training, we can pretrain a set of N_C predictors using Eq. (5.31) and assign a new test point \mathbf{x}^* to the closest centroid at prediction time. We summarize the complete procedure in Algorithm 7. The computational cost is now

$\mathcal{O}(N_C ND^2)$ and is thus much lower than $\mathcal{O}(N_* ND^2)$ if $N_C \ll N_*$.

5.5.3 Empirical Evaluation on Real-World Data

In this section, we evaluate our proposed method on real-world data of different dimensions for which we employ the real-valued polynomial sketches summarized in Table 5.1. As for the synthetic example in Fig. 5.8, we use the Gaussian kernel with hyperparameters l and $\sigma^2 > 0$ that are found through gradient based optimization of the log marginal likelihood of a reference GPR, along with σ_{noise}^2 .

We compare our method against a random Fourier features baseline Rahimi and Recht (2007) as well as its structured extension Yu et al. (2016) when the data is sufficiently high dimensional⁸. We also add the vanilla optimized Maclaurin method without recentering the data to this comparison. It is equivalent to Algorithm 7 using only a single cluster with its centroid being the training mean. We measure the approximation quality with respect to the reference GPR using the KL-divergence between the predictive means and variances in Eq. (5.27) and Eq. (5.31). We measure downstream regression performance using the root mean squared error (RMSE).

UK Apartment Price Data

We downloaded the monthly property sales data for England and Wales from the HM land registry⁹. We filtered for sold apartments for the month of January 2022 leading to a data set with 24 553 observations. Matching the post codes for each apartment sold with a database of latitudes and longitudes¹⁰ allowed us to obtain a two-dimensional data set (latitude, longitude) that we could regress against the logarithm of the sales prices. We randomly split the data into 10 000 training points and kept the rest for testing.

In our first experiment, we aim to recover the reference GPR predictive distribution on a regular grid of latitudes (between $+50^\circ$ and $+55^\circ$) and longitudes (between -6° and $+2^\circ$) of size 100 by 100. Fig. 5.9 shows the results of this experiment. As for the sinc-example in Fig. 5.8, random Fourier features struggle to recover the predictive distribution using $D = 100$ random features and the vanilla Maclaurin method using \hat{k}_p (5.28) suffers from vanishing kernels due to the short (compared to the scaling of the data) length scale of $l = 0.25$. Our proposed kernel \hat{k}_p^* (5.29) improves predictions considerably leading to the lowest KL divergence with respect to the reference GP predictive distribution. It also converges to the prior for test points far from the training data.

In our second experiment, we evaluate the use of Algorithm 7 to precompute the matrix inversion in Eq. (5.31) on a set of pseudo test inputs. This time we

⁸Otherwise, the structured random Fourier features induce a large bias.

⁹<https://www.gov.uk/government/statistical-data-sets/price-paid-data-downloads>

¹⁰<https://www.freemaptools.com/download-uk-postcode-lat-lng.htm>

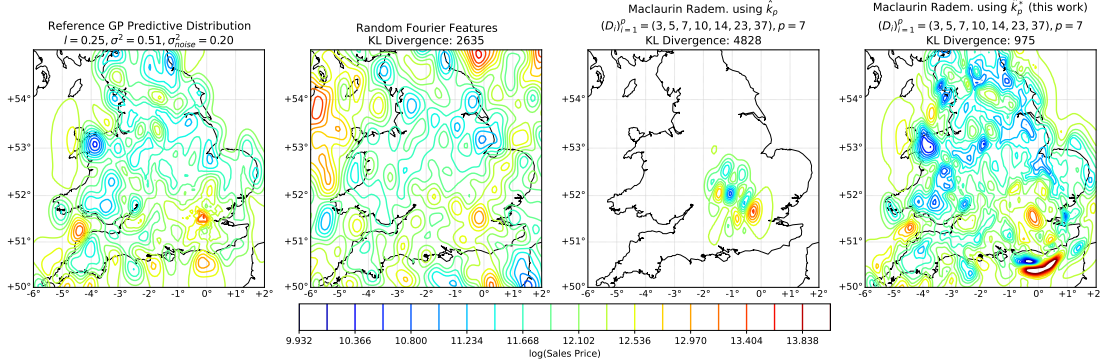


Figure 5.9: Approximating the predictive mean of a reference GP using different approximation methods and $D = 100$. KL Divergences also include predictive variances. The Maclaurin method using \hat{k}_p makes predictions centered around the mean of the training data (close to London).

report results on the left-out test data instead of a regular grid. The left part of Fig. 5.10 shows these results. We can see that the KL divergence falls off considerably (top plot) as we add more clusters until reaching 57 clusters. From then on the KL divergence remains roughly the same indicating that 57 clusters give a good trade-off between efficiency and performance. In the bottom plot we show a comparison of RMSE values for these 57 clusters, where the Maclaurin method outperforms random Fourier features, in particular for small D .

UCI Data Sets: Yacht and kin8nm

In the following, we repeat the evaluation of Algorithm 7 for two higher dimensional data sets that are taken from the UCI machine learning repository (Dua and Graff, 2017) in Fig. 5.10.

We obtain very similar results for the UCI Yacht data set as for the UK apartment price data set. Adding more clusters gives large gains initially but these diminish when setting $\theta > 2.0l$ in Algorithm 7 determining a good trade-off between performance and computational cost. This time we included structured orthogonal random Fourier features (Yu et al., 2016) that are also outperformed by the Maclaurin method.

For the UCI kin8nm data set, results look quite different. Adding more clusters *increases* the KL divergence towards the reference GP predictor, which is why only a single centroid, the mean of the training data, is chosen for the RMSE comparison in the plot below. This corresponds to the vanilla optimized Maclaurin method (Algorithm 5).

We explain this observation as follows. The length scales obtained for the sinc example, the UK apartment price data and for UCI Yacht are short. For the UCI Yacht data set it is 0.32, i.e., much less than 3.28, the median pairwise Euclidean

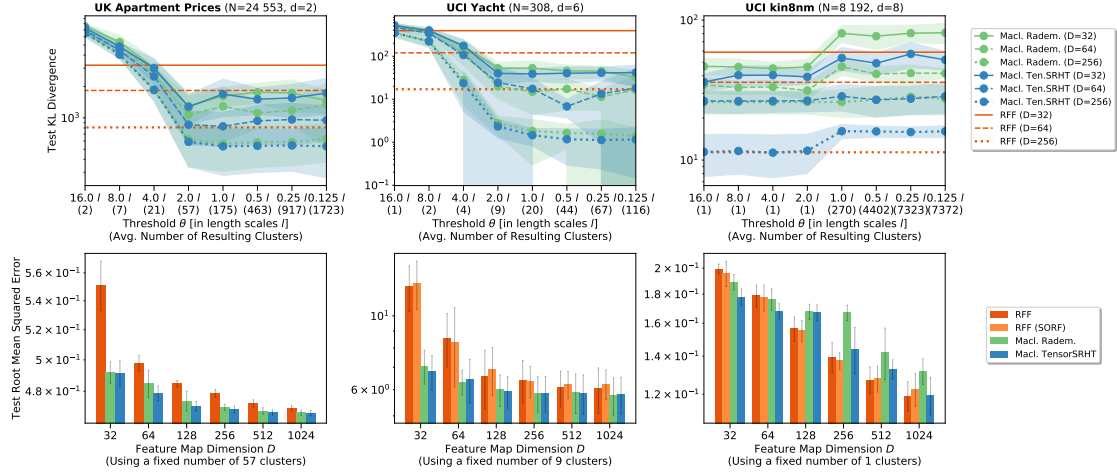


Figure 5.10: (Top) Test KL Divergence for a given number of clusters obtained from the training data using Algorithm 7 with threshold θ . (Bottom) RMSE over feature map dimension D for a fixed number of clusters.

distance of the training data, indicating that the data is fit by a reference GP of high frequency. For kin8nm the length scale is 2.15 compared to 3.93 (median heuristic) indicating a much smoother GP than the ones before.

In this case, the values of k_{f^*} and k_{**} in Eq. (5.27) are less affected by the vanishing kernels Theorem 5.5.1 due to a longer length scale. However, the approximation of \mathbf{K}_{ff} in Eq. (5.27) is more accurate when using the vanilla Maclaurin approximation (5.28) because the data is centered around the training mean. This shows that the clusters need to be chosen depending on the smoothness of the target GP. In this work, we have provided a *generalization* of the vanilla Maclaurin method that (with an appropriate choice of clusters) can fit both, high and low-frequency data.

5.6 Discussion

Based on the derived variance formulas in Chapter 3, we developed a novel optimization algorithm for data-driven random feature approximations of dot product kernels in this chapter, which is also applicable to the Gaussian kernel. Our approach uses a truncated Maclaurin approximation of the kernel, which approximates the kernel as a finite sum of polynomial kernels of different degrees. Given a total number of random features, our optimization algorithm determines how many random features should be used for each polynomial degree in the Maclaurin approximation. We defined the objective function of this optimization algorithm as the mean squared error w.r.t. data distribution, and used the variance formulas derived in Chapter 3 to evaluate the objective in closed form. We empirically demonstrated that our optimized Maclaurin approach achieves state-of-the-

art performance on a variety of data sets both in terms of the kernel approximation accuracy and downstream task performance.

However, we also showed that Maclaurin approximations suffer from a pathology when approximating the Gaussian kernel. In this case, the GP predictive distribution collapses at locations with large input norms, which happens particularly fast if the length scale is set to a small value relative to the scaling of the data, leading to poor predictions. We proposed a cure for this pathology by recentering the data around the test point to predict, yielding once again state-of-the-art performance on a variety of data sets.

The methods presented in this chapter improve the empirical performance of Maclaurin approximations of dot product kernels and Gaussian kernels considerably (e.g. compared to the random Maclaurin approach (Kar and Karnick, 2012)). However, [Algorithm 5](#) requires the evaluation of the exact kernel on a subset of the training data, where the size of the subset needs to be specified by the user. Moreover, the scaling of the data, or its rescaling via a length scale parameter, has a significant impact on the quality of the Maclaurin approximation. Future research could therefore investigate how the random feature allocation between polynomial degrees can be derived from the scale of the input data alone, without requiring pairwise kernel evaluations. Similarly, the number of clusters and their positions in [Algorithm 7](#) could be further improved with such considerations, leading to a simpler implementation of the methods described in this chapter.

Chapter 6

The Optical Processing Unit

So far, the contents of this thesis have revolved around random feature approximations for polynomial kernels and more general dot product kernels. This research was originally inspired by working with an optical hardware, the so-called Optical Processing Unit (OPU).

The OPU (Saade et al., 2016; Ohana et al., 2020) computes random projections literally at the speed of light and without having to store the projection matrix in memory, and it does so at a low power consumption. Therefore, applying the OPU in machine learning applications could potentially contribute to a more ecologically responsible data science infrastructure by saving time and energy. The OPU has already been applied in applications like reservoir computing (Dong et al., 2018, 2019) and anomaly detection (Keriven et al., 2018) demonstrating its usefulness in practice.

Building on the milestone work of Saade et al. (2016), the goal of this chapter is threefold: a) we derive in full generality the kernel to which the random features computed by the OPU converge, generalizing the earlier computation of Saade et al. (2016) to a larger class of kernels; b) we present new examples and a benchmark of applications for the kernel of the OPU; and c) we give a detailed comparison of the running time and energy consumption between the OPU and a last generation GPU.

In fact, we show in this chapter that the kernel implicitly computed by the OPU has a close relation to the polynomial kernel, making its derivation closely related to the previous contents of this thesis.

6.1 The Optical Processing Unit (OPU)

We begin by describing the inner workings of the OPU giving rise to the random features it computes. Fig. 6.1 shows a schematic drawing of the physical apparatus.

The principle of the random projections performed by the OPU is based on the use of a semi-transparent “scattering” material that refracts the laser light passing through it. As discussed in Liutkus et al. (2014), light passing through such

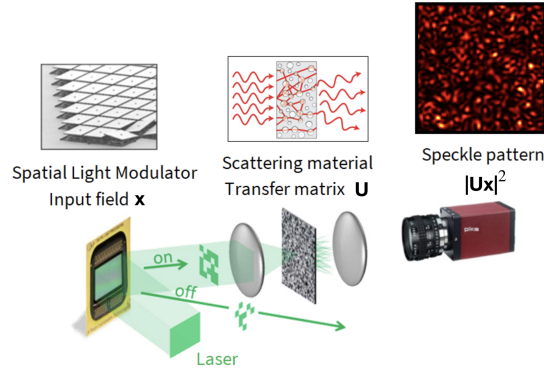


Figure 6.1: Experimental setup of the Optical Processing Unit (modified with permission from [Saade et al. \(2016\)](#)). The data vector is encoded in the coherent light from a laser using a DMD. Light then goes through a scattering medium and a speckle pattern is measured by a camera.

a scattering medium follows many complex paths that depend on the refractive index inhomogeneities at random positions. For a fixed scattering medium, the resulting process is still linear, deterministic, and reproducible.

We thus obtain linear projections of a numerical data vector $\mathbf{x} \in \mathbb{R}^d$ by filtering the light beam at different positions using a *digital micromirror device (DMD)*. This encoded light then passes through the heterogeneous medium, implicitly performing a random matrix multiplication, where the random matrix is inherently determined by the scattering medium. The scattered light results in a so-called *speckle pattern* that is recorded by a camera, where the light intensity at each measuring point is the result of a superposition of the light beams filtered by the DMD. This superposition is modelled mathematically by a sum of the components of \mathbf{x} weighted by random coefficients. Measuring the intensity of light at an arbitrary pixel location recorded by the camera induces a non-linear transformation of this sum. D arbitrary pixel values taken together then constitute the random feature map

$$\Phi(\mathbf{x}) = |\mathbf{U}\mathbf{x}|^2/\sqrt{D}, \quad \text{with } \mathbf{U} \in \mathbb{C}^{D \times d} \quad \text{and} \quad (\mathbf{U})_{i,j} \stackrel{i.i.d.}{\sim} \mathcal{CN}(0,1), \quad (6.1)$$

where $|\cdot|^2$ is the squared absolute value taken element-wise.

6.2 Computing the Kernel

In this section, we compute the kernel associated with the random feature map Φ (6.1), which is our first contribution. The goal is to derive the limiting kernel that the inner product $\hat{k}_2(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$ for two data points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ converges to as D goes to infinity.

Let $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_D)^\top$. We then write

$$\hat{k}_2(\mathbf{x}, \mathbf{y}) = \frac{1}{D} \sum_{\ell=1}^D |\mathbf{u}_\ell^\top \mathbf{x}|^2 |\mathbf{u}_\ell^\top \mathbf{y}|^2 \rightarrow \mathbb{E}[|\mathbf{u}^\top \mathbf{x}|^2 |\mathbf{u}^\top \mathbf{y}|^2] =: k_2(\mathbf{x}, \mathbf{y}) \quad (D \rightarrow \infty),$$

with $\mathbf{u} = (u_1, \dots, u_d)^\top \in \mathbb{C}^d$ and $u_1, \dots, u_d \stackrel{i.i.d.}{\sim} \mathcal{CN}(0, 1)$. The convergence here is due to the law of large numbers and defined as a convergence in probability, where $k_2(\mathbf{x}, \mathbf{y})$ is the limit of this convergence. Thus, we work out $k_2(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\hat{k}_2(\mathbf{x}, \mathbf{y})]$ and summarize the result in the following theorem.

Theorem 6.2.1. *The kernel k_2 approximated by the dot product of optical random features of Eq. (6.1) is given by*

$$k_2(\mathbf{x}, \mathbf{y}) = \|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 \quad (6.2)$$

Proof. Using the contents of this thesis that we developed so far, it is straightforward to solve $k_2(\mathbf{x}, \mathbf{y}) = \mathbb{E}[|\mathbf{u}^\top \mathbf{x}|^2 |\mathbf{u}^\top \mathbf{y}|^2]$. In fact, it can be directly derived from the variance formula of complex Gaussian Product-Sketches in Eq. (3.19) in Chapter 3. Recall from Section 3.3 that the variance of a complex Product-Sketch

$$\hat{k}_C(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^p (\mathbf{z}_i^\top \mathbf{x}) (\overline{\mathbf{z}_i^\top \mathbf{y}})$$

with $\mathbf{z}_1, \dots, \mathbf{z}_p \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_d)$ is computed as $\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = \mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2] - (\mathbf{x}^\top \mathbf{y})^{2p}$.

Hence, computing $\mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2] = \mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] + (\mathbf{x}^\top \mathbf{y})^{2p}$ and setting $p = 1$ yields

$$\mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2] = \|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 = k_2(\mathbf{x}, \mathbf{y}),$$

where we used the variance formula for complex Gaussian Product-Sketches in Eq. (3.19). $\mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2]$ corresponds exactly to $k_2(\mathbf{x}, \mathbf{y})$ because we have

$$\mathbb{E}[|\hat{k}_C(\mathbf{x}, \mathbf{y})|^2] = \mathbb{E}[|\mathbf{z}_1^\top \mathbf{x}|^2 |\mathbf{z}_1^\top \mathbf{y}|^2] = \mathbb{E}[|\mathbf{u}^\top \mathbf{x}|^2 |\mathbf{u}^\top \mathbf{y}|^2].$$

■

The result of Theorem 6.2.1 is remarkable because it shows that the optical random features computed by the OPU naturally give rise to a limiting kernel that is similar to a degree-2 polynomial kernel. In fact, when appending the norm of the inputs to this feature map, i.e., by setting $\Phi'(\mathbf{x}) := (\mathbf{i} \|\mathbf{x}\|, \Phi(\mathbf{x})^\top)^\top$ and $\Phi'(\mathbf{y}) := (\mathbf{i} \|\mathbf{y}\|, \Phi(\mathbf{y})^\top)^\top$ with $\mathbf{i} := \sqrt{-1}$, we obtain precisely $\mathbb{E}[\Phi'(\mathbf{x})^\top \Phi'(\mathbf{y})] = (\mathbf{x}^\top \mathbf{y})^2$. We will thus show in Section 6.3 that optical random features are useful for image classification tasks for which polynomial kernels are known to work well.

Higher degree features. In the following, we show that one can obtain higher degree kernels when numerically changing the exponent of the feature map in Eq. (6.1). We focus here on the case $m = 2s$ for some $s \in \mathbb{N}$, i.e., by setting:

$$\Phi_m(\mathbf{x}) := \frac{1}{\sqrt{D}} |\mathbf{U}\mathbf{x}|^m \quad (6.3)$$

In this case, we obtain the kernel $k_{2s}(\mathbf{x}, \mathbf{y})$ described through the following theorem.

Theorem 6.2.2. *When the exponent m in Eq. (6.3) is even, i.e., $m = 2s$, for some $s \in \mathbb{N}$, the dot product of feature maps of Eq. (6.3) tends to the kernel k_{2s} (for $D \rightarrow \infty$) with*

$$k_{2s}(\mathbf{x}, \mathbf{y}) := \|\mathbf{x}\|^m \|\mathbf{y}\|^m \sum_{i=0}^s (s!)^2 \binom{s}{i}^2 \frac{(\mathbf{x}^\top \mathbf{y})^{2i}}{\|\mathbf{x}\|^{2i} \|\mathbf{y}\|^{2i}}. \quad (6.4)$$

The proof in this case is more involved and we move it to [Appendix E.1](#).

Eq. (6.4) is interesting for two reasons. On the one hand, it shows that numerically changing the exponent of the optical random features (6.1) leads to a positively weighted sum of normalized polynomial kernels, similar to the expansion of inhomogeneous polynomial kernels

$$(\nu + \mathbf{x}^\top \mathbf{y})^p = \sum_{i=0}^p \binom{p}{i} \nu^{p-i} (\mathbf{x}^\top \mathbf{y})^i \quad (6.5)$$

with $\nu \geq 0$ and $p \in \mathbb{N}$. The difference between Eq. (6.5) and Eq. (6.4) is that the polynomial kernels in the sum of Eq. (6.4) use inputs normalized to unit length. The entire kernel (6.4) is then rescaled by the norms $\|\mathbf{x}\|^m \|\mathbf{y}\|^m$ after computing this sum.

On the other hand, Eq. (6.4) allows us to compute the variance of the kernel estimate obtained through optical random features. Let $s \in \mathbb{N}$ and assume $D = 1$ for now. Then the following relationship holds:

$$\begin{aligned} \mathbb{V}[\hat{k}_s(\mathbf{x}, \mathbf{y})] &= \mathbb{E}[(\Phi_s(\mathbf{x})\Phi_s(\mathbf{y}))^2] - \mathbb{E}[\Phi_s(\mathbf{x})\Phi_s(\mathbf{y})]^2 \\ &= \mathbb{E}[\Phi_{2s}(\mathbf{x})\Phi_{2s}(\mathbf{y})] - \mathbb{E}[\Phi_s(\mathbf{x})\Phi_s(\mathbf{y})]^2 \\ &= k_{2s}(\mathbf{x}, \mathbf{y}) - k_s(\mathbf{x}, \mathbf{y})^2 \end{aligned}$$

For $D > 1$, we simply divide the variance $\mathbb{V}[\hat{k}_s(\mathbf{x}, \mathbf{y})]$ by D , as the random features are sampled i.i.d. The variance of $\hat{k}_s(\mathbf{x}, \mathbf{y})$ can thus be expressed in terms of $k_{2s}(\mathbf{x}, \mathbf{y})$.

The variance formula for optical random features cannot be directly used to derive a Bernstein-type exponential concentration bound as we did in [Theorem 3.4.1](#) for Rademacher Product-Sketches. This is because both \mathbf{U} and \mathbf{x}, \mathbf{y} are

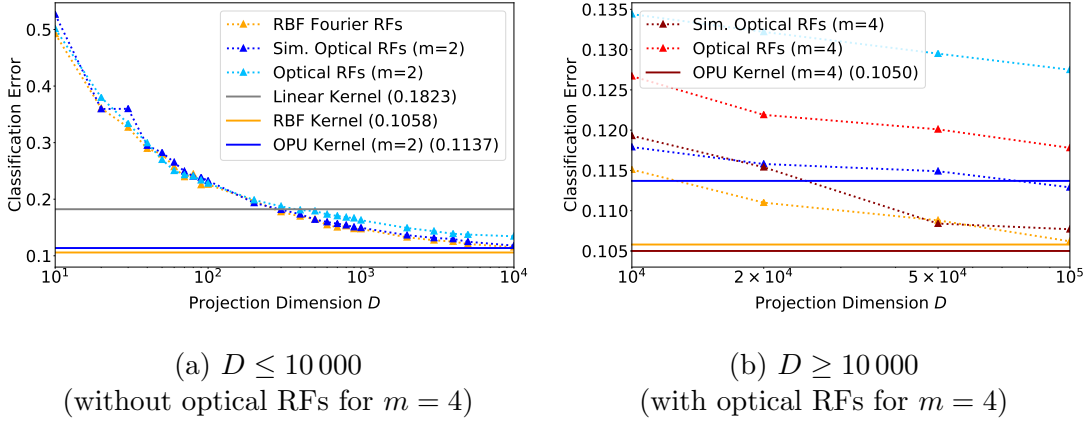


Figure 6.2: Ridge Regression test error on Fashion MNIST for different RFs and projection dimensions D . Horizontal lines show the test error using the limiting kernel. Standard deviations for different seeds are negligibly small and not shown in the plot. Plot (a) compares optical RFs of degree $m = 2$ to RFFs for the Gaussian kernel. Higher degree optical RFs are left out for better readability. The more slowly converging optical RFs for $m = 4$ are added for larger D in plot (b).

unbounded here. Therefore, we used a Bernstein inequality for sub-exponential random variables (Vershynin, 2018, Theorem 2.8.2) in the associated work (Ohana et al., 2020) to show exponentially decaying errors of the kernel estimate under the simplifying assumption that the weights in Eq. (6.1) are real-valued, i.e., $\mathbf{u}_1, \dots, \mathbf{u}_D \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. In this case, we obtain the following theorem:

Theorem 6.2.3. *Let $\epsilon > 0$ and $\hat{k}_m(\mathbf{x}, \mathbf{y}) = \Phi_m(\mathbf{x})^\top \Phi_m(\mathbf{y})$ with Φ_m being defined through Eq. (6.3). Then we can provide the following exponential tail bound for the kernel approximation error:*

$$\Pr \left[\hat{k}_m(\mathbf{x}, \mathbf{y}) - k_m(\mathbf{x}, \mathbf{y}) \geq \epsilon \right] \leq D \exp(-c(D\epsilon)^{1/m}) \quad (6.6)$$

with $c > 0$ being a constant.

The bound given by Theorem 6.2.3 is not tight, but it shows that the kernel approximation error decays exponentially for a sufficiently large number D of random features. Moreover, this decay is slowed down by an increasing exponent m of the feature map. We skip a detailed derivation of this result here and refer the reader to the original work.

6.3 Classification Experiments

In this section, we assess the usefulness of optical random features (6.3) for $m = 2$ and $m = 4$ for different settings and datasets. The model of our choice in

each case is Ridge Regression. OPU experiments were performed remotely on the OPU prototype "Vulcain", running in the LightOn Cloud with library `LightOnOPU` v1.0.2. Since the current version only supports binary input data, we binarize inputs for all experiments using a threshold binarizer, where the threshold is chosen using cross-validation. The code of the experiments is publicly available¹.

6.3.1 Optical Random Features for Fashion MNIST

We compare optical random features (simulated as well as physical) to a random Fourier features (RFF) baseline (Rahimi and Recht, 2007) for different projection dimensions D on Fashion MNIST (Xiao et al., 2017). We use individually optimized hyperparameters for all random feature methods that are found for $D = 10\,000$ using an extensive grid search on a held-out validation set. The same hyperparameters are also used for the precise kernel limit.

Fig. 6.2 shows how the overall classification error decreases as D increases. Part (b) shows that simulated optical RFs for $m = 2$ and RFF (RBF Fourier RFs in the plot) reach the respective kernel test score at $D = 10^5$. Simulated optical RFs for $m = 4$ converge more slowly but outperform $m = 2$ features from $D = 2 \times 10^4$. They perform similarly well as RFF at $D = 10^5$. The performance gap between $m = 2$ and $m = 4$ also increases for the real optical RFs with increasing D . This gap is larger than for the simulated optical RFs due to an increase in regularization for the $m = 2$ features that was needed to add numerical stability when solving linear systems for large D .

The real OPU loses around 1.5% accuracy for $m = 2$ and 1.0% for $m = 4$ for $D = 100\,000$, which is due slightly suboptimal hyperparameters to improve numerical stability for large dimensions. Moreover, there is a small additional loss due to the quantization of the analog signal when the OPU camera records the visual projection.

6.3.2 Transfer Learning on CIFAR-10

An interesting use case for the OPU is transfer learning for image classification. For this purpose, we extract a diverse set of features from the CIFAR-10 (Krizhevsky et al., 2009) image classification dataset using three different convolutional neural networks (ResNet34 (He et al., 2015), AlexNet (Krizhevsky et al., 2009) and VGG16 (Chatfield et al., 2014)). The networks were pretrained on the well-known ImageNet classification benchmark (Russakovsky et al., 2015). For transfer learning, we can either fine-tune these networks and therefore the convolutional features to the data at hand, or we can directly apply a classifier on them assuming that they generalize well enough to the data. The latter case requires much less computational resources while still producing considerable performance gains over the use of the original features. This light-weight approach can be carried out on a

¹<https://github.com/joneswack/opu-kernel-experiments>

Architecture	ResNet34				AlexNet			VGG16				
Layer	L1	L2	L3	Final	MP1	MP2	Final	MP2	MP3	MP4	MP5	Final
Dimension d	4 096	2 048	1 024	512	576	192	9 216	8 192	4 096	2 048	512	25 088
Sim. Opt. RFs	30.4	24.7	28.9	11.6	38.1	41.9	19.6	28.2	20.5	20.7	29.8	15.2 (12.9)
Optical RFs	31.1	25.7	29.7	12.3	39.2	42.6	20.8	30.9	23.3	21.5	30.2	16.4
RBF Four. RFs	30.1	25.2	30.0	12.3	39.4	41.9	19.1	28.0	20.7	20.7	30.1	14.8 (13.0)
No RFs	31.3	26.7	33.5	14.7	44.6	48.8	19.6	27.1	21.0	22.5	34.8	13.3

Table 6.1: Test errors (in %) on CIFAR-10 using $D = 10^4$ RFs for each kernel (except linear). Features were extracted from intermediate layers when using the original input size (32x32). Final convolutional layers were used with upscaled inputs (224x224). L(i) refers to the i th ResNet34 layer and MP(i) to the i th Max-Pool layer of VGG16/AlexNet. Values for the kernel limit are shown in parenthesis (last column).

CPU in a short amount of time where the classification error can be improved with random features.

We compare optical random features and RFF (RBF Four. RFs) to a simple baseline that directly works with the provided convolutional features (no RFs). Table 6.1 shows the test errors achieved on CIFAR-10. Each column corresponds to convolutional features extracted from a specific layer of one of the three networks.

Since the projection dimension $D = 10^4$ was left constant throughout the experiments, it can be observed that all random features perform particularly well compared to a linear kernel when $D \gg d$, where d is the input dimension. For the opposite case $D \ll d$, the lower dimensional projection leads to an increasing test error. This effect can be observed in particular in the last column where the test error of the random feature approximation is higher than without random features. The contrary can be achieved with large enough D as indicated by the values for the true kernel in parenthesis.

A big drawback here is that the computation of sufficiently large dimensional random features may be very costly, especially when d is large as well. This is a regime where the OPU outperforms CPU and GPU by a large margin (see Fig. 6.3) since its computation time is invariant to d and D .

In general, the simulated as well as the physical optical random features yield similar performances as the RBF Fourier RFs on the provided convolutional data.

6.4 Projection Time and Energy Consumption

The main advantage of the OPU compared to a traditional CPU/GPU setup is that the OPU takes a constant time for computing the random features in Eq. (6.1) of arbitrary dimension D (up to $D = 10^6$ on current hardware) for a single input. Moreover, its power consumption stays below 30 W independently of the workload. We therefore compare the computation time of the random feature map in Eq. (6.1)

taken by the OPU against the time it takes to carry out the same computation on a GPU. We simultaneously compare the energy consumption of both devices.

Fig. 6.3 shows the computation time and the energy consumption over time for GPU and OPU for different projection dimensions D . In both cases, the time and energy spending do not include matrix building and loading. For the GPU, only the calls to the PyTorch (Paszke et al., 2019) function `torch.matmul` are measured and energy consumption is the integration over time of power values given by the `nvidia-smi` command.

For the OPU, the energy consumption is constant w.r.t. D and equal to 45 Joules (30 W multiplied by 1.5 seconds). The GPU computation time and energy consumption are monotonically increasing except for an irregular energy development between $D = 45\,000$ and $D = 56\,000$. This exact irregularity was observed throughout all simulations we performed and can most likely be attributed to an optimization routine that the GPU carries out internally. The GPU consumes more than 10 times as much energy as the OPU for $D = 58\,000$ (GPU memory limit). The GPU starts to use more energy than the OPU from $D = 18\,000$. The exact crossover points may change in future hardware versions. The relevant point we make here is that the OPU has a better scalability in D with respect to computation time and energy consumption.

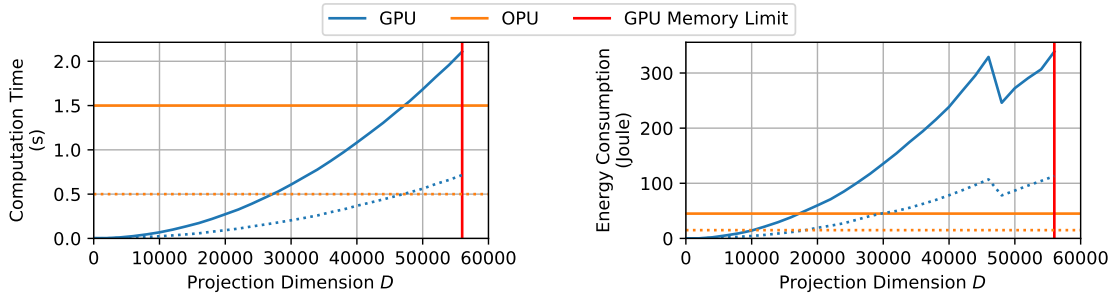


Figure 6.3: Time and energy spent for computing a matrix multiplication $(n, D) \times (D, D)$. The batchsize n is 3000 (solid line) or 1000 (dotted). The curves cross each other at the same D independent from n . We verified more precisely that time and energy are linear with n for both OPU and GPU (experiments were run on an NVIDIA P100).

6.5 Discussion

In this chapter, we have derived the limiting kernel of the OPU and demonstrated its usefulness for practical machine learning applications. We further showed that the feature map can be numerically modified to yield higher order features, and carried out an energy consumption and projection time benchmark, where a modern GPU was outperformed when computing a large number of random features.

Although a large number of random features is needed to give the OPU a practical advantage over the GPU, these results show the *potential* of the OPU to be used for scalable machine learning applications. Future versions of the OPU could accelerate the device such that it becomes beneficial from a lower number of random features D .

Concurrently to the development of the results in this chapter, [Gupta et al. \(2019\)](#) have found a way to obtain an equivalent random projection as the one in [Eq. \(6.1\)](#) using the OPU, but *without* the non-linearity $|\cdot|^2$ that is inherently obtained when recording the light intensity. They thus managed to use the OPU for *linear* random projections. This result makes it possible to use the OPU for random feature approximations and numerical linear algebra methods using Gaussian random projections and thus widens the range of applications immensely. It will be interesting to see in the future how further applications can be accelerated by the OPU and how this can result in energy savings compared to using GPUs.

Chapter 7

Conclusions and Future Work

In this thesis, we carried out an in-depth analysis of random feature approximations for polynomial kernels and more general dot product kernels. In particular, we focused on Product-Sketches originally proposed by [Kar and Karnick \(2012\)](#) and [Hamid et al. \(2014\)](#), i.e., random feature approximations of the form

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p (\mathbf{w}_{i,\ell}^\top \mathbf{x})(\mathbf{w}_{i,\ell}^\top \mathbf{y}), \quad (7.1)$$

where $\hat{k}(\mathbf{x}, \mathbf{y})$ is an unbiased estimate of the polynomial kernel $(\mathbf{x}^\top \mathbf{y})^p$.

Using a detailed variance analysis (Chapter 3), we successively optimized these sketches in terms of the sampling distribution over the weights $\{\mathbf{w}_{i,\ell}\}_{\ell=1}^D$ that we assume to be independent over $i = 1, \dots, p$. In particular, we showed that a variance lower bound can be attained when $\mathbf{w}_{i,\ell}$ are independent Rademacher vectors for $i = 1, \dots, p$ and $\ell = 1, \dots, D$. We further showed that structured polynomial sketches, i.e., sketches where the $\{\mathbf{w}_{i,\ell}\}_{\ell=1}^D$ are *dependent* over $\ell = 1, \dots, D$ can further improve these variances. To the best of our knowledge, our work is the first one to report statistical advantages of structured polynomial sketches.

Moreover, we developed complex-valued generalizations for both structured and non-structured Product-Sketches and elucidated conditions under which such generalizations yield lower variances. As complex-valued random feature maps require the downstream application to handle complex data, we proposed a Complex-to-Real (CtR) transformation for them (Chapter 4) that makes it possible to use these sketches as a drop-in replacement for classical real-valued random features. We studied the impact of this transformation on the approximation variances and showed that CtR sketches keep the variance advantage of the complex sketches introduced before. We also derived a novel CtR-TensorSRHT sketch that outperforms the state-of-the-art in terms of computation speed and kernel approximation errors for low-degree polynomial kernels.

In Chapter 5, we proposed a technique to combine the aforementioned polynomial sketches in order to approximate more general dot product kernels. We

developed a novel feature distribution algorithm that exploits the previously derived variances in order to allocate random features among polynomial sketches of different degrees. We showed that these allocations can improve kernel approximations considerably, which confirms the importance of this subject that has been neglected in related works. At the same time, we uncovered a pathology associated with Maclaurin-based approximations of the Gaussian kernel and proposed a cure for this pathology, leading to state-of-the-art performance for the task of Gaussian kernel approximation across different length scale parameterizations.

Lastly, using the techniques developed in this thesis, we derived the limiting kernel for the random feature approximations obtained by an optical hardware, the Optical Processing Unit. We discussed relations of this kernel to the polynomial kernel and demonstrated its usefulness in image classification tasks.

With this thesis, we have set the foundation for a principled way to improve random feature approximations for dot product kernels. Yet, there is still space for improvements and we summarize some remaining challenges in the following.

1) Exploration of alternatives to complex weights for variance reduction.

We have shown that Product-Sketches (7.1) using complex weights $\mathbf{w}_{i,\ell}$ can achieve considerable variance reductions, in particular for high-degree polynomial kernels. This is because complex random features approximate the linear kernel $(\mathbf{x}^\top \mathbf{y})$ with a lower variance than their real-valued analogs, and this effect is amplified through the product in Eq. (7.1), as we show in Chapter 4.

However, there is no guarantee that our proposed complex-valued estimate of the linear kernel is optimal. We have only proved the optimality of complex Rademacher vectors under the assumption that the elements of $\mathbf{w}_{i,\ell}$ are sampled i.i.d. and that the vectors are independent over $i = 1, \dots, p$ and $\ell = 1, \dots, D$. Future work should therefore explore whether there exist more efficient estimates of the linear kernel that would have an even larger effect on the variance reduction when approximating degree- p polynomial kernels using the estimator in Eq. (7.1).

It would also be interesting to see if an improvement can be achieved by removing the independence assumption over $i = 1, \dots, p$ as it is done for some polynomial sketches in the literature such as TensorSketch (Pham and Pagh, 2013) and TensorSRHT by Ahle et al. (2020). In this case however, it becomes difficult to derive the associated approximation variances.

2) Improving the optimized Maclaurin method.

Although we have shown that our proposed optimized Maclaurin approach optimally allocates the number of random features for a given *stochastic* optimization objective, it would be desirable to make this approach simpler to use in practice. The current algorithm relies on the evaluation of the variance formulas on pairs of data points from a given data set, where the user needs to specify the number of pairs being evaluated.

It would be desirable to achieve a feature allocation without the need for pairwise evaluations, simply by looking at the scale of the data at hand using a one-pass

algorithm with a linear time complexity. We showed in our empirical campaign that the scale of the data is the main driver for the distribution of random features between polynomial degrees. E.g., when the data is modelled using a large length scale, pairwise dot products are distributed closely around the origin and most random features are attributed to low polynomial degrees. The opposite is the case for high-frequency data modelled with a short length scale compared to the scale of the data. There may thus be ways to “shift” the random feature allocation between degrees according to the distribution of the norms of the data (taking the length scale being used into account).

Future work should further investigate the use of alternative basis functions for the approximation of general dot product kernels, e.g., Gegenbauer polynomials as they have lately been proposed by [Han et al. \(2022\)](#). In this case, the ideas developed in this thesis for Maclaurin expansions could be transferred to Gegenbauer expansions, potentially requiring only lower order expansions for the same approximation quality. These may eventually lead to lower approximation variances.

3) Approximating extreme-frequency data. Recently, a random feature approximation of the exponential kernel ([Choromanski et al., 2021](#)) has been proposed for Transformer architectures ([Vaswani et al., 2017](#)). Such a method has the potential of making Transformers scalable to long input sequences and to incur memory savings as explained in Chapter 1. This topic is closely related to the contents of this thesis because the exponential kernel $\exp(\mathbf{x}^\top \mathbf{y}/l^2)$ for $l > 0$ is a dot product kernel.

However, such random feature approximations require a significant amount of fine-tuning of entire Transformer architectures to work well in practice, and they still result in performance degradation compared to using the exact attention mechanism as shown in a recent benchmark ([Tay et al., 2021](#)). This is mainly because the attention patterns that typically arise in Transformer models are extremely sparse. This sparsity in turn stems from high-magnitude inputs being passed to the softmax function that is an integral part of attention layers.

We demonstrate this effect in [Fig. 7.1](#), which shows unnormalized attention patterns, i.e., softmax evaluations without the softmax denominator, of the first layer of a pretrained Bert-based ([Devlin et al., 2019](#)) model when processing an IMDB movie review ([Maas et al., 2011](#)) as an example input. Every pair of rows in [Fig. 7.1](#) corresponds to attention patterns obtained from a single attention head. The first column shows the first 50x50 out of 512x512 values contained in the matrix $\exp(\mathbf{Q}\mathbf{K}^\top)$ (see Chapter 1 for a definition), where the rows of \mathbf{Q} and \mathbf{K} are normalized to unit length in every second row of the figure.

Columns 2-4 show the random feature approximations of $\exp(\mathbf{Q}\mathbf{K}^\top)$, where approximation errors (relative Frobenius norm errors) are shown in parenthesis in the titles. Here we compare a modified random Fourier feature map ([Rahimi and Recht, 2007](#)) using structured projections (see [Yu et al. \(2016\)](#)) against FAVOR+

(Choromanski et al., 2021) and our optimized Maclaurin method (Chapter 5). It becomes clear that the approximation quality of all methods is quite poor, but improves considerably when approximating attention maps with normalized \mathbf{Q} and \mathbf{K} .

This is because the distribution of the entries of $\mathbf{Q}\mathbf{K}^\top$ (right column) has a much larger spread in the unnormalized case producing extreme values for $\exp(\mathbf{Q}\mathbf{K}^\top)$ that ultimately lead to sparse attention patterns in practice. As shown in this thesis and in Choromanski et al. (2021), approximation variances typically scale with the norms of the inputs, \mathbf{Q} and \mathbf{K} in this case, and it is therefore very challenging to approximate the original sparse attention patterns. Although the optimized Maclaurin method achieves very good results when using normalized \mathbf{Q} and \mathbf{K} , it still struggles when approximating the original attention matrix (the same holds for the RFF and FAVOR+ approaches proposed in the literature (Choromanski et al., 2021)).

We therefore conclude that the current state-of-the-art of random feature approximations for dot product kernels is not mature enough to approximate sparse attention patterns. We believe that more efficient approximations would be a significant step towards more efficient modern neural network architectures like Transformers and bilinear neural networks (Lin et al., 2018) (see Chapter 1). This thesis has laid the foundations for further developments in this direction.

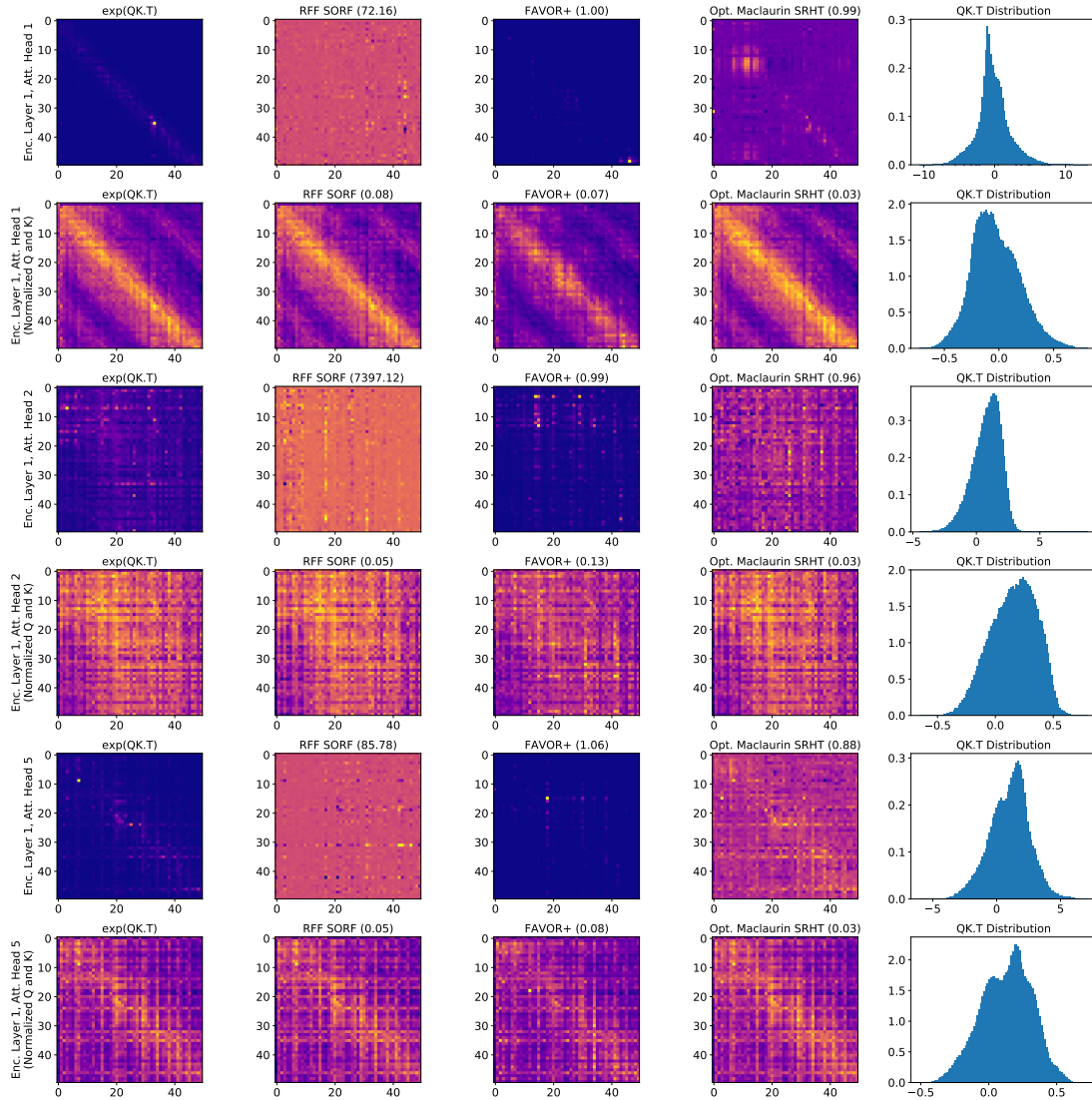


Figure 7.1: Unnormalized attention maps extracted from a single text fragment of length 512 from the IMDB movie reviews (Maas et al., 2011) using a pretrained Bert-base model (Devlin et al., 2019). We show the first 50x50 out of 512x512 pairwise attention evaluations of $\exp(QK^\top)$ along with their approximations through different random feature methods (rel. Frob. error shown in parenthesis in the title). The rightmost column shows a histogram of the entries of QK^\top . We show three different attention heads with the rows of Q and K being normalized to unit-length in every second row.

Appendix A

Appendix of Chapter 2

A.1 Sum of Independent and Uniformly Distributed Variables

In Section 2.2.2, we use the following result. Let $A, B \sim U(\{0, \dots, D-1\})$ be independent. Then $C = (A + B) \bmod D \sim U(\{0, \dots, D-1\})$.

Proof. Without the modulus operation, we have:

$$\begin{aligned} P(A + B = c) &= \sum_{b=0}^{D-1} P(A = c - b | B = b) P(B = b) \\ &= \frac{1}{D} \sum_{b=0}^{D-1} P(A = c - b | B = b) = \frac{1}{D} \sum_{b=0}^{D-1} P(A = c - b) \end{aligned}$$

Now we use the relation $(A + B) \bmod D = c \iff A = (c - B) \bmod D$. So it follows that $((c - b) \bmod D) \in \{0, \dots, D-1\}$ and hence $P(A = (c - b) \bmod D) = 1/D$. Therefore, $C \sim U(\{0, \dots, D-1\})$. ■

Appendix B

Appendix of Chapter 3

B.1 Proofs for Sections 3.1-3.4

B.1.1 Proof of Theorem 3.3.1

We first show

$$\begin{aligned} \mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] &= \left(\sum_{k=1}^d \mathbb{E}[|z_k|^4] x_k^2 y_k^2 + \|\mathbf{x}\|^2 \|\mathbf{y}\|^2 - 2 \sum_{k=1}^d x_k^2 y_k^2 + (\mathbf{x}^\top \mathbf{y})^2 \right. \\ &\quad \left. + \sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d \mathbb{E}[z_i^2] \mathbb{E}[\bar{z}_j^2] x_i x_j y_i y_j \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p}. \end{aligned} \quad (\text{B.1})$$

where $z_i^2 := z_i z_i$ and $\bar{z}_i^2 := \bar{z}_i \bar{z}_i$ are in general different from $|z_i|^2 = z_i \bar{z}_i$. We have

$$\begin{aligned} \mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] &= \mathbb{E}[|\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})|^2] - |\mathbb{E}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})]|^2 = \mathbb{E}\left[\left| \prod_{i=1}^p z_i^\top \mathbf{x} \bar{z}_i^\top \mathbf{y} \right|^2 \right] - (\mathbf{x}^\top \mathbf{y})^{2p} \\ &= \prod_{i=1}^p \mathbb{E}[|z_i^\top \mathbf{x} \bar{z}_i^\top \mathbf{y}|^2] - (\mathbf{x}^\top \mathbf{y})^{2p} = (\mathbb{E}[|\mathbf{z}^\top \mathbf{x} \bar{\mathbf{z}}^\top \mathbf{y}|^2])^p - (\mathbf{x}^\top \mathbf{y})^{2p}. \end{aligned} \quad (\text{B.2})$$

Henceforth we focus on $\mathbb{E}[|\mathbf{z}^\top \mathbf{x} \bar{\mathbf{z}}^\top \mathbf{y}|^2]$ in the last expression (B.2). Write $\mathbf{z} = (z_1, \dots, z_d)^\top$, $\mathbf{x} = (x_1, \dots, x_d)^\top$, and $\mathbf{y} = (y_1, \dots, y_d)^\top$. Since $\mathbb{E}[\mathbf{z} \bar{\mathbf{z}}^\top] = \mathbf{I}_d$, we have $\mathbb{E}[z_i \bar{z}_j] = 1$ if $i = j$ and $\mathbb{E}[z_i \bar{z}_j] = 0$ if $i \neq j$. Recall also that $z_1, \dots, z_d \in \mathbb{C}$ are i.i.d, and $\mathbb{E}[z_i] = 0$ for $i = 1, \dots, d$. Then

$$\begin{aligned} \mathbb{E}[|\mathbf{z}^\top \mathbf{x} \bar{\mathbf{z}}^\top \mathbf{y}|^2] &= \mathbb{E} \left[\left(\sum_{i=1}^d z_i x_i \right) \left(\sum_{j=1}^d \bar{z}_j y_j \right) \left(\sum_{k=1}^d \bar{z}_k x_k \right) \left(\sum_{l=1}^d z_l y_l \right) \right] \\ &= \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d \sum_{l=1}^d \mathbb{E}[z_i \bar{z}_j \bar{z}_k z_l] x_i y_j x_k y_l. \end{aligned} \quad (\text{B.3})$$

The expected value $\mathbb{E}[z_i \bar{z}_j \bar{z}_k z_l]$ is different from 0, only if:

- (a) $i = j = k = l$, for which there are d terms and $\mathbb{E}[z_i \overline{z_j z_k z_l}] x_i y_j x_k y_l = \mathbb{E}[|z_i|^4] x_i^2 y_i^2$.
- (b) $i = j \neq k = l$, for which there are $d(d-1)$ terms and $\mathbb{E}[z_i \overline{z_j z_k z_l}] x_i y_j x_k y_l = \mathbb{E}[|z_i|^2] \mathbb{E}[|z_k|^2] x_i x_k y_i y_k$.
- (c) $i = k \neq j = l$, for which there are $d(d-1)$ terms and $\mathbb{E}[z_i \overline{z_j z_k z_l}] x_i y_j x_k y_l = \mathbb{E}[|z_i|^2] \mathbb{E}[|z_j|^2] x_i^2 y_j^2$.
- (d) $i = l \neq j = k$, for which there are $d(d-1)$ terms and $\mathbb{E}[z_i \overline{z_j z_k z_l}] x_i y_j x_k y_l = \mathbb{E}[z_i^2] \mathbb{E}[\overline{z_j^2}] x_i x_j y_i y_j$.

Therefore,

$$\begin{aligned}
(\text{B.3}) &= \underbrace{\sum_{i=1}^d \mathbb{E}[|z_i|^4] x_i^2 y_i^2}_{\text{case (a)}} + \underbrace{\sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d \mathbb{E}[|z_i|^2] \mathbb{E}[|z_j|^2] x_i^2 y_j^2}_{\text{case (c)}} + \underbrace{\sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d \mathbb{E}[|z_i|^2] \mathbb{E}[|z_j|^2] x_i x_j y_i y_j}_{\text{case (b)}} \\
&+ \underbrace{\sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d \mathbb{E}[z_i^2] \mathbb{E}[\overline{z_j^2}] x_i x_j y_i y_j}_{\text{case (d)}} \\
&= \sum_{i=1}^d \mathbb{E}[|z_i|^4] x_i^2 y_i^2 + \sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d x_i^2 y_j^2 + \sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d x_i x_j y_i y_j + \sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d \mathbb{E}[z_i^2] \mathbb{E}[\overline{z_j^2}] x_i x_j y_i y_j \\
&= \sum_{i=1}^d \mathbb{E}[|z_i|^4] x_i^2 y_i^2 + \left[\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 - \sum_{i=1}^d x_i^2 y_i^2 \right] + \left[(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right] \\
&+ \sum_{i=1}^d \sum_{\substack{j=1 \\ j \neq i}}^d \mathbb{E}[z_i^2] \mathbb{E}[\overline{z_j^2}] x_i x_j y_i y_j
\end{aligned}$$

The proof of Eq. (B.1) completes by using this expression of $\mathbb{E}[(\mathbf{z}^\top \overline{\mathbf{x} \mathbf{z}^\top \mathbf{y}})^2]$ in Eq. (B.2).

Eq. (3.16) follows from Eq. (B.1) and $\mathbb{E}[z_k^2] = \mathbb{E}[\overline{z_k^2}] = 2q - 1$, which uses Eq. (3.15). The lower bound Eq. (3.17) follows from Jensen's inequality $\mathbb{E}[|z_k|^4] \geq (\mathbb{E}[|z_k|^2])^2 = 1$.

B.1.2 Proof of Theorem 3.4.1

We make use of Bernstein's inequality (e.g., Vershynin, 2018, Theorem 2.8.4): For independent random variables $X_1, \dots, X_D \in \mathbb{R}$ such that $\mathbb{E}[X_i] = 0$ and $|X_i| \leq R$

almost surely for a constant $R > 0$, we have for any $t > 0$:

$$\Pr \left[\left| \sum_{i=1}^D X_i \right| \geq t \right] \leq 2 \exp \left(\frac{-t^2/2}{\sum_{i=1}^D \mathbb{V}[X_i] + Rt/3} \right) \quad (\text{B.4})$$

We define $X_i := \Phi_{\mathcal{C}}(\mathbf{x})_i \overline{\Phi_{\mathcal{C}}(\mathbf{y})_i} - (\mathbf{x}^\top \mathbf{y})^p / D \in \mathbb{R}$, where $\Phi_{\mathcal{C}}(\mathbf{x}) \in \mathbb{C}^D$ is defined in Eq. (3.10): $\Phi_{\mathcal{C}}(\mathbf{x}) = \frac{1}{\sqrt{D}} [(\prod_{i=1}^p \mathbf{z}_{i,1}^\top \mathbf{x}), \dots, (\prod_{i=1}^p \mathbf{z}_{i,D}^\top \mathbf{x})]^\top$. Then we have $\mathbb{E}[X_i] = 0$. Moreover,

$$\begin{aligned} |X_i| &\leq |\Phi_{\mathcal{C}}(\mathbf{x})_i \overline{\Phi_{\mathcal{C}}(\mathbf{y})_i}| + |(\mathbf{x}^\top \mathbf{y})^p / D| = \frac{1}{D} \left(\prod_{j=1}^p |\mathbf{z}_j^\top \mathbf{x}| |\overline{\mathbf{z}_j^\top \mathbf{y}}| + |(\mathbf{x}^\top \mathbf{y})^p| \right) \\ &\leq \frac{1}{D} (\|\mathbf{x}\|_1^p \|\mathbf{y}\|_1^p + \|\mathbf{x}\|_2^p \|\mathbf{y}\|_2^p) \leq \frac{2}{D} \|\mathbf{x}\|_1^p \|\mathbf{y}\|_1^p =: R \end{aligned}$$

where the first inequality is the triangle inequality. The second inequality uses Hölder's inequality (and that the absolute value of each element of \mathbf{z}_j is 1) as well as the upper bound $\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$. Furthermore, by assumption we have

$$\mathbb{V}[X_i] = \frac{\sigma^2 \|\mathbf{x}\|_2^{2p} \|\mathbf{y}\|_2^{2p}}{D^2} \leq \frac{\sigma^2 \|\mathbf{x}\|_1^{2p} \|\mathbf{y}\|_1^{2p}}{D^2}$$

for some $\sigma^2 \geq 0$. Therefore, using Eq. (B.4) and setting $t := \|\mathbf{x}\|_1^p \|\mathbf{y}\|_1^p \epsilon$, we have

$$\Pr \left[\left| \sum_{i=1}^D X_i \right| \geq \epsilon \|\mathbf{x}\|_1^p \|\mathbf{y}\|_1^p \right] \leq 2 \exp \left(\frac{-D\epsilon^2/2}{\frac{2}{3}\epsilon + \sigma^2} \right)$$

Setting $D \geq 2(\frac{2}{3\epsilon} + \frac{\sigma^2}{\epsilon^2}) \log(\frac{2}{\delta})$ and taking the complementary probability gives the desired result.

B.2 Proofs for Section 3.5

B.2.1 Key Lemma

First, we state a key lemma that is needed for deriving the variance of real and complex TensorSRHT. This result is essentially given in Choromanski et al. (2017, Proof of Proposition 8.2). However, their proof contains a typo missing the negative sign, and they use a different definition of the Hadamard matrix from ours. Therefore, for completeness, we state the result formally and provide a proof. We also provide an extension to stacked Hadamard matrices that is used for our upsampled TensorSRHT algorithm.

Lemma B.2.1. *Let $d = 2^m$ for some $m \in \mathbb{N}$ and $\mathbf{H}_d = (\mathbf{h}_1, \dots, \mathbf{h}_d) \in \{1, -1\}^{d \times d}$ be the unnormalized Hadamard matrix defined in Eq. (2.21), where $\mathbf{h}_\ell = (h_{\ell,1}, \dots, h_{\ell,d})^\top \in$*

$\{1, -1\}^d$ for $\ell \in \{1, \dots, d\}$. Let $\pi : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$ be a uniformly random permutation. Then for any $\ell, \ell' \in \{1, \dots, d\}$ with $\ell \neq \ell'$ and $t, u \in \{1, \dots, d\}$ with $t \neq u$, we have

$$\mathbb{E}[h_{\pi(\ell),t}h_{\pi(\ell'),t}h_{\pi(\ell),u}h_{\pi(\ell'),u}] = -\frac{1}{d-1},$$

where the expectation is with respect to the random permutation π .

Let $B := \lceil D/d \rceil$. We can then extend this result to a concatenated matrix $\mathbf{H}_d^B := \underbrace{(\mathbf{H}_d, \dots, \mathbf{H}_d)}_{B \text{ times}} \in \{1, -1\}^{d \times Bd}$ whose columns we name as $\mathbf{h}_1^B, \dots, \mathbf{h}_{Bd}^B$. In

this case, we use the uniform permutation $\pi^B : \{1, \dots, Bd\} \rightarrow \{1, \dots, Bd\}$. Then for any $\ell, \ell' \in \{1, \dots, Bd\}$ with $\ell \neq \ell'$ and $t, u \in \{1, \dots, d\}$ with $t \neq u$, we have

$$\mathbb{E}[h_{\pi(\ell),t}h_{\pi(\ell'),t}h_{\pi(\ell),u}h_{\pi(\ell'),u}] = -\frac{1}{\lceil D/d \rceil d - 1}.$$

Proof. We first derive a few key identities needed for our proof and we focus on using \mathbf{H}_d instead of \mathbf{H}_d^B for simplicity for now. For ease of notation, define

$$\alpha_\ell := h_{\ell,t}h_{\ell,u}, \quad \ell \in \{1, \dots, d\}.$$

Since any two distinct rows (and any two distinct columns) of \mathbf{H}_d are orthogonal, we have

$$\sum_{\ell=1}^d \alpha_\ell = \sum_{\ell=1}^d h_{\ell,t}h_{\ell,u} = 0.$$

Since $\alpha_\ell \in \{-1, 1\}$, this identity implies that exactly $d/2$ elements in $\{\alpha_1, \dots, \alpha_d\}$ are 1, and the rest are -1 . Note that for each $\ell \in \{1, \dots, d\}$ the randomly permuted index $\pi(\ell)$ takes values in $\{1, \dots, d\}$ with equal probabilities. Therefore, the probability of $\alpha_{\pi(\ell)}$ being 1 and that of $\alpha_{\pi(\ell)}$ being -1 are equal:

$$\Pr(\alpha_{\pi(\ell)} = 1) = \Pr(\alpha_{\pi(\ell)} = -1) = 0.5.$$

Note that $\pi^b(\ell) \neq \pi^b(\ell')$ since $\ell \neq \ell'$ and π is a (random) permutation. Therefore, we have the following conditional probabilities:

$$\Pr(\alpha_{\pi(\ell')} = a \mid \alpha_{\pi(\ell)} = b) = \begin{cases} \frac{d/2-1}{d-1} & \text{if } a = b = 1 \text{ or } a = b = -1 \\ \frac{d/2}{d-1} & \text{if } a = 1, b = -1 \text{ or } a = -1, b = -1 \end{cases}$$

Using the above identities, we now prove the assertion:

$$\begin{aligned} \mathbb{E}[h_{\pi(\ell),t}h_{\pi(\ell'),t}h_{\pi(\ell),u}h_{\pi(\ell'),u}] &= \mathbb{E}[\alpha_{\pi(\ell)}\alpha_{\pi(\ell')}] \\ &= \Pr(\alpha_{\pi(\ell)} = 1)\mathbb{E}[\alpha_{\pi(\ell)}\alpha_{\pi(\ell')} \mid \alpha_{\pi(\ell)} = 1] + \Pr(\alpha_{\pi(\ell)} = -1)\mathbb{E}[\alpha_{\pi(\ell)}\alpha_{\pi(\ell')} \mid \alpha_{\pi(\ell)} = -1] \\ &= \frac{1}{2}\mathbb{E}[\alpha_{\pi(\ell')} \mid \alpha_{\pi(\ell)} = 1] - \frac{1}{2}\mathbb{E}[\alpha_{\pi(\ell')} \mid \alpha_{\pi(\ell)} = -1] \\ &= \frac{1}{2} \left(\frac{d/2-1}{d-1} - \frac{d/2}{d-1} \right) - \frac{1}{2} \left(\frac{d/2}{d-1} - \frac{d/2-1}{d-1} \right) = -\frac{1}{d-1}. \end{aligned}$$

Extending this analysis to the case of using \mathbf{H}_d^B and π^B instead of \mathbf{H}_d and π is straightforward. In this case, the marginal probabilities $\Pr(\alpha_{\pi^B(\ell)} = 1)$ and $\Pr(\alpha_{\pi^B(\ell)} = -1)$ remain the same because for two fixed rows t, u of \mathbf{H}_d^B with $t \neq u$, we still have

$$\sum_{\ell=1}^{Bd} \alpha_\ell = \sum_{\ell=1}^d h_{\ell,t}^B h_{\ell,u}^B = 0$$

and thus $Bd/2$ of the elements in $\{\alpha_1, \dots, \alpha_d\}$ must be 1 and -1 , respectively. The conditional probabilities however change, and become

$$\Pr(\alpha_{\pi^B(\ell')} = a \mid \alpha_{\pi(\ell)} = b) = \begin{cases} \frac{(\lceil D/d \rceil d)/2 - 1}{\lceil D/d \rceil d - 1} & \text{if } a = b = 1 \text{ or } a = b = -1 \\ \frac{(\lceil D/d \rceil d)/2}{\lceil D/d \rceil d - 1} & \text{if } a = 1, b = -1 \text{ or } a = -1, b = -1 \end{cases}$$

This leads to

$$\mathbb{E}[h_{\pi^B(\ell'),t}^B h_{\pi^B(\ell'),t}^B h_{\pi^B(\ell'),u}^B h_{\pi^B(\ell'),u}^B] = \mathbb{E}[\alpha_{\pi^B(\ell')} \alpha_{\pi^B(\ell')}] = -\frac{1}{\lceil D/d \rceil d - 1}$$

in this case. ■

B.2.2 Proof of Theorem 3.5.9

We first clarify the notation we use. Recall that our feature map $\Phi(\mathbf{x}) \in \mathbb{C}^D$ is given by

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{D}} \left[\left(\prod_{i=1}^p \mathbf{s}_{i,1}^\top \mathbf{x} \right), \dots, \left(\prod_{i=1}^p \mathbf{s}_{i,D}^\top \mathbf{x} \right) \right]^\top \in \mathbb{C}^D.$$

The random vectors $\mathbf{s}_{i,\ell} \in \mathbb{C}^d$ are independently generated blockwise, and there are $B := \lceil D/d \rceil$ blocks in total (and note that $D = (B-1)d + \text{mod}(D, d)$): For each $i = 1, \dots, p$,

$$\begin{aligned} & \underbrace{(\mathbf{s}_{i,1}, \dots, \mathbf{s}_{i,d})}_{\text{Block 1}}, \underbrace{(\mathbf{s}_{i,d+1}, \dots, \mathbf{s}_{i,2d})}_{\text{Block 2}}, \dots, \\ & \underbrace{(\mathbf{s}_{i,(B-2)d+1}, \dots, \mathbf{s}_{i,(B-1)d})}_{\text{Block } B-1}, \underbrace{(\mathbf{s}_{i,(B-1)d+1}, \dots, \mathbf{s}_{i,(B-1)d+\text{mod}(D,d)})}_{\text{Block } B} \\ & =: \underbrace{(\mathbf{s}_{i,1}^1, \dots, \mathbf{s}_{i,d}^1)}_{\text{Block 1}}, \underbrace{(\mathbf{s}_{i,1}^2, \dots, \mathbf{s}_{i,d}^2)}_{\text{Block 2}}, \dots, \underbrace{(\mathbf{s}_{i,1}^{B-1}, \dots, \mathbf{s}_{i,d}^{B-1})}_{\text{Block } B-1}, \underbrace{(\mathbf{s}_{i,1}^B, \dots, \mathbf{s}_{i,\text{mod}(D,d)}^B)}_{\text{Block } B}, \end{aligned}$$

where we introduced in the second line a new notation:

$$\mathbf{s}_{i,\ell}^b := \mathbf{s}_{i,(b-1)d+\ell} \quad (b = 1, \dots, B, \ell = 1, \dots, d).$$

Here b serves as the indicator of the b -th block. Thus, using this notation,

$$\mathbf{s}_{i,\ell}^b = \mathbf{z}_i^b \circ \mathbf{h}_{\pi^b(\ell)} \in \mathbb{C}^d \quad (\ell = 1, \dots, d),$$

where $\mathbf{z}_i^b = (z_{i,1}^b, \dots, z_{i,d}^b)^\top \in \mathbb{C}^d$ is a random vector whose elements $z_{i,1}^b, \dots, z_{i,d}^b$ are i.i.d., and $\pi^b: \{1, \dots, d\} \rightarrow \{1, \dots, d\}$ is a random permutation of the indices. Note that \mathbf{z}_i^b and π^b are generated independently for each $b \in \{1, \dots, B\}$. Therefore, the random vectors $\mathbf{s}_{i,\ell}^b$ and $\mathbf{s}_{i,\ell'}^{b'}$ are statistically independent if they are from different blocks, i.e., if $b \neq b'$.

For each $b = 1, \dots, B$, define $\mathbf{z}^b = (z_1^b, \dots, z_d^b)^\top \in \mathbb{C}^d$ as a random vector independently and identically distributed as $\mathbf{z}_1^b, \dots, \mathbf{z}_p^b$. Define

$$\mathbf{s}_\ell^b := \mathbf{z}^b \circ \mathbf{h}_{\pi^b(\ell)} = (z_1^b h_{\pi^b(\ell),1}, \dots, z_d^b h_{\pi^b(\ell),d})^\top =: (s_{\ell,1}^b, \dots, s_{\ell,d}^b)^\top \in \mathbb{C}^d. \quad (\text{B.5})$$

Then \mathbf{s}_ℓ^b is independently and identically distributed as $\mathbf{s}_{1,\ell}^b, \dots, \mathbf{s}_{p,\ell}^b$. Moreover, given the permutation π^b fixed, \mathbf{s}_ℓ^b is identically distributed as \mathbf{z}^b . This is because 1) z_1^b, \dots, z_d^b are i.i.d., 2) each z_t^b is symmetrically distributed ($t = 1, \dots, d$), and 3) $h_{\pi^b(\ell),1}, \dots, h_{\pi^b(\ell),d} \in \{1, -1\}$.

Now let us start proving the assertion. We first have

$$\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})] = \mathbb{E}[|\hat{k}(\mathbf{x}, \mathbf{y})|^2] - |\mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})]|^2 = \mathbb{E}[|\hat{k}(\mathbf{x}, \mathbf{y})|^2] - (\mathbf{x}^\top \mathbf{y})^{2p},$$

where the second identity follows from the approximate kernel being unbiased for both real and complex TensorSRHT. Thus, from now on we study the term $\mathbb{E}[|\hat{k}(\mathbf{x}, \mathbf{y})|^2]$.

For simplicity of notation, define $I_b := \{1, \dots, d\}$ for $b = 1, \dots, B-1$ and $I_b := \{1, \dots, \text{mod}(D, d)\}$ for $b = B$. Since the approximate kernel can be written as

$$\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi(\mathbf{x})^\top \overline{\Phi(\mathbf{y})} = \frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p \mathbf{s}_{i,\ell}^\top \mathbf{x} \overline{\mathbf{s}_{i,\ell}^\top \mathbf{y}} = \frac{1}{D} \sum_{b=1}^B \sum_{\ell \in I_b} \prod_{i=1}^p (\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) \overline{(\mathbf{s}_{i,\ell}^{b\top} \mathbf{y})}$$

its second moment can be written as

$$\begin{aligned} \mathbb{E}[|\hat{k}(\mathbf{x}, \mathbf{y})|^2] &= \frac{1}{D^2} \sum_{b,b'=1}^B \sum_{\ell \in I_b} \sum_{\ell' \in I_{b'}} \mathbb{E} \left[\prod_{i=1}^p (\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) \overline{(\mathbf{s}_{i,\ell}^{b\top} \mathbf{y})} (\mathbf{s}_{i,\ell'}^{b'\top} \mathbf{x}) \overline{(\mathbf{s}_{i,\ell'}^{b'\top} \mathbf{y})} \right] \\ &= \frac{1}{D^2} \sum_{b,b'=1}^B \sum_{\ell \in I_b} \sum_{\ell' \in I_{b'}} \prod_{i=1}^p \mathbb{E} \left[(\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) \overline{(\mathbf{s}_{i,\ell}^{b\top} \mathbf{y})} (\mathbf{s}_{i,\ell'}^{b'\top} \mathbf{x}) \overline{(\mathbf{s}_{i,\ell'}^{b'\top} \mathbf{y})} \right] \\ &= \frac{1}{D^2} \sum_{b,b'=1}^B \sum_{\ell \in I_b} \sum_{\ell' \in I_{b'}} \left(\mathbb{E} \left[(\mathbf{s}_\ell^{b\top} \mathbf{x}) \overline{(\mathbf{s}_\ell^{b\top} \mathbf{y})} (\mathbf{s}_{\ell'}^{b'\top} \mathbf{x}) \overline{(\mathbf{s}_{\ell'}^{b'\top} \mathbf{y})} \right] \right)^p. \quad (\text{B.6}) \end{aligned}$$

Now we study individual terms in (B.6), categorizing the indices $b, b' \in \{1, \dots, B\}$ and $\ell, \ell' \in \{1, \dots, d\}$ of indices into the following 3 cases:

1. $b = b'$ and $\ell = \ell'$ (D terms): As mentioned earlier, conditioned on the permutation π^b , \mathbf{s}_ℓ^b is identically distributed as \mathbf{z}^b (see the paragraph following Eq. (B.5)). Thus,

$$\begin{aligned} \mathbb{E} \left[(\mathbf{s}_\ell^{b\top} \mathbf{x})^2 (\mathbf{s}_\ell^{b\top} \mathbf{y})^2 \right] &= \mathbb{E}_{\pi^b} \left[\mathbb{E} \left[(\mathbf{s}_\ell^{b\top} \mathbf{x})^2 (\mathbf{s}_\ell^{b\top} \mathbf{y})^2 \mid \pi^b \right] \right] \\ &= \mathbb{E}_{\pi^b} \left[\mathbb{E} \left[(\mathbf{z}^{b\top} \mathbf{x})^2 (\mathbf{z}^{b\top} \mathbf{y})^2 \right] \right] = \mathbb{E} \left[(\mathbf{z}^{b\top} \mathbf{x})^2 (\mathbf{z}^{b\top} \mathbf{y})^2 \right] = \mathbb{E} \left[(\mathbf{z}^\top \mathbf{x})^2 (\mathbf{z}^\top \mathbf{y})^2 \right], \end{aligned}$$

where \mathbb{E}_{π^b} denotes the expectation with respect to π^b and $\mathbf{z} \in \mathbb{C}^d$ is a random vector identically distributed as $\mathbf{z}^1, \dots, \mathbf{z}^B$.

2. $b = b'$ and $\ell \neq \ell'$ ($c(D, d)$ terms, where $c(D, d)$ is defined in Eq. (3.29)): This case requires a detailed analysis, which we will do below.
3. $b \neq b'$ (The rest of terms = $D^2 - D - c(D, d)$ terms): Since \mathbf{s}_ℓ^b and $\mathbf{s}_{\ell'}^{b'}$ are independent in this case, we have

$$\begin{aligned} \mathbb{E} \left[(\mathbf{s}_\ell^{b\top} \mathbf{x}) \overline{(\mathbf{s}_\ell^{b\top} \mathbf{y})} (\mathbf{s}_{\ell'}^{b'\top} \mathbf{x}) (\mathbf{s}_{\ell'}^{b'\top} \mathbf{y}) \right] &= \mathbb{E} \left[(\mathbf{s}_\ell^{b\top} \mathbf{x}) \overline{(\mathbf{s}_\ell^{b\top} \mathbf{y})} \right] \mathbb{E} \left[(\mathbf{s}_{\ell'}^{b'\top} \mathbf{x}) (\mathbf{s}_{\ell'}^{b'\top} \mathbf{y}) \right] \\ &= \mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})] \mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})] = (\mathbf{x}^\top \mathbf{y})^2, \end{aligned}$$

where the last equality follows from the approximate kernel being unbiased.

We now analyze the case 2:

$$\begin{aligned} \mathbb{E} \left[(\mathbf{s}_\ell^{b\top} \mathbf{x}) \overline{(\mathbf{s}_\ell^{b\top} \mathbf{y})} (\mathbf{s}_{\ell'}^{b'\top} \mathbf{x}) (\mathbf{s}_{\ell'}^{b'\top} \mathbf{y}) \right] &= \sum_{t,u,w,v=1}^d \mathbb{E}[s_{\ell,t}^b \overline{s_{\ell,u}^b} s_{\ell',v}^b s_{\ell',w}^b] x_t y_u x_v y_w \\ &= \sum_{t,u,w,v=1}^d \mathbb{E}[z_t^b \overline{z_u^b} z_v^b z_w^b] \underbrace{\mathbb{E}[h_{\pi^b(\ell),t} h_{\pi^b(\ell),u} h_{\pi^b(\ell),v} h_{\pi^b(\ell),w}]}_{=:E} x_t y_u x_v y_w \end{aligned}$$

Note that we have $\mathbb{E}[z_t^b \overline{z_u^b} z_v^b z_w^b] = 0$ unless:

- (a) $t = u = v = w$: $\mathbb{E}[z_t^b \overline{z_u^b} z_v^b z_w^b] = \mathbb{E}[|z_t^b|^4] = 1$ and $E = \mathbb{E}[h_{\pi^b(\ell),t}^2 h_{\pi^b(\ell),t}^2] = 1$.
- (b) $t = u \neq v = w$: $\mathbb{E}[z_t^b \overline{z_u^b} z_v^b z_w^b] = \mathbb{E}[|z_t^b|^2 |z_v^b|^2] = 1$ and $E = \mathbb{E}[h_{\pi^b(\ell),t}^2 h_{\pi^b(\ell),v}^2] = 1$.
- (c) $t = v \neq u = w$: $\mathbb{E}[z_t^b \overline{z_u^b} z_v^b z_w^b] = \mathbb{E}[|z_t^b|^2 |z_u^b|^2] = 1$ and $E = \mathbb{E}[h_{\pi^b(\ell),t} h_{\pi^b(\ell),u} h_{\pi^b(\ell),t} h_{\pi^b(\ell),u}]$.
- (d) $t = w \neq u = v$: $\mathbb{E}[z_t^b \overline{z_u^b} z_v^b z_w^b] = \mathbb{E}[(z_t^b)^2 (\overline{z_u^b})^2] = (2q - 1)^2$ and $E = \mathbb{E}[h_{\pi^b(\ell),t} h_{\pi^b(\ell),u} h_{\pi^b(\ell),u} h_{\pi^b(\ell),t}]$.

Therefore, we have

$$\begin{aligned}
& \mathbb{E} [\mathbf{s}_\ell^\top \mathbf{x} \mathbf{s}_\ell^\top \mathbf{y} \mathbf{s}_{\ell'}^\top \mathbf{x} \mathbf{s}_{\ell'}^\top \mathbf{y}] \\
&= \sum_{t=1}^d x_t^2 y_t^2 + \sum_{t \neq v} x_t y_t x_v y_v + \sum_{t \neq u} \mathbb{E}[h_{\pi^b(\ell),t} h_{\pi^b(\ell'),t} h_{\pi^b(\ell),u} h_{\pi^b(\ell'),u}] (x_t^2 y_u^2 + (2q-1)x_t y_t x_u y_u) \\
&= (\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{d-1} \sum_{t \neq u} (x_t^2 y_u^2 + (2q-1)x_t y_t x_u y_u) \quad (\because \text{Lemma B.2.1}) \\
&= (\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1},
\end{aligned}$$

where $V_q^{(1)} := \sum_{t \neq u} (x_t^2 y_u^2 + (2q-1)x_t y_t x_u y_u)$ is Eq. (3.17) with $p = 1$, which is the variance of the unstructured polynomial sketch (3.10) with a single feature.

Now, using these identities in Eq. (B.6), the variance of the approximate kernel can be expanded as

$$\begin{aligned}
\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})] &= \mathbb{E}[\hat{k}(\mathbf{x}, \mathbf{y})^2] - (\mathbf{x}^\top \mathbf{y})^{2p} \\
&= \frac{1}{D} \left(\mathbb{E} \left[(\mathbf{z}^\top \mathbf{x})^2 (\mathbf{z}^\top \mathbf{y})^2 \right] \right)^p + \frac{c(D, d)}{D^2} \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^p \\
&\quad + \frac{D^2 - D - c(D, d)}{D^2} (\mathbf{x}^\top \mathbf{y})^{2p} - (\mathbf{x}^\top \mathbf{y})^{2p} \\
&= \frac{1}{D} \left[\left(\mathbb{E} \left[(\mathbf{z}^\top \mathbf{x})^2 (\mathbf{z}^\top \mathbf{y})^2 \right] \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right] - \frac{c(D, d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^p \right] \\
&= \frac{1}{D} V_q^{(p)} - \frac{c(D, d)}{D^2} \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^p \right]
\end{aligned}$$

where $V_q^{(p)} \geq 0$ is Eq. (3.17) with the considered value of the polynomial degree p , which is the variance of the unstructured polynomial sketch (3.10) with a single feature. This completes the proof.

Extending the proof to upsampled TensorSRHT

For the case of upsampled TensorSRHT, we use a single block of size $\lceil D/d \rceil d$. Therefore, we only need to distinguish the first two cases in the analysis of Eq. (B.6), 1. $\ell = \ell'$ (D terms) and 2. $\ell \neq \ell'$, for which there are now $D^2 - D$ terms.

Case 1. remains identical to stacked TensorSRHT. Case 2. is almost identical, the difference being that the number of non-zero terms for which $\ell \neq \ell'$ is now $D^2 - D$ instead of $c(D, d)$. Moreover, we do not permute the columns of a single matrix \mathbf{H}_d , but the ones of \mathbf{H}_d^B , which is why we need to make use of the second part of Lemma B.2.1. We thus substitute $\mathbb{E}[h_{\pi^b(\ell),t} h_{\pi^b(\ell'),t} h_{\pi^b(\ell),u} h_{\pi^b(\ell'),u}]$ by $-\frac{1}{\lceil D/d \rceil d - 1}$ instead of $-\frac{1}{d-1}$.

This finally leads to

$$\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})] = \frac{1}{D} V_q^{(p)} - \left(1 - \frac{1}{D}\right) \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{\lceil D/d \rceil d - 1} \right)^p \right].$$

B.3 Gaussian Processes with Complex Random Features

We describe here how to use complex random features in Gaussian process (GP) regression and classification. Since real random features are special cases of complex random features, all derivations for the complex case also hold for the real case as well.

For GP classification, we employ the framework of [Milios et al. \(2018\)](#), which formulates GP classification using GP regression and provides a solution in closed form. Therefore, closed form solutions are available for both GP regression and classification, and this enables us to compare different random feature approximations directly.¹

Notation and definitions. For a matrix $\mathbf{A} \in \mathbb{C}^{n \times m}$ with $n, m \in \mathbb{N}$, denote by $\mathbf{A}^H := \overline{\mathbf{A}}^\top \in \mathbb{C}^{m \times n}$ be its conjugate transpose. Note that if $\mathbf{A} \in \mathbb{R}^{n \times m}$, then $\mathbf{A}^H = \mathbf{A}^\top \in \mathbb{R}^{m \times n}$. For $n \in \mathbb{N}$, $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ be the identity matrix.

For $\boldsymbol{\mu} \in \mathbb{C}^n$ and positive semi-definite² $\boldsymbol{\Sigma} \in \mathbb{C}^{n \times n}$ with $n \in \mathbb{N}$, we denote by $\mathcal{CN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ the n -dimensional *proper* complex Gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, whose density function is given by (e.g., [Neeser and Massey, 1993](#), Theorem 1)

$$\mathcal{CN}(\mathbf{v}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) := \frac{1}{\pi^n \sqrt{|\boldsymbol{\Sigma}|}} \exp\left(-(\mathbf{v} - \boldsymbol{\mu})^H \boldsymbol{\Sigma}^{-1} (\mathbf{v} - \boldsymbol{\mu})\right), \quad \mathbf{v} \in \mathbb{C}^n,$$

where $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$. If a random vector $\mathbf{f} \in \mathbb{C}^n$ follows $\mathcal{CN}(\mathbf{v}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$, we have $\mathbb{E}[\mathbf{f}] = \boldsymbol{\mu}$, $\mathbb{E}[(\mathbf{f} - \boldsymbol{\mu})(\mathbf{f} - \boldsymbol{\mu})^H] = \boldsymbol{\Sigma}$, and $\mathbb{E}[(\mathbf{f} - \boldsymbol{\mu})(\mathbf{f} - \boldsymbol{\mu})^\top] = 0$, where the last property is the definition of \mathbf{f} being a proper complex random variable ([Neeser and Massey, 1993](#), Definition 1).

B.3.1 Complex GP Regression

We first describe the approach of *complex GP regression* ([Boloix-Tortosa et al., 2018](#)), a Bayesian nonparametric approach to complex-valued regression.

¹If we use a formulation of GP classification that requires an optimization procedure, comparisons of random feature approximations become more involved, as we need to perform convergence verification for the optimization procedure.

²A Hermitian matrix $\boldsymbol{\Sigma} \in \mathbb{C}^{n \times n}$ is called positive semi-definite, if for all $\mathbf{v} \in \mathbb{C}^n$, we have $\mathbf{v}^H \boldsymbol{\Sigma} \mathbf{v} \geq 0$.

Suppose that there are training data $(\mathbf{x}_i, y_i)_{i=1}^N \subset \mathbb{R}^d \times \mathbb{C}$ for a complex-valued regression problem with $N \in \mathbb{N}$, and let $\mathbf{X} := (\mathbf{x}_1, \dots, \mathbf{x}_N)^\top \in \mathbb{R}^{N \times d}$ and $\mathbf{y} := (y_1, \dots, y_N)^\top \in \mathbb{C}^N$. We assume the following model for the training data:

$$y_i = f(x_i) + \varepsilon_i, \quad (i = 1, \dots, N), \quad (\text{B.7})$$

where $f : \mathbb{R}^d \rightarrow \mathbb{C}$ is an unknown complex-valued function, and $\varepsilon_i \sim \mathcal{CN}(0, \sigma_i^2)$ is an independent complex Gaussian noise with variance $\sigma_i^2 > 0$. Let $\boldsymbol{\sigma}^2 := (\sigma_1^2, \dots, \sigma_N^2)^\top \in \mathbb{R}^N$.

The task of complex-valued function is to estimate the unknown complex-valued function f in Eq. (B.7) from the training data $(\mathbf{x}_i, y_i)_{i=1}^N \subset \mathbb{R}^d \times \mathbb{C}$. In complex GP regression, one defines a *complex GP prior distribution* for the unknown function f , and derives a *complex GP posterior distribution* of f , given the data $(\mathbf{x}_i, y_i)_{i=1}^N \subset \mathbb{R}^d \times \mathbb{C}$ and the likelihood function given by Eq. (B.7). For the prior, we focus on a *proper* complex GP (Boloix-Tortosa et al., 2018, Section II-C), which we describe below.

Proper complex Gaussian processes. A complex-valued function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ is called *positive definite kernel*, if 1) $k(\mathbf{x}, \mathbf{x}') = \overline{k(\mathbf{x}', \mathbf{x})}$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$; and ii) for all $n \in \mathbb{N}$ and all $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, the matrix $\mathbf{K} \in \mathbb{C}^{n \times n}$ with $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ satisfies $\mathbf{v}^H \mathbf{K} \mathbf{v} \geq 0$.

Let $f : \mathbb{R}^d \rightarrow \mathbb{C}$ be a zero-mean complex-valued stochastic process, and $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ be a positive definite kernel. We call f a (zero-mean) *proper complex GP* with covariance kernel k , if for all $n \in \mathbb{N}$ and all $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, the random vector $\mathbf{f} := (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^\top \in \mathbb{C}^n$ follows the proper complex Gaussian distribution $\mathcal{CN}(\mathbf{0}, \mathbf{K})$ with covariance matrix $\mathbf{K} \in \mathbb{C}^{n \times n}$ with $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$. If f is a zero-mean proper complex GP with covariance kernel k , we write $f \sim \mathcal{CGP}(0, k)$.

We now describe the approach of complex GP regression. For the unknown f in Eq. (B.7), we define a proper complex GP prior with kernel k , assuming that

$$f \sim \mathcal{CGP}(0, k) \quad (\text{B.8})$$

Then the observation model (B.7) and the prior (B.8) induce a joint distribution of the unknown function f and the training observations $\mathbf{y} = (y_1, \dots, y_N)^\top$. Conditioned on \mathbf{y} , we obtain the *posterior distribution* of f , which is also a proper complex GP (Boloix-Tortosa et al., 2018, Section II-C):

$$f \mid \mathbf{y} \sim \mathcal{CGP}(\mu_N, k_N), \quad (\text{B.9})$$

where $\mu_N : \mathbb{R}^d \rightarrow \mathbb{C}$ is the *posterior mean function* and $k_N : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ is the *posterior covariance function* given by

$$\mu_N(\mathbf{x}) := \mathbf{k}(\mathbf{x})^H (\mathbf{K} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{y}, \quad \mathbf{x} \in \mathbb{R}^d \quad (\text{B.10})$$

$$k_N(\mathbf{x}, \mathbf{x}') := k(\mathbf{x}, \mathbf{x}') - \mathbf{k}(\mathbf{x})^H (\mathbf{K} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{k}(\mathbf{x}), \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d, \quad (\text{B.11})$$

where $\mathbf{k}(\mathbf{x}) := (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_N))^\top \in \mathbb{C}^N$, $\mathbf{K} \in \mathbb{C}^{N \times N}$ with $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$, and $\text{diag}(\boldsymbol{\sigma}^2) \in \mathbb{R}^{d \times d}$ is the diagonal matrix with diagonal elements $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_N^2)^\top$. Notice that, if the kernel k is real-valued and so are the observations \mathbf{y} , Eq. (B.10) and Eq. (B.11) reduce to the posterior mean and covariance functions of standard real-valued GP regression (e.g., Rasmussen and Williams, 2006, Chapter 2). In this sense, complex GP regression with a proper GP prior is a natural complex extension of standard GP regression.

B.3.2 GP Regression with Complex Features

We next describe how to use complex features in GP regression. Let $\Phi : \mathbb{R}^d \rightarrow \mathbb{C}^D$ be a complex-valued (random) feature map,³ and let $\hat{k}(\mathbf{x}, \mathbf{x}') := \Phi(\mathbf{x})^\top \overline{\Phi(\mathbf{x}')}$ be the approximate kernel. Define

$$\Phi(\mathbf{X}) := (\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_N))^\top \in \mathbb{C}^{N \times D}, \quad \hat{\mathbf{K}} := \Phi(\mathbf{X})\Phi(\mathbf{X})^H \in \mathbb{C}^{N \times N}, \quad (\text{B.12})$$

where $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ are training inputs. Note that $\hat{\mathbf{K}}_{i,j} = \Phi(\mathbf{x}_i)^\top \overline{\Phi(\mathbf{x}_j)}$ = $\hat{k}(\mathbf{x}_i, \mathbf{x}_j)$, i.e., $\hat{\mathbf{K}}$ is the kernel matrix with kernel \hat{k} .

The approximate kernel $\hat{k} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ is complex-valued, and thus induces a proper complex GP, $f \sim \mathcal{CGP}(0, \hat{k})$. Using this GP as a prior for the unknown function f in the observation model (B.7), and conditioning on the observations $\mathbf{y} = (y_1, \dots, y_N)^\top$, we obtain the following approximate complex GP posterior:

$$f \mid \mathbf{y} \sim \mathcal{CGP}(\hat{\boldsymbol{\mu}}_N, \hat{k}_N), \quad (\text{B.13})$$

where $\hat{\boldsymbol{\mu}}_N : \mathbb{R}^d \rightarrow \mathbb{C}$ is an approximate posterior mean function and $\hat{k}_N : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ is an approximate posterior covariance function, defined as

$$\hat{\boldsymbol{\mu}}_N(\mathbf{x}) := \hat{\mathbf{k}}(\mathbf{x})^H (\hat{\mathbf{K}} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{y}, \quad \mathbf{x} \in \mathbb{R}^d \quad (\text{B.14})$$

$$\hat{k}_N(\mathbf{x}, \mathbf{x}') := \hat{k}(\mathbf{x}, \mathbf{x}') - \hat{\mathbf{k}}(\mathbf{x})^H (\hat{\mathbf{K}} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \hat{\mathbf{k}}(\mathbf{x}'), \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d, \quad (\text{B.15})$$

where $\hat{\mathbf{k}}(\mathbf{x}) := (\hat{k}(\mathbf{x}, \mathbf{x}_1), \dots, \hat{k}(\mathbf{x}, \mathbf{x}_N))^\top \in \mathbb{C}^N$, and $\hat{\mathbf{K}} \in \mathbb{C}^{N \times N}$ with $\hat{\mathbf{K}}_{i,j} = \hat{k}(\mathbf{x}_i, \mathbf{x}_j)$.

Finally, we define a real-valued approximate GP posterior using the real parts of Eq. (B.14) and Eq. (B.15). That is, define $\hat{\boldsymbol{\mu}}_{N,\mathbb{R}} : \mathbb{R}^d \rightarrow \mathbb{R}$ as the real part of the approximate posterior mean function in Eq. (B.14), and $\hat{k}_{N,\mathbb{R}}$ as the real part of the approximate covariance function in Eq. (B.15):

$$\hat{\boldsymbol{\mu}}_{N,\mathbb{R}}(\mathbf{x}) := \mathcal{R} \{ \hat{\boldsymbol{\mu}}_N(\mathbf{x}) \}, \quad \mathbf{x} \in \mathbb{R}^d, \quad (\text{B.16})$$

$$\hat{k}_{N,\mathbb{R}}(\mathbf{x}, \mathbf{x}') := \mathcal{R} \{ \hat{k}_N(\mathbf{x}, \mathbf{x}') \}, \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d. \quad (\text{B.17})$$

³Again, this subsumes the case of real-valued feature maps.

Then, we define a real-valued GP with mean function $\hat{\mu}_{N,\mathbb{R}}$ and covariance function $\hat{k}_{N,\mathbb{R}}$:

$$f|\mathbf{y} \sim \mathcal{GP}(\hat{\mu}_{N,\mathbb{R}}, \hat{k}_{N,\mathbb{R}}).$$

We use this approximate GP for prediction tasks in our experiments.

Note that naive computations of Eq. (B.14) and Eq. (B.15) require $\mathcal{O}(N^3 + N^2D)$ complexity, and thus do not leverage the computational advantage of random features. We will show next how to reformulate Eq. (B.14) and Eq. (B.15) to compute them in $\mathcal{O}(D^3 + ND^2)$, which is linear in the number of training data points N .

B.3.3 Computationally Efficient Implementation

We describe how to efficiently compute the approximate posterior mean and covariance functions in Eq. (B.14) and Eq. (B.15), respectively. To this end, recall the notation in Eq. (B.12). Let $\sigma^{-1} := (\sigma_1^{-1}, \dots, \sigma_N^{-1})^\top \in \mathbb{R}^N$ and $\sigma^{-2} := (\sigma_1^{-2}, \dots, \sigma_N^{-2})^\top \in \mathbb{R}^N$.

First we deal with Eq. (B.14). For a matrix $\mathbf{A} \in \mathbb{C}^{N \times D}$, we have $(\mathbf{A}^H \mathbf{A} + \mathbf{I}_N) \mathbf{A}^H = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H + \mathbf{I}_D)$, and thus $\mathbf{A}^H (\mathbf{A} \mathbf{A}^H + \mathbf{I}_N)^{-1} = (\mathbf{A}^H \mathbf{A} + \mathbf{I}_D)^{-1} \mathbf{A}^H$. By using this last identity with $A = \text{diag}(\sigma^{-1}) \Phi(\mathbf{X}) \in \mathbb{C}^{N \times D}$, we can rewrite Eq. (B.14) as

$$\begin{aligned} \hat{\mu}_N(\mathbf{x}) &= \hat{\mathbf{k}}(\mathbf{x})^H (\hat{\mathbf{K}} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{y}, \\ &= \Phi(\mathbf{x})^\top \Phi(\mathbf{X})^H (\Phi(\mathbf{X}) \Phi(\mathbf{X})^H + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{y} \\ &= \Phi(\mathbf{x})^\top \Phi(\mathbf{X})^H \text{diag}(\boldsymbol{\sigma}^{-1}) (\text{diag}(\boldsymbol{\sigma}^{-1}) \Phi(\mathbf{X}) \Phi(\mathbf{X})^H \text{diag}(\boldsymbol{\sigma}^{-1}) + \mathbf{I}_N)^{-1} \text{diag}(\boldsymbol{\sigma}^{-1}) \mathbf{y} \\ &= \Phi(\mathbf{x})^\top (\Phi(\mathbf{X})^H \text{diag}(\boldsymbol{\sigma}^{-2}) \Phi(\mathbf{X}) + \mathbf{I}_D)^{-1} \Phi(\mathbf{X})^H \text{diag}(\boldsymbol{\sigma}^{-2}) \mathbf{y}. \end{aligned} \quad (\text{B.18})$$

Next we deal with Eq. (B.15). For matrices A, C, U, V of appropriate sizes with A invertible, the Woodbury matrix identity states that $A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1} = (A + UCV)^{-1}$. By using the Woodbury identity with $A = \mathbf{I}_D$, $C = \text{diag}(\boldsymbol{\sigma}^{-2})$, $U = \Phi(\mathbf{X})^H$ and $V = \Phi(\mathbf{X})$, we can rewrite Eq. (B.15) as

$$\begin{aligned} \hat{k}_N(\mathbf{x}, \mathbf{x}') &= \hat{k}(\mathbf{x}, \mathbf{x}') - \hat{\mathbf{k}}(\mathbf{x})^H (\hat{\mathbf{K}} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \hat{\mathbf{k}}(\mathbf{x}) \\ &= \Phi(\mathbf{x})^\top \overline{\Phi(\mathbf{x}')} - \Phi(\mathbf{x})^\top \Phi(\mathbf{X})^H (\Phi(\mathbf{X}) \Phi(\mathbf{X})^H + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \Phi(\mathbf{X}) \overline{\Phi(\mathbf{x}')} \\ &= \Phi(\mathbf{x})^\top \left(\mathbf{I}_D - \Phi(\mathbf{X})^H (\text{diag}(\boldsymbol{\sigma}^2) + \Phi(\mathbf{X}) \Phi(\mathbf{X})^H)^{-1} \Phi(\mathbf{X}) \right) \overline{\Phi(\mathbf{x}')} \\ &= \Phi(\mathbf{x})^\top (\mathbf{I}_D + \Phi(\mathbf{X})^H \text{diag}(\boldsymbol{\sigma}^{-2}) \Phi(\mathbf{X}))^{-1} \overline{\Phi(\mathbf{x}')} \end{aligned} \quad (\text{B.19})$$

We now study the costs of computing Eq. (B.18) and Eq. (B.19). For both Eq. (B.18) and Eq. (B.19), the bottleneck is the computation of the inverse of the following matrix.

$$\mathbf{B} := \Phi(\mathbf{X})^H \text{diag}(\boldsymbol{\sigma}^{-2}) \Phi(\mathbf{X}) + \mathbf{I}_D \in \mathbb{C}^{D \times D}. \quad (\text{B.20})$$

The time complexity of computing \mathbf{B} is $\mathcal{O}(ND^2)$, and that of the inverse \mathbf{B}^{-1} is $\mathcal{O}(D^3)$, the latter being the complexity of computing the Cholesky decomposition $\mathbf{B} = \mathbf{L}\mathbf{L}^H$ with $\mathbf{L} \in \mathbb{C}^{D \times D}$ being a lower triangular matrix. Thus, the overall cost of computing \mathbf{B}^{-1} is $\mathcal{O}(ND^2 + D^3)$.

We next conduct a more detailed analysis of the costs of \mathbf{B} and its Cholesky decomposition, and compare them with the computational costs for the corresponding matrix inversion using real-valued features (i.e., when $\Phi(\mathbf{X}) \in \mathbb{R}^{N \times D}$). Below we use the shorthand $\tilde{\Phi}(\mathbf{X}) := \text{diag}(\boldsymbol{\sigma}^{-1})\Phi(\mathbf{X})$ so that $\mathbf{B} = \tilde{\Phi}(\mathbf{X})^H \tilde{\Phi}(\mathbf{X}) + \mathbf{I}_D$. Then the real and imaginary parts of \mathbf{B} can be written as

$$\begin{aligned} \mathcal{R}\{\mathbf{B}\} &= \mathcal{R}\{\tilde{\Phi}(\mathbf{X})^H \tilde{\Phi}(\mathbf{X})\} + \mathbf{I}_D = \mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{R}\{\tilde{\Phi}(\mathbf{X})\} + \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\} + \mathbf{I}_D \\ \mathcal{I}\{\mathbf{B}\} &= \mathcal{I}\{\tilde{\Phi}(\mathbf{X})^H \tilde{\Phi}(\mathbf{X})\} = \mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\} - \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}. \end{aligned}$$

Since $(\mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\})^\top = \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}$, one can compute $\mathcal{I}\{\mathbf{B}\}$ by only computing $\mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}$. Therefore, the computation of \mathbf{B} requires the computations of the three real D -by- D matrices (i.e., $\mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}$, $\mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}$, and $\mathcal{R}\{\tilde{\Phi}(\mathbf{X})\}^\top \mathcal{I}\{\tilde{\Phi}(\mathbf{X})\}$). Thus, the total number of operations for computing \mathbf{B} is $3 \cdot (ND^2) + 2 \cdot D^2$, where $3 \cdot (ND^2)$ is operations for the matrix products and $2 \cdot D^2$ for the addition and subtraction inside $\mathcal{R}\{\mathbf{B}\}$ and $\mathcal{I}\{\mathbf{B}\}$, respectively. Hence, assuming $N \gg D$, the computational cost for \mathbf{B} is roughly 3 times more expensive than the corresponding cost for when Φ is real-valued.

Computing the Cholesky decomposition of a D by D matrix requires roughly $\frac{1}{6}D^3$ subtractions and $\frac{1}{6}D^3$ multiplications (e.g., [Trefethen and Bau, 1997](#), p. 175). Therefore, when Φ is real-valued (and thus \mathbf{B} is real-valued), the Cholesky decomposition of \mathbf{B} requires $\frac{1}{6}D^3 + \frac{1}{6}D^3 = \frac{1}{3}D^3$ FLOPS. On the other hand, when Φ is complex-valued, the Cholesky decomposition of \mathbf{B} require $\frac{4}{3}D^3$ FLOPS: one complex subtraction requires 2 real subtractions, and thus subtractions in total require $\frac{1}{6}D^3 \times 2 = \frac{1}{3}D^3$ FLOPS; one complex multiplication requires 4 real multiplications and 2 real subtractions, and thus multiplications in total require $\frac{1}{6}D^3 \times 6 = D^3$ FLOPS; thus $\frac{1}{3}D^3 + D^3 = \frac{4}{3}D^3$ FLOPS in total. Thus, the cost for computing the Cholesky decomposition of \mathbf{B} when Φ is complex-valued is 4 times more expensive than the real-valued case.

The memory requirement for the complex case is 2 times as large as the real case, since the complex case requires storing both real and imaginary parts.

Note that, if one uses a $2D$ -dimensional *real* feature map (i.e., $\Phi(\mathbf{X}) \in \mathbb{R}^{N \times 2D}$), then this requires 4 times as much memory, 4 times as many operations to compute the matrix \mathbf{B} , and 8 times as many operations for the Cholesky decomposition of \mathbf{B} as those required for a D -dimensional real feature map. Therefore, using a $2D$ -dimensional real feature map is computationally more expensive than using a D -dimensional complex feature map, since the latter only requires 2 times as much memory, 3 times as many operations for computing \mathbf{B} , and 4 times as many operations for computing the Cholesky decomposition of \mathbf{B} as those required for a

D -dimensional real feature map, as shown above. Note also that the performance improvement from using a D -dimensional complex feature map is typically larger than using a $2D$ -dimensional real feature map; see the experiments in Section 3.5.4 in Chapter 3.

B.3.4 GP Classification as Closed-form Multi-output Regression

We now describe the GP classification approach of [Milios et al. \(2018\)](#), and how to use approximate posteriors for GP regression in this approach.

We assume that there are $C \in \mathbb{N}$ classes and that output labels are expressed by one-hot-encoding. Thus, for each class $c \in \{1, \dots, C\}$ and each training input $\mathbf{x}_i \in \mathbb{R}^d$ with $i = 1, \dots, N$, there exist an output $y_{c,i} \in \{0, 1\}$ such that $y_{c,i} = 1$ if \mathbf{x}_i belongs to class c and $y_{c,i} = 0$ otherwise.

The approach of [Milios et al. \(2018\)](#). Let $\alpha > 0$ be a constant. For each class $c \in \{1, \dots, C\}$, [Milios et al. \(2018\)](#) define transformed versions $\tilde{y}_{c,1}, \dots, \tilde{y}_{c,N} \in \mathbb{R}$ of the training outputs $y_{c,1}, \dots, y_{c,N}$ as

$$\tilde{y}_{c,i} := \log(y_{c,i} + \alpha) - \sigma_{c,i}^2/2, \quad \text{where} \quad \sigma_{c,i}^2 := \log((y_{c,i} + \alpha)^{-1} + 1), \quad i = 1, \dots, N.$$

[Milios et al. \(2018\)](#) then define an observation model of $\tilde{y}_{c,1}, \dots, \tilde{y}_{c,N}$ as

$$\tilde{y}_{c,i} = f_c(\mathbf{x}_i) + \varepsilon_{c,i}, \quad i = 1, \dots, N, \quad (\text{B.21})$$

where $f_c : \mathbb{R}^d \rightarrow \mathbb{R}$ is a latent function and $\varepsilon_{c,i} \sim \mathcal{N}(0, \sigma_{c,i}^2)$ is an independent Gaussian noise with variance $\sigma_{c,i}^2$. [Milios et al. \(2018\)](#) propose to model f_c for each $c \in \{1, \dots, C\}$ independently as a GP:

$$f_c \sim \mathcal{GP}(0, k), \quad (\text{B.22})$$

where $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a kernel. [Eq. \(B.21\)](#) and [Eq. \(B.22\)](#) define the joint distribution of the latent function f_c and the transformed labels $\tilde{y}_{c,1}, \dots, \tilde{y}_{c,N}$. Thus, conditioning on $\tilde{y}_{c,1}, \dots, \tilde{y}_{c,N}$, one obtains a GP posterior of f_c . In other words, one can obtain a GP posterior of f_c by performing GP regression for each class $c \in \{1, \dots, C\}$ using $(\mathbf{x}_i, \tilde{y}_{c,i})_{i=1}^N$ as training data.

The constant α is a hyperparameter, which [Milios et al. \(2018\)](#) propose to choose by cross validation, using the Mean Negative Log Likelihood (MNLL) (e.g., [Rasmussen and Williams, 2006](#), p. 23) as an evaluation criterion.

Using approximate GP posteriors. We now explain how to use approximate posteriors for GP regression in the above approach: For each class $c \in \{1, \dots, C\}$, we perform approximate GP regression using $(\mathbf{x}_i, \tilde{y}_{c,i})_{i=1}^N$ as training data, to obtain an approximate GP posterior for the latent function f_c in [Eq. \(B.21\)](#). For instance,

with our approach on approximate GP regression using complex random features in [Appendix B.3.2](#), we obtain a GP posterior $f_c \sim \mathcal{GP}(\hat{\mu}_{N,\mathbb{R},c}, \hat{k}_{N,\mathbb{R},c})$ for each class $c \in \{1, \dots, C\}$, where $\hat{\mu}_{N,\mathbb{R},c} : \mathbb{R}^d \rightarrow \mathbb{R}$ and $\hat{k}_{N,\mathbb{R},c} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ are the approximate GP posterior mean and covariance functions in [Eq. \(B.16\)](#) and [Eq. \(B.17\)](#), respectively, with $\mathbf{y} := (\tilde{y}_{c,1}, \dots, \tilde{y}_{c,N})^\top$ and $\boldsymbol{\sigma}^2 := (\sigma_{c,1}^2, \dots, \sigma_{c,N}^2)^\top$.

For a given test input $\mathbf{x} \in \mathbb{R}^d$, one can obtain its posterior predictive probabilities over the C classes in the following way. For each class $c \in \{1, \dots, C\}$, we first generate a sample $z_c \in \mathbb{R}$ from the posterior distribution of the latent function value $f_c(\mathbf{x})$. We then apply the softmax transformation to z_1, \dots, z_C to obtain probabilities $p_1, \dots, p_C \geq 0$ over the C class labels: $p_c := \exp(z_c) / \sum_{j=1}^C \exp(z_j)$. [Milios et al. \(2018\)](#) show that these probabilities p_1, \dots, p_C are approximately a sample from a Dirichlet distribution, yielding well-calibrated predictions.

B.3.5 Kullback-Leibler (KL) Divergence

In the experiments in [Chapters 3](#) and [5](#), we use the Kullback-Leibler (KL) divergence between the exact and approximate GP posteriors, to evaluate the quality of each approximation approach. Let $\mu_{\text{exact}}(\mathbf{x})$ and $\sigma_{\text{exact}}^2(\mathbf{x})$ be the posterior mean and variance at $\mathbf{x} \in \mathbb{R}^d$ from the exact GP posterior, and let $\mu_{\text{appr}}(\mathbf{x})$ and $\sigma_{\text{appr}}^2(\mathbf{x})$ be those from an approximate GP posterior. Let $\mathbf{x}_{*,1}, \dots, \mathbf{x}_{*,m_*} \in \mathbb{R}^d$ be test input points. Define $\boldsymbol{\mu}_{\text{exact}} := (\mu_{\text{exact}}(\mathbf{x}_{*,1}), \dots, \mu_{\text{exact}}(\mathbf{x}_{*,m_*}))^\top$, $\boldsymbol{\sigma}_{\text{exact}}^2 := (\sigma_{\text{exact}}^2(\mathbf{x}_{*,1}), \dots, \sigma_{\text{exact}}^2(\mathbf{x}_{*,m_*}))^\top$, $\boldsymbol{\mu}_{\text{appr}} := (\mu_{\text{appr}}(\mathbf{x}_{*,1}), \dots, \mu_{\text{appr}}(\mathbf{x}_{*,m_*}))^\top$, and $\boldsymbol{\sigma}_{\text{appr}}^2 := (\sigma_{\text{appr}}^2(\mathbf{x}_{*,1}), \dots, \sigma_{\text{appr}}^2(\mathbf{x}_{*,m_*}))^\top$.

We then measure the KL divergence between two diagonal Gaussian distributions, $\mathcal{N}(\boldsymbol{\mu}_{\text{appr}}, \text{diag}(\boldsymbol{\sigma}_{\text{appr}}^2))$ and $\mathcal{N}(\boldsymbol{\mu}_{\text{exact}}, \text{diag}(\boldsymbol{\sigma}_{\text{exact}}^2))$:

$$\begin{aligned} & KL [\mathcal{N}(\boldsymbol{\mu}_{\text{appr}}, \text{diag}(\boldsymbol{\sigma}_{\text{appr}}^2)) \parallel \mathcal{N}(\boldsymbol{\mu}_{\text{exact}}, \text{diag}(\boldsymbol{\sigma}_{\text{exact}}^2))] \\ &= \frac{1}{2} \sum_{i=1}^{m_*} \left(\frac{\sigma_{\text{exact}}^2(\mathbf{x}_{*,i})}{\sigma_{\text{appr}}^2(\mathbf{x}_{*,i})} + \log \frac{\sigma_{\text{exact}}^2(\mathbf{x}_{*,i})}{\sigma_{\text{appr}}^2(\mathbf{x}_{*,i})} - 1 + \frac{(\mu_{\text{exact}}(\mathbf{x}_{*,i}) - \mu_{\text{appr}}(\mathbf{x}_{*,i}))^2}{\sigma_{\text{appr}}^2(\mathbf{x}_{*,i})} \right), \end{aligned} \tag{B.23}$$

We consider these diagonal Gaussian distributions, since the focus of our experiments in [Chapters 3](#) and [5](#) is the prediction performance at test input points $\mathbf{x}_{*,1}, \dots, \mathbf{x}_{*,m_*} \in \mathbb{R}^d$.

Appendix C

Appendix of Chapter 4

C.1 Pseudo-Variances of Gaussian and Rademacher Polynomial Sketches

In this section, we work out the pseudo-variances for Gaussian and Rademacher Product-Sketches. Let

$$\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y}) = \frac{1}{D} \sum_{\ell=1}^D \prod_{i=1}^p (\mathbf{w}_{i,\ell}^\top \mathbf{x})(\mathbf{w}_{i,\ell}^\top \mathbf{y})$$

be the corresponding complex kernel estimate. As $\{\mathbf{w}_{i,\ell}\}_{\ell=1}^D$ are i.i.d. for Gaussian and Rademacher sketches, the variance of $\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}) = \text{Re}\{\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})\}$ is proportional to $1/D$ and we can assume $D = 1$ here for ease of presentation and drop the index ℓ . We then rescale the variances by $1/D$ later.

Since the CtR-variance can be written as $\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] = \frac{1}{2} \text{Re}\{\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] + \mathbb{P}\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})]\}$, the only thing left to work out is the pseudo-variance $\mathbb{P}\mathbb{V}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})]$.

As our estimator is unbiased, we further have $\mathbb{E}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})] = k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y})^p$. Thus, we only need to work out $\mathbb{E}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})^2]$.

$$\mathbb{E}[\hat{k}_{\mathcal{C}}(\mathbf{x}, \mathbf{y})^2] = \mathbb{E} \left[\left(\prod_{i=1}^p \mathbf{w}_i^\top \mathbf{x} \overline{\mathbf{w}_i^\top \mathbf{y}} \right)^2 \right] = \prod_{i=1}^p \mathbb{E} \left[(\mathbf{w}_i^\top \mathbf{x})^2 (\overline{\mathbf{w}_i^\top \mathbf{y}})^2 \right] \quad (\text{C.1})$$

$$= \mathbb{E} \left[(\mathbf{w}^\top \mathbf{x})^2 (\overline{\mathbf{w}^\top \mathbf{y}})^2 \right]^p = \left(\sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d \sum_{l=1}^d \mathbb{E}[w_i w_j \overline{w_k w_l}] x_i x_j y_k y_l \right)^p \quad (\text{C.2})$$

$\mathbb{E}_{ij\bar{k}\bar{l}} := \mathbb{E}[w_i w_j \overline{w_k w_l}] \neq 0$, only if:

1. $i = j = k = l$: there are d terms $(\mathbb{E}_{ij\bar{k}\bar{l}}) x_i x_j y_k y_l = \mathbb{E}[|w_i|^4] x_i^2 y_i^2$.
2. $i = k \neq j = l$: there are $d(d-1)$ terms $(\mathbb{E}_{ij\bar{k}\bar{l}}) x_i x_j y_k y_l = \mathbb{E}[|w_i|^2] x_i y_i \mathbb{E}[|w_j|^2] x_j y_j = x_i y_i x_j y_j$.

3. $i = l \neq j = k$: there are $d(d-1)$ terms $(\mathbb{E}_{i,j\bar{k}l})x_i x_j y_k y_l = \mathbb{E}[|w_i|^2]x_i y_i \mathbb{E}[|w_j|^2]x_j y_j = x_i y_i x_j y_j$.

As for both the Gaussian and the Rademacher sketch, we have $\mathbb{E}[|w_i|^2] = 1$ for all $\{w_i\}_{i=1}^d$, we obtain:

$$\mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] = \left(\sum_{i=1}^d \mathbb{E}[|w_i|^4]x_i^2 y_i^2 + 2 \sum_{i=1}^d \sum_{j \neq i}^d x_i y_i x_j y_j \right)^p \quad (\text{C.3})$$

We have $\mathbb{E}[|w_i|^4] = 2$ and $\mathbb{E}[|w_i|^4] = 1$ for the Gaussian and Rademacher case, respectively. So the pseudo-variances $\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] - \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})]^2$ are given by the following real-valued expressions:

$$\mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = \frac{1}{D} \left((2(\mathbf{x}^\top \mathbf{y})^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right) \quad (\text{Gaussian}) \quad (\text{C.4})$$

$$\mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = \frac{1}{D} \left(\left(2(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right) \quad (\text{Rademacher}) \quad (\text{C.5})$$

where we added the $1/D$ scaling that we left out before. Note that $\mathbb{E}[|w_i|^4] \geq (\mathbb{E}[|w_i|^2])^2 = 1$ by Jensen's inequality, which is why the Rademacher sketch yields the lowest possible pseudo-variance for the estimator studied in Eq. (3.11) in Section 3.2.

C.2 Gaussian and Rademacher CtR Variance Advantage over their Real-Valued Analogs

In the following, we compare Gaussian and Rademacher CtR-sketches against their real-valued analogs assuming that the corresponding feature maps have equal dimensions. Thus, we assign $2D$ random features to the real feature map $\Phi_R : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$ and only D random features to the CtR feature map $\Phi_{\text{CtR}} : \mathbb{R}^d \rightarrow \mathbb{R}^{2D}$ as it has twice as many dimensions for the same number of features.

We call the corresponding kernel estimates $\hat{k}_R(\mathbf{x}, \mathbf{y}) = \Phi_R(\mathbf{x})^\top \Phi_R(\mathbf{y})$ and $\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y}) = \Phi_{\text{CtR}}(\mathbf{x})^\top \Phi_{\text{CtR}}(\mathbf{y})$. $\mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})]$ is given in Eq. (3.8) and Eq. (3.7).

We further have $\mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] = \frac{1}{2}(\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] + \mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})])$ as shown in Section 4.2. $\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})]$ is given in Eq. (3.19) and Eq. (3.18). $\mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})]$ is given in Eq. (C.4) and Eq. (C.5), respectively.

We start with the simpler Gaussian case and study the Rademacher case after.

Gaussian case: Proof of Theorem (4.2.4).

Proof. Taking into account the different number of random features for Φ_R and Φ_{CtR} to have equal dimensions, the variance difference of their kernel estimates

yields:

$$\begin{aligned}
& \mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] \\
&= \frac{1}{2D} \left((\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + 2(\mathbf{x}^\top \mathbf{y})^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right) \\
&\quad - \frac{1}{2D} \left((\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2)^p + (2(\mathbf{x}^\top \mathbf{y})^2)^p - 2(\mathbf{x}^\top \mathbf{y})^{2p} \right) \\
&= \frac{1}{2D} \left((\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + 2(\mathbf{x}^\top \mathbf{y})^2)^p - (2(\mathbf{x}^\top \mathbf{y})^2)^p \right) - \frac{1}{2D} \left((\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right) \\
&= \frac{1}{2D} \sum_{k=0}^{p-1} \binom{p}{k} (2(\mathbf{x}^\top \mathbf{y})^2)^k (\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^{p-k} - \frac{1}{2D} \sum_{k=0}^{p-1} \binom{p}{k} (\mathbf{x}^\top \mathbf{y})^{2k} (\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^{p-k} \\
&= \frac{1}{2D} \sum_{k=0}^{p-1} \binom{p}{k} (2^k - 1) (\mathbf{x}^\top \mathbf{y})^{2k} (\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^{p-k} \geq 0
\end{aligned}$$

■

Thus, the Gaussian CtR-estimator is always better regardless of the choice of \mathbf{x}, \mathbf{y} and p and despite using only half the random feature samples. Note that the variance difference is zero if $p = 1$ and increases as p increases. Moreover, the difference is maximized for parallel \mathbf{x} and \mathbf{y} . In this case, we have $(\mathbf{x}^\top \mathbf{y}) = \|\mathbf{x}\| \|\mathbf{y}\|$ and the difference becomes

$$\begin{aligned}
\mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] &= \frac{1}{2D} \sum_{k=0}^{p-1} \binom{p}{k} (2^k - 1) (\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^p \\
&= \frac{1}{2D} \|\mathbf{x}\|^{2p} \|\mathbf{y}\|^{2p} (3^p - 2^{p+1} + 1)
\end{aligned}$$

We analyze the more difficult Rademacher case next.

Rademacher case: Proof of Theorem (4.2.3).

Proof. Taking into account the different number of random features for Φ_R and Φ_{CtR} to have equal dimensions, the variance difference of their kernel estimates yields:

$$\begin{aligned}
& \mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CtR}}(\mathbf{x}, \mathbf{y})] \\
&= \frac{1}{2D} \left(\left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + 2 \sum_{i=1}^d \sum_{j \neq i} x_i x_j y_i y_j \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right) \\
&\quad - \frac{1}{2D} \left\{ \left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + \sum_{i=1}^d \sum_{j \neq i} x_i x_j y_i y_j \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} + \left(2(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right\}
\end{aligned}$$

Next, we write $(\mathbf{x}^\top \mathbf{y})^{2p} = ((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 + \sum_{i=1}^d x_i^2 y_i^2)^p$. In this way, we can factor out the term $a := (\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 = \sum_{i=1}^d \sum_{j \neq i} x_i x_j y_i y_j$ and apply the binomial theorem to *all* addends. This gives:

$$\begin{aligned} \mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CTR}}(\mathbf{x}, \mathbf{y})] &= \frac{1}{2D} \sum_{k=0}^p \binom{p}{k} a^{p-k} \\ &\quad \left(\left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^k - \left((\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^k + (\mathbf{x}^\top \mathbf{y})^{2k} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^k \right) \right) \end{aligned}$$

We now show that the following term is always non-negative:

$$B := \left(\left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^k - \left((\|\mathbf{x}\|^2 \|\mathbf{y}\|^2)^k + (\mathbf{x}^\top \mathbf{y})^{2k} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^k \right) \right) \quad (\text{C.6})$$

For $k = 0$ and $k = 1$, $B = 0$. For $k \geq 2$, we have:

$$\left(\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 + (\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^k = \sum_{j=0}^k \binom{k}{j} \|\mathbf{x}\|^{2j} \|\mathbf{y}\|^{2j} \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^{k-j}$$

Plugging this expression into B and cancelling out the addend for $j = k$ yields:

$$B = \sum_{j=0}^{k-1} \binom{k}{j} \|\mathbf{x}\|^{2j} \|\mathbf{y}\|^{2j} \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^{k-j} - \left((\mathbf{x}^\top \mathbf{y})^{2k} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^k \right)$$

Next, we refactor $(\mathbf{x}^\top \mathbf{y})^{2k} - (\sum_{i=1}^d x_i^2 y_i^2)^k$:

$$\begin{aligned} (\mathbf{x}^\top \mathbf{y})^{2k} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^k &= ((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 + \sum_{i=1}^d x_i^2 y_i^2)^k - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^k \\ &= \sum_{j=0}^k \binom{k}{j} \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^j ((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2)^{k-j} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^k \\ &= \sum_{j=0}^{k-1} \binom{k}{j} \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^j ((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2)^{k-j} \end{aligned}$$

Plugging this expression into B yields:

$$B = \sum_{j=0}^{k-1} \binom{k}{j} \left(\|\mathbf{x}\|^{2j} \|\mathbf{y}\|^{2j} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^j \right) \left((\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^{k-j}$$

Finally, we insert B back into the original variance difference $\mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CTR}}(\mathbf{x}, \mathbf{y})]$ (remember that $B = 0$ if $k < 2$):

$$\begin{aligned} \mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CTR}}(\mathbf{x}, \mathbf{y})] &= \frac{1}{2D} \sum_{k=2}^p \binom{p}{k} a^{p-k} B \\ &= \frac{1}{2D} \sum_{k=2}^p \sum_{j=0}^{k-1} \binom{p}{k} \binom{k}{j} a^{p-j} \left(\|\mathbf{x}\|^{2j} \|\mathbf{y}\|^{2j} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^j \right) \end{aligned}$$

Finally, we note that $b := \|\mathbf{x}\|^{2j} \|\mathbf{y}\|^{2j} - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^j = \left(\sum_{i=1}^d \sum_{\ell=1}^d x_i^2 y_\ell^2 \right)^j - \left(\sum_{i=1}^d x_i^2 y_i^2 \right)^j \geq 0$ and $\mathbb{V}[\hat{k}_R(\mathbf{x}, \mathbf{y})] - \mathbb{V}[\hat{k}_{\text{CTR}}(\mathbf{x}, \mathbf{y})] \geq 0$ if $a = \sum_{i=1}^d \sum_{j \neq i} x_i x_j y_i y_j^2 \geq 0$. \blacksquare

C.3 Pseudo-Variance of TensorSRHT

In the following, we work out the pseudo-variance of TensorSRHT. We focus on upsampled TensorSRHT here. The pseudo-variance of stacked TensorSRHT can be derived analogously by applying the same modifications as the ones used in Section B.2.2.

For the pseudo-variance $\mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})] = \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] - \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})]^2$, we need to work out $\mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2]$:

$$\begin{aligned} \mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] &= \frac{1}{D^2} \sum_{\ell=1}^D \sum_{\ell'=1}^D \prod_{i=1}^p \mathbb{E} \left[(\mathbf{w}_{i,\ell}^\top \mathbf{x}) (\overline{\mathbf{w}_{i,\ell}^\top \mathbf{y}}) (\mathbf{w}_{i,\ell'}^\top \mathbf{x}) (\overline{\mathbf{w}_{i,\ell'}^\top \mathbf{y}}) \right] \\ &= \frac{1}{D^2} \sum_{\ell=1}^D \sum_{\ell'=1}^D \underbrace{\mathbb{E} \left[(\mathbf{w}_\ell^\top \mathbf{x}) (\overline{\mathbf{w}_\ell^\top \mathbf{y}}) (\mathbf{w}_{\ell'}^\top \mathbf{x}) (\overline{\mathbf{w}_{\ell'}^\top \mathbf{y}}) \right]^p}_{e(\ell, \ell')} \end{aligned}$$

We dropped the index i in the last equality for ease of notation, as all $\{\mathbf{w}_{i,\ell}\}_{i=1}^p$ are i.i.d. samples and the expectation is thus the same for any i . To work out the expectation $e(\ell, \ell')$, we need to distinguish different cases for ℓ and ℓ' .

1. $\ell = \ell'$ (D terms): $e(\ell, \ell')^p = \mathbb{E} \left[(\mathbf{w}_\ell^\top \mathbf{x})^2 (\overline{\mathbf{w}_\ell^\top \mathbf{y}})^2 \right]^p = \left(2(\mathbf{x}^\top \mathbf{y}^2) - \sum_{i=1}^d x_i^2 y_i^2 \right)^p$
(taken from Eq. (C.3) for the Rademacher case)
2. $\ell \neq \ell'$ ($D(D-1)$ terms):

$$\begin{aligned} e(\ell, \ell')^p &= \left(\sum_{q=1}^d \sum_{r=1}^d \sum_{s=1}^d \sum_{t=1}^d \mathbb{E} [w_{\ell,q} \overline{w_{\ell,r}} w_{\ell',s} \overline{w_{\ell',t}}] x_q y_r x_s y_t \right)^p \\ &= \left(\sum_{q=1}^d \sum_{r=1}^d \sum_{s=1}^d \sum_{t=1}^d \mathbb{E} [d_q \overline{d_r} d_s \overline{d_t}] \mathbb{E} [h_{p\ell,q} h_{p\ell,r} h_{p\ell',s} h_{p\ell',t}] x_q y_r x_s y_t \right)^p \end{aligned}$$

d_q, d_r, d_s, d_t are uniform samples from $\{1, -1, i, -i\}$, i.e., complex Rademacher samples, that are independent from the index samples $p_{\ell,q}, p_{\ell,r}, p_{\ell,s}, p_{\ell,t}$, which is why we can factor out the two expectations. We will simplify the above sum by studying when $\mathbb{E}[d_q \bar{d}_r d_s \bar{d}_t] \neq 0$.

We have to distinguish three non-zero cases for $\mathbb{E}[d_q \bar{d}_r d_s \bar{d}_t]$:

1. $q = r = s = t$ (d terms): $\mathbb{E}[d_q \bar{d}_r d_s \bar{d}_t] = \mathbb{E}[|d_q|^4] = 1$
2. $q = r \neq s = t$ ($d(d-1)$ terms): $\mathbb{E}[d_q \bar{d}_r d_s \bar{d}_t] = \mathbb{E}[|d_q|^2] \mathbb{E}[|d_s|^2] = 1$
3. $q = t \neq r = s$ ($d(d-1)$ terms): $\mathbb{E}[d_q \bar{d}_r d_s \bar{d}_t] = \mathbb{E}[|d_q|^2] \mathbb{E}[|d_r|^2] = 1$

because $\mathbb{E}[|d_q|^4] = \mathbb{E}[|d_q|^2] = 1$.

In the second part of [Lemma B.2.1](#), we show that for $\ell \neq \ell'$ and $q \neq r$, $\mathbb{E}[h_{p_{\ell,q}} h_{p_{\ell',r}} h_{p_{\ell',r}} h_{p_{\ell,q}}] = -\frac{1}{\lceil D/d \rceil d - 1}$ holds. Therefore, $e(\ell, \ell')^p$ for $\ell \neq \ell'$ yields:

$$\begin{aligned} e(\ell, \ell')^p &= \left(\sum_{i=1}^d x_i^2 y_i^2 + \sum_{i=1}^d \sum_{j \neq i}^d x_i y_i x_j y_j - \frac{1}{\lceil D/d \rceil d - 1} \sum_{i=1}^d \sum_{j \neq i}^d x_i y_i x_j y_j \right)^p \\ &= \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{\lceil D/d \rceil d - 1} \left[(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right] \right)^p \end{aligned}$$

In fact, $e(\ell, \ell')^p$ does not depend on ℓ and ℓ' anymore after working out the expectations involved. Plugging $e(\ell, \ell')^p$ back into $\mathbb{E}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2]$ yields the following pseudo-variance for TensorSRHT:

$$\begin{aligned} \mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] &= \frac{1}{D} \left[\left(2(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right] \\ &\quad + \left(1 - \frac{1}{D} \right) \left[\left((\mathbf{x}^\top \mathbf{y})^2 - \frac{1}{\lceil D/d \rceil d - 1} \left[(\mathbf{x}^\top \mathbf{y})^2 - \sum_{i=1}^d x_i^2 y_i^2 \right] \right)^p - (\mathbf{x}^\top \mathbf{y})^{2p} \right] \\ &= \frac{1}{D} \mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)} - \left(1 - \frac{1}{D} \right) \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)}}{\lceil D/d \rceil d - 1} \right)^p \right] \end{aligned} \tag{C.7}$$

$\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)}$ and $\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)}$ are the Rademacher pseudo-variance ([C.5](#)) for a given degree p and $p = 1$, respectively.

For stacked TensorSRHT, the pseudo-variance becomes

$$\mathbb{P}\mathbb{V}[\hat{k}_C(\mathbf{x}, \mathbf{y})^2] = \frac{1}{D} \mathbb{P}\mathbb{V}_{\text{Rad.}}^{(p)} - \left(\frac{c(D, d)}{D^2} \right) \left[(\mathbf{x}^\top \mathbf{y})^{2p} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{\mathbb{P}\mathbb{V}_{\text{Rad.}}^{(1)}}{d-1} \right)^p \right].$$

with $c(D, d) = \lfloor D/d \rfloor d(d-1) + \text{mod}(D, d)(\text{mod}(D, d) - 1)$. It can be derived using the same changes to the above derivation as presented in [Section B.2.2](#).

Appendix D

Appendix of Chapter 5

D.1 Convex Surrogate Functions for Stacked TensorSRHT Variances

To extend the applicability of the Incremental Algorithm in [Algorithm 4](#) to TensorSRHT, we derive here convex surrogate functions for the variances of *stacked* TensorSRHT. To this end, we first analyze the variances of stacked TensorSRHT in [Appendix D.1.1](#). We then derive convex surrogate functions in [Appendix D.1.2](#).

D.1.1 Analyzing the Variances of stacked TensorSRHT

We first derive another form of the variance of stacked TensorSRHT given in [Eq. \(3.35\)](#) of [Theorem 3.5.9](#), which we will use in a later analysis. Let $\Phi_n : \mathbb{R}^d \rightarrow \mathbb{C}^D$ be a complex TensorSRHT sketch of degree $n \in \mathbb{N}$ satisfying the assumptions in [Theorem 3.5.9](#) with $0 \leq q \leq 1$. For $q = 1$ we recover the real TensorSRHT and for $q = 1/2$ the complex one.

As shown in [Appendix B.2.2](#), the approximate kernel of the complex TensorSRHT can be written as

$$\hat{k}(\mathbf{x}, \mathbf{y}) := \Phi_n(\mathbf{x})^\top \overline{\Phi_n(\mathbf{y})} = \frac{1}{D} \sum_{b=1}^B \sum_{\ell \in I_b} \prod_{i=1}^n (\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) \overline{(\mathbf{s}_{i,\ell}^{b\top} \mathbf{y})},$$

where $B := \lceil D/d \rceil$, $I_b := \{1, \dots, d\}$ for $b = 1, \dots, B-1$ and $I_b := \{1, \dots, \text{mod}(D, d)\}$ for $b = B$, and $\mathbf{s}_{i,\ell}^b \in \mathbb{C}^d$ are the structured random weights defined in [Eq. \(B.5\)](#).

We can then write the variance of the approximate kernel as

$$\begin{aligned}
\mathbb{V}[\hat{k}(\mathbf{x}, \mathbf{y})] &= \frac{1}{D^2} \sum_{b=1}^B \mathbb{V} \left[\sum_{\ell \in I_b} \prod_{i=1}^n (\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) (\mathbf{s}_{i,\ell}^{b\top} \mathbf{y}) \right] \\
&= \frac{1}{D^2} \sum_{b=1}^B \sum_{\ell \in I_b} \underbrace{\mathbb{V} \left[\prod_{i=1}^n (\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) (\mathbf{s}_{i,\ell}^{b\top} \mathbf{y}) \right]}_{= V_q^{(n)}} \\
&\quad + \frac{1}{D^2} \sum_{b=1}^B \sum_{\substack{\ell, \ell' \in I_b, \\ \ell \neq \ell'}} \underbrace{\text{Cov} \left(\prod_{i=1}^n (\mathbf{s}_{i,\ell}^{b\top} \mathbf{x}) (\mathbf{s}_{i,\ell}^{b\top} \mathbf{y}), \prod_{i=1}^n (\mathbf{s}_{i,\ell'}^{b\top} \mathbf{x}) (\mathbf{s}_{i,\ell'}^{b\top} \mathbf{y}) \right)}_{=: \text{Cov}_q^{(n)}} \\
&= \frac{V_q^{(n)}}{D} + \frac{c(D, d)}{D^2} \text{Cov}_q^{(n)} = \text{Eq. (3.35)}, \tag{D.1}
\end{aligned}$$

where $c(D, d) = \lfloor D/d \rfloor d(d-1) + \text{mod}(D, d)(\text{mod}(D, d) - 1)$ and the last line follows from that the values of $V_q^{(n)}$ and $\text{Cov}_q^{(n)}$ do not depend on the choice of ℓ, ℓ' and b (which can be shown from the arguments in [Appendix B.2.2](#)). Here, $V_q^{(n)}$ is the variance of the unstructured Rademacher sketch with a single feature in [Eq. \(3.17\)](#) with $p = n$, and $\text{Cov}_q^{(n)}$ is the covariance for distinct indices ℓ, ℓ' inside each block b . By comparing [Eq. \(3.35\)](#) and [Eq. \(D.1\)](#), the concrete form of $\text{Cov}_q^{(n)}$ is given by

$$\text{Cov}_q^{(n)} = - \left[(\mathbf{x}^\top \mathbf{y})^{2n} - \left((\mathbf{x}^\top \mathbf{y})^2 - \frac{V_q^{(1)}}{d-1} \right)^n \right]$$

[Eq. \(D.1\)](#) is a useful representation of the variance of TensorSRHT in [Eq. \(3.35\)](#) for studying its (non-)convexity with respect to D . The following result shows a range of values of D for which [Eq. \(3.35\)](#) is convex.

Theorem D.1.1. *The variance of the TensorSRHT sketch in [Eq. \(3.35\)](#) is convex and monotonically decreasing with respect to $D \in \{1, \dots, d\}$ and with respect to $D \in \{kd \mid k \in \mathbb{N}\}$.*

Proof. If $D \in \{1, \dots, d\}$, we have $c(D, d) = D(D-1)$ in [Eq. \(D.1\)](#). Therefore, [Eq. \(D.1\)](#) is equal to

$$\frac{1}{D} V_q^{(n)} + \left(1 - \frac{1}{D}\right) \text{Cov}_q^{(n)} = \frac{1}{D} \left(V_q^{(n)} - \text{Cov}_q^{(n)} \right) + \text{Cov}_q^{(n)}. \tag{D.2}$$

For two random variables X, Y it generally holds that $|\text{Cov}(X, Y)| \leq \sqrt{\mathbb{V}[X]\mathbb{V}[Y]}$ by the Cauchy-Schwarz inequality. Hence, we have $|\text{Cov}_q^{(n)}| \leq V_q^{(n)}$ and thus $V_q^{(n)} - \text{Cov}_q^{(n)} \geq 0$. Therefore, [Eq. \(D.2\)](#) is proportional to $1/D$ with a non-negative coefficient, and thus it is convex and monotonically decreasing for $D \in \{1, \dots, d\}$.

Next, suppose $D = kd$ for some $k \in \mathbb{N}$, in which case we have $c(D, d) = kd(d-1)$ in Eq. (D.1). Therefore Eq. (D.1) is equal to

$$\frac{1}{kd} \left(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \right) = \frac{1}{D} \left(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \right). \quad (\text{D.3})$$

The term in the parenthesis is non-negative, because (D.1) is the variance of TensorSRHT and thus non-negative. Therefore, (D.1) is convex and monotonically decreasing with respect to $D \in \{kd \mid k \in \mathbb{N}\}$. ■

As we do next, Theorem D.1.1 is useful for designing a convex surrogate function for Eq. (3.35), as it shows the range of D on which Eq. (3.35) is already convex and does not need to be modified.

D.1.2 Convex Surrogate Functions

Based on Eq. (3.35), we now propose a convex surrogate function for the variance of TensorSRHT in Eq. (3.35). We consider the following two cases separately: i) $\text{Cov}_q^{(n)} \leq 0$ and ii) $\text{Cov}_q^{(n)} > 0$. For each case, we propose a convex surrogate function.

i) Case $\text{Cov}_q^{(n)} \leq 0$. We define a surrogate function of Eq. (3.35) by concatenating the two expressions of Eq. (D.1) for $D \in \{1, \dots, d\}$ and $D \in \{kd \mid k \in \mathbb{N}\}$ given in Eq. (D.2) and Eq. (D.3), respectively, and extend their ranges to the entire domain $D \in \mathbb{N}$:

$$V_{\text{Surr.}}^{(n)}(D) := \begin{cases} \frac{1}{D} \left(V_q^{(n)} - \text{Cov}_q^{(n)} \right) + \text{Cov}_q^{(n)} & \text{if } D \leq d \\ \frac{1}{D} \left(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \right) & \text{if } D > d. \end{cases} \quad (\text{D.4})$$

ii) Case $\text{Cov}_q^{(n)} > 0$. We use the expression (D.3) to define a surrogate function on $D \in \mathbb{N}$:

$$V_{\text{Surr.}}^{(n)}(D) := \frac{1}{D} \left(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \right) \quad (\text{D.5})$$

The convexity of Eq. (D.5) immediately follows from $V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \geq 0$, which holds as we show in the proof of Theorem D.1.1. Note that $\text{Cov}_q^{(n)} > 0$ can only occur when n is even, as shown in Corollary 3.5.10 of Section 3.5.

We defined the surrogate function in Eq. (D.4) by interpolating the variances of TensorSRHT in Eq. (3.35) for $D \in \{1, \dots, d\}$ and $D \in \{kd \mid k \in \mathbb{N}\}$ and extending the domain to \mathbb{N} . In fact, for $D \in \{1, \dots, d\}$ and $D \in \{kd \mid k \in \mathbb{N}\}$, Eq. (D.4) is equal to Eq. (3.35), as shown in the proof of Theorem D.1.1. Fig. D.1 illustrates

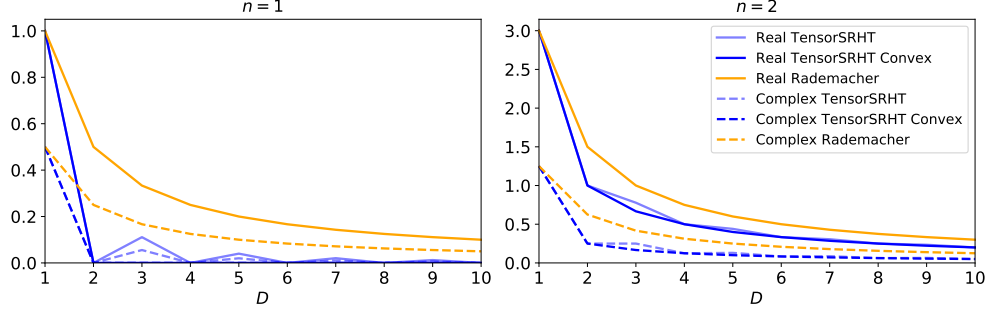


Figure D.1: Convex surrogate functions in Eq. (D.4) and the variances of TensorSRHT in (3.35) as a function of the number of random features D , with polynomial degrees $n = 1, 2$ and input vectors $\mathbf{x} = \mathbf{y} = [\sqrt{1/2}, \sqrt{1/2}]^\top$ ($d = 2$). For a comparison, we also plot the variances of the real Rademacher sketch in Eq. (3.7) and the complex Rademacher sketch in Eq. (3.18).

the convex surrogate function in Eq. (D.4) and the variance of TensorSRHT in (3.35) when $\text{Cov}_q^{(n)} \leq 0$ holds.

Note that, as mentioned later in Remark D.1.3, the surrogate function in Eq. (D.4) may not be convex over $D \in \mathbb{N}$ if the condition $\text{Cov}_q^{(n)} \leq 0$ does not hold. This is why we defined another convex surrogate function as in Eq. (D.5) for the case $\text{Cov}_q^{(n)} > 0$.

The following theorem shows that the surrogate function in Eq. (D.4) is convex in the considered case of i) $\text{Cov}_q^{(n)} \leq 0$.

Theorem D.1.2. *If $\text{Cov}_q^{(n)} \leq 0$, Eq. (D.4) is convex with respect to $D \in \mathbb{N}$.*

Proof. As shown in Theorem D.1.1, $V_{\text{Surr.}}^{(n)}(D) = \frac{1}{D}(V_q^{(n)} - \text{Cov}_q^{(n)}) + \text{Cov}_q^{(n)}$ is convex over $D \in \{1, \dots, d\}$. Likewise, $V_{\text{Surr.}}^{(n)}(D) = \frac{1}{D}(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)})$ is convex over $D \in [d, \infty) \cap \mathbb{N}$, since $V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \geq 0$ holds as we show in the proof of Theorem D.1.1.

Therefore, the proof completes by showing that the concatenated function $V_{\text{Surr.}}^{(n)}(D)$ in Eq. (D.4) is also convex over $D \in \{d-1, d, d+1\}$, i.e.,

$$\frac{1}{2} \left(V_{\text{Surr.}}^{(n)}(d-1) + V_{\text{Surr.}}^{(n)}(d+1) \right) \geq V_{\text{Surr.}}^{(n)}(d). \quad (\text{D.6})$$

By using the definition in Eq. (D.4), this inequality is equivalent to

$$\begin{aligned} & \frac{1}{2} \left(\frac{1}{d-1} \left(V_q^{(n)} + (d-2)\text{Cov}_q^{(n)} \right) + \frac{1}{d+1} \left(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \right) \right) \\ & \geq \frac{1}{d} \left(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \right). \end{aligned} \quad (\text{D.7})$$

Note that we have $V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} \geq 0$, as mentioned earlier. If $V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} = 0$ holds, then we have $V_{\text{Surr.}}^{(n)}(D) = 0$ for $D \geq d$ by the definition in

Eq. (D.4), and thus Eq. (D.6) holds (which concludes the proofs). Therefore, we assume that the inequality to be strict, i.e., $V_q^{(n)} + (d-1)\text{Cov}_q^{(n)} > 0$. Dividing the both sides of Eq. (D.7) by $(V_q^{(n)} + (d-1)\text{Cov}_q^{(n)})$, we obtain

$$\frac{1}{2} \left(\frac{1}{d-1} \frac{V_q^{(n)} + (d-2)\text{Cov}_q^{(n)}}{V_q^{(n)} + (d-1)\text{Cov}_q^{(n)}} + \frac{1}{d+1} \right) \geq \frac{1}{d},$$

which after some rearrangement gives

$$\frac{V_q^{(n)} + (d-2)\text{Cov}_q^{(n)}}{V_q^{(n)} + (d-1)\text{Cov}_q^{(n)}} \geq 1 - \frac{2}{d^2 + d}. \quad (\text{D.8})$$

This inequality holds because we have $(d-2)\text{Cov}_q^{(n)} \geq (d-1)\text{Cov}_q^{(n)}$, which follows from our assumption $\text{Cov}_q^{(n)} \leq 0$. Therefore Eq. (D.7) holds. ■

Remark D.1.3. *Theorem D.1.2 shows the convexity of the surrogate function in Eq. (D.4), assuming $\text{Cov}_q^{(n)} \leq 0$. If this condition does not hold, i.e., if $\text{Cov}_q^{(n)} > 0$, then the surrogate function in Eq. (D.4) may not be convex. To see this, let $d = 2$, $\mathbf{x} = (a, 0)^\top$ with $a > 0$, $\mathbf{y} = (0, b)^\top$ with $b > 0$, and n be even; then we have $V_q^{(n)} = \text{Cov}_q^{(n)} = a^{2n}b^{2n} > 0$, and the inequality in Eq. (D.8) in the proof of Theorem D.1.2 does not hold, which implies that the surrogate function in Eq. (D.4) is not convex.*

As mentioned in Section 3.5, the variance of TensorSRHT in Eq. (3.35) becomes zero if $n = 1$ and $D \in \{kd \mid k \in \mathbb{N}\}$, i.e., $\mathbb{V}[\Phi_1(\mathbf{x})^\top \Phi_1(\mathbf{y})] = 0$ holds. Therefore, because the convex surrogate functions in Eq. (D.4) and Eq. (D.5) are equal to the variance of TensorSRHT in Eq. (3.35) for $D \in \{kd \mid k \in \mathbb{N}\}$, these surrogate functions also become zero for $n = 1$ and $D \in \{kd \mid k \in \mathbb{N}\}$. Thus, the Incremental Algorithm (Algorithm 4), when used with the surrogate functions in Eq. (D.4) and Eq. (D.5), will not assign more than $D = d$ random features to the polynomial degree $n = 1$. Note that assigning $D = d$ random features is equivalent to appending the input vectors \mathbf{x} and \mathbf{y} to the approximate kernel (5.10), which is called *H0/1 heuristic* in Kar and Karnick (2012). Therefore, the Incremental Algorithm with the surrogate functions in Eq. (D.4) and Eq. (D.5) automatically achieve the H0/1 heuristic.

Finally, we describe briefly how to use the convex surrogate functions in Eq. (D.4) and Eq. (D.5) in the Incremental Algorithm in Algorithm 4. To this end, we rewrite Eq. (5.20) using the surrogate functions as follows: (Here, we make the dependence of $V_q^{(n)}$ and $\text{Cov}_q^{(n)}$ on the input vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ explicit and write

them as $V_q^{(n)}(\mathbf{x}, \mathbf{y})$ and $\text{Cov}_q^{(n)}(\mathbf{x}, \mathbf{y})$, respectively.)

$$\text{Eq. (5.20)} = \begin{cases} \frac{a_n^2}{D_n} \left(\sum_{i \neq j} V_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j) + (d-1) \sum_{i \neq j} \text{Cov}_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j) \right) \\ \quad \text{if } \sum_{i \neq j} \text{Cov}_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j) > 0 \text{ or } D_n > d, \\ \frac{a_n^2}{D_n} \left(\sum_{i \neq j} V_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i \neq j} \text{Cov}_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j) \right) + a_n^2 \sum_{i \neq j} \text{Cov}_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j) \\ \quad \text{otherwise.} \end{cases}$$

After precomputing the constants $\sum_{i \neq j} V_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j)$ and $\sum_{i \neq j} \text{Cov}_q^{(n)}(\mathbf{x}_i, \mathbf{x}_j)$ for each $n \in \{1, \dots, p\}$, which can be done in $\mathcal{O}(m^2)$ time, one can directly use the above modification of Eq. (5.20) in the objective function in Eq. (5.19). In this way, we adapt the objective function in (5.19) to be convex, so that the Incremental Algorithm in Algorithm 4 is directly applicable.

D.2 Additional Experiments

We present here additional experimental results, supplementing those in Chapter 5. Table D.1 shows the effects of applying zero-centering to input vectors in the polynomial kernel approximation experiments. Fig. D.2 and Fig. D.3 show the results of additional experiments on GP regression. Fig. D.4 and Fig. D.5 show the results of additional experiments on GP classification.

Dataset	MNLL						Rel. Prob. Error			
	Non-centred			Centred			Non-centred		Centred	
	SRF Gaus.	Opt. Macl. Rad.	SRF Gaus.	Opt. Macl. Rad.	SRF Gaus.	Opt. Macl. Rad.	SRF Gaus.	Opt. Macl. Rad.		
Boston	3.410±0.37	3.447±0.38	3.449±0.62	3.161±0.28	0.044±0.02	0.212±0.15	0.356±0.05	0.421±0.06		
Concrete	3.779±0.07	3.811±0.04	3.660±0.12	3.542±0.07	0.019±0.01	0.276±0.17	0.610±0.07	0.482±0.03		
Energy	6.090±0.12	6.090±0.12	5.116±0.20	5.012±0.13	0.003±0.00	0.222±0.14	0.507±0.08	0.484±0.05		
kin8nm	-0.203±0.07	-0.310±0.03	-0.203±0.07	-0.323±0.03	0.946±0.04	0.525±0.04	0.947±0.03	0.521±0.03		
Naval	-6.069±0.03	-6.066±0.03	-8.083±0.04	-7.788±0.10	0.040±0.03	0.183±0.06	0.112±0.04	0.384±0.11		
Powerplant	3.064±0.03	3.061±0.06	3.282±0.14	3.400±0.73	0.001±0.00	0.062±0.04	0.609±0.09	0.527±0.10		
Protein	3.233±0.01	3.233±0.01	3.072±0.02	3.060±0.02	0.000±0.00	0.002±0.00	0.277±0.05	0.429±0.14		
Yacht	4.317±0.45	4.478±0.45	3.773±0.21	3.844±0.28	0.028±0.01	0.276±0.11	0.512±0.04	0.484±0.03		
Cod_rna	0.307±0.00	0.308±0.00	0.288±0.06	0.151±0.01	0.022±0.01	0.087±0.05	0.641±0.05	0.467±0.05		
Coverttype	0.821±0.01	-	0.650±0.01	0.639±0.01	0.024±0.01	-	0.361±0.01	0.300±0.01		
Drive	1.446±0.02	1.453±0.03	0.677±0.02	0.497±0.01	0.068±0.02	0.135±0.05	0.348±0.01	0.312±0.02		
FashionMNIST	0.353±0.00	0.364±0.00	0.364±0.00	0.361±0.00	0.029±0.00	0.062±0.01	0.099±0.00	0.104±0.01		
Magic	0.453±0.01	0.452±0.01	0.381±0.02	0.350±0.01	0.068±0.01	0.147±0.05	0.430±0.03	0.418±0.04		
Miniboo	0.253±0.01	-	0.239±0.01	0.213±0.01	0.027±0.01	-	0.214±0.01	0.229±0.02		
MNIST	0.076±0.00	0.074±0.00	0.290±0.02	0.353±0.09	0.073±0.00	0.082±0.00	0.085±0.00	0.089±0.01		
Mocap	0.360±0.01	0.334±0.01	0.357±0.02	0.289±0.01	0.115±0.01	0.187±0.04	0.414±0.01	0.290±0.02		

Table D.1: GP regression (top) and classification (bottom) for centred vs. non-centred data with $D = 5d$ features. Non-centred Miniboo and Coverttype led to numerical issues for Maclaurin (no scores reported).

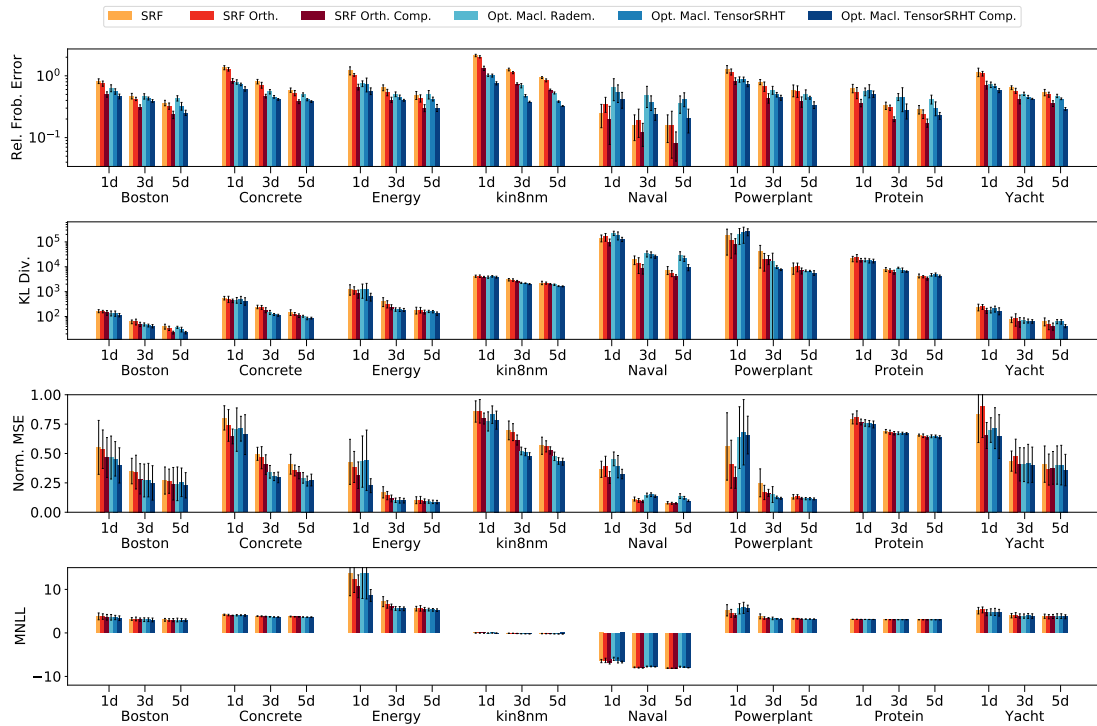


Figure D.2: Additional results of the experiments in Section 5.4.3 on approximate GP regression with a high-degree polynomial kernel. Lower values are better for all the metrics. For each dataset, we show the number of random features $D \in \{1d, 3d, 5d\}$ used in each method on the horizontal axis, with d being the input dimensionality of the dataset. We put the legend labels and the bars in the same order.

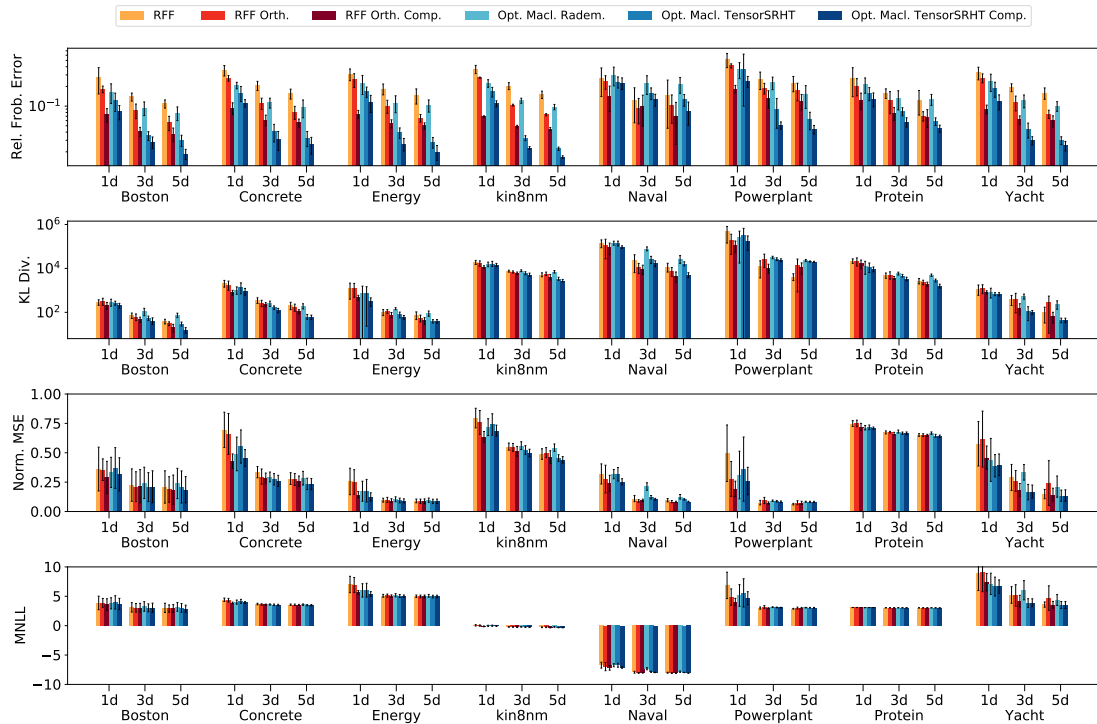


Figure D.3: Additional results of the experiments in Section 5.4.3 on approximate GP regression with a Gaussian kernel. Lower values are better for all the metrics. For each dataset, we show the number of random features $D \in \{1d, 3d, 5d\}$ used in each method on the horizontal axis, with d being the input dimensionality of the dataset. We put the legend labels and the bars in the same order.

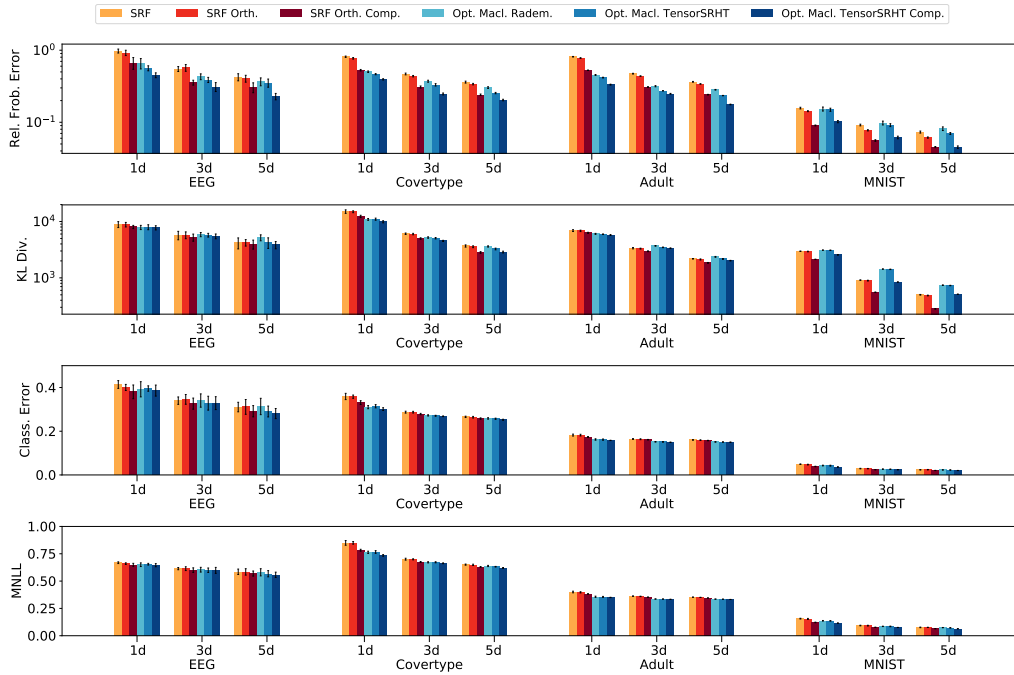


Figure D.4: Additional results of the experiments in Section 5.4.3 on approximate GP classification with a high-degree polynomial kernel. Lower values are better for all the metrics. For each dataset, we show the number of random features $D \in \{1d, 3d, 5d\}$ used in each method on the horizontal axis, with d being the input dimensionality of the dataset. We put the legend labels and the bars in the same order.

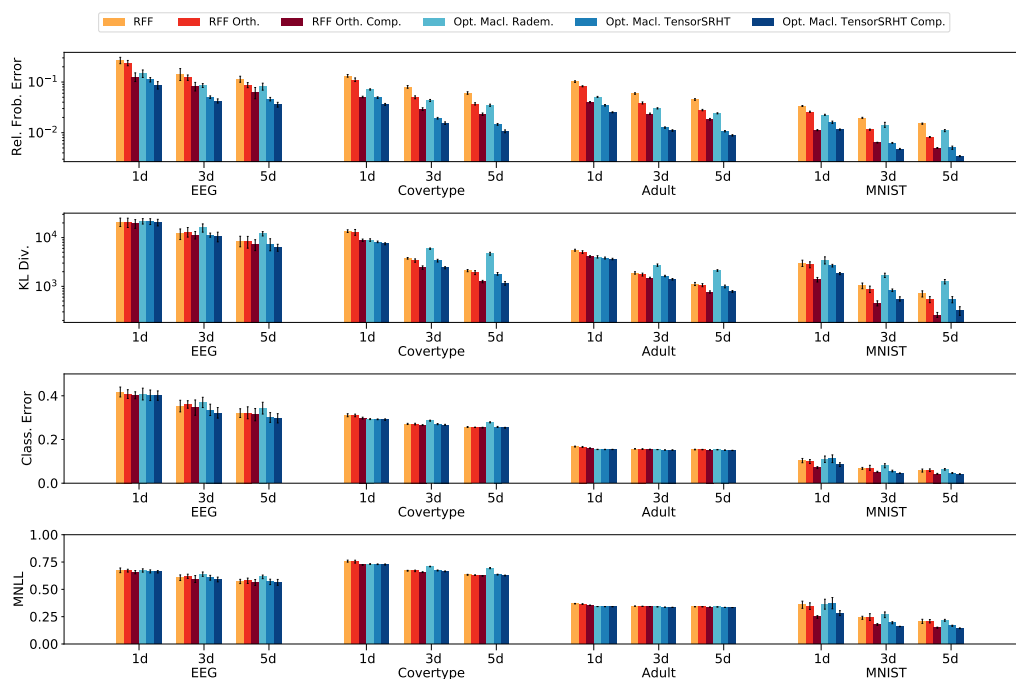


Figure D.5: Additional results of the experiments in Section 5.4.3 on approximate GP classification with a Gaussian kernel. Lower values are better for all the metrics. For each dataset, we show the number of random features $D \in \{1d, 3d, 5d\}$ used in each method on the horizontal axis, with d being the input dimensionality of the dataset. We put the legend labels and the bars in the same order.

Appendix E

Appendix of Chapter 6

E.1 Proof of Theorem 6.2.2

We now derive the limiting kernel when the optical random feature map is given by $\Phi_m(\mathbf{x}) = \frac{1}{\sqrt{D}}|\mathbf{U}\mathbf{x}|^m$ for some $m = 2s$ and $s \in \mathbb{N}$. Thus, we work out $k_m(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\Phi_m(\mathbf{x})^\top \Phi_m(\mathbf{y})]$ and obtain:

$$\begin{aligned} k_m(\mathbf{x}, \mathbf{y}) &= \int |\mathbf{x}^\top \mathbf{u}|^m |\mathbf{y}^\top \mathbf{u}|^m p(\mathbf{u}) d\mathbf{u} \\ &= \int \|\mathbf{x}\|^m \|\mathbf{y}\|^m |u_1|^m |u_1 \cos \theta + u_2 \sin \theta|^m p(\mathbf{u}) d\mathbf{u} \\ &= \|\mathbf{x}\|^m \|\mathbf{y}\|^m \int |u_1|^m |u_1 \cos \theta + u_2 \sin \theta|^m p(\mathbf{u}) d\mathbf{u} \end{aligned}$$

with $p(\mathbf{u}) = \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I})$.

Now we work out the term $A = |u_1 \cos \theta + u_2 \sin \theta|^{2s}$ (with $m = 2s$) as we focus on even powers:

$$\begin{aligned} A &= (\bar{u}_1 \cos \theta + \bar{u}_2 \sin \theta)^s (u_1 \cos \theta + u_2 \sin \theta)^s \\ &= \left(\sum_{j=0}^s \binom{s}{j} (\bar{u}_1 \cos \theta)^j (\bar{u}_2 \sin \theta)^{s-j} \right) \left(\sum_{i=0}^s \binom{s}{i} (u_1 \cos \theta)^i (u_2 \sin \theta)^{s-i} \right) \\ &= \sum_i \sum_j \binom{s}{j} \binom{s}{i} \bar{u}_1^j u_1^i \bar{u}_2^{s-j} u_2^{s-i} (\cos \theta)^{i+j} (\sin \theta)^{2s-i-j} \end{aligned}$$

One can notice that $\mathbb{E}[A] = \sum_i \sum_j \binom{s}{j} \binom{s}{i} \mathbb{E}[\bar{u}_1^j u_1^i \bar{u}_2^{s-j} u_2^{s-i}]$ has its cross terms equal to 0: when $i \neq j$, the expectation inside the sum is equal to zero because u_1 and u_2 are i.i.d (their correlation is then equal to 0) or we can see this by rotational invariance (any complex random variable with a phase uniform between 0 and 2π

has a mean equal to zero). Therefore, we obtain:

$$\mathbb{E}[A] = \sum_{i=0}^s \binom{s}{i}^2 |u_1|^{2i} |\cos \theta|^{2i} |u_2|^{2(s-i)} |\sin \theta|^{2(s-i)}$$

Using the computation of $\mathbb{E}[A]$, we can therefore deduce the analytical formula for the kernel:

$$\begin{aligned} \frac{k_{2s}(\mathbf{x}, \mathbf{y})}{\|\mathbf{x}\|^m \|\mathbf{y}\|^m} &= \int |u_1|^{2s} \sum_{i=0}^s \binom{s}{i}^2 |u_1|^{2i} |\cos \theta|^{2i} |u_2|^{2(s-i)} |\sin \theta|^{2(s-i)} p(\mathbf{u}) d\mathbf{u} \\ &= \sum_{i=0}^s \binom{s}{i}^2 \mathbb{E} \left[|u_1|^{2(s+i)} |\cos \theta|^{2i} \right] \mathbb{E} \left[|u_2|^{2(s-i)} |\sin \theta|^{2(s-i)} \right] \\ &= \sum_{i=0}^s \binom{s}{i}^2 |\cos \theta|^{2i} 2^{s+i} \sigma^{*2(s+i)} \Gamma(s+i+1) |\sin \theta|^{2(s-i)} 2^{s-i} \sigma^{*2(s-i)} \Gamma(s-i+1) \\ &= \sum_{i=0}^s \binom{s}{i}^2 2^{2s} \sigma^{*4s} |\cos \theta|^{2i} |\sin \theta|^{2(s-i)} \Gamma(s+i+1) \Gamma(s-i+1), \end{aligned}$$

where $\Gamma(\cdot)$ is the Gamma function and $\sigma^{*2} = \frac{1}{2}\sigma^2$ being the variance of the real and imaginary parts of \mathbf{u} . W.l.o.g. we will use $\sigma^2 = 1$ and thus $\sigma^{*2} = 0.5$ in the following.

We can simplify this formula, by noticing that $\sin^2 \theta = 1 - \cos^2 \theta$. The latter formula then becomes:

$$\frac{k_{2s}(\mathbf{x}, \mathbf{y})}{\|\mathbf{x}\|^m \|\mathbf{y}\|^m} = \sum_{i=0}^s \binom{s}{i}^2 |\cos \theta|^{2i} (1 - \cos^2 \theta)^{(s-i)} 2^{2s} \sigma^{*4s} \Gamma(s+i+1) \Gamma(s-i+1)$$

The new term can be expanded using a binomial expansion: $(1 - \cos^2 \theta)^{(s-i)} = \sum_{t=0}^{s-i} (-\cos^2 \theta)^t$, leading to

$$\frac{k_{2s}(\mathbf{x}, \mathbf{y})}{\|\mathbf{x}\|^m \|\mathbf{y}\|^m} = \sum_{i=0}^s \sum_{t=0}^{s-i} \binom{s}{i}^2 (s+i)!(s-i)! \binom{s-i}{t} (-1)^t \cos^{2(i+t)} \theta$$

Using the change of variable $a = i + t$ and keeping $i = i$, we obtain:

$$\frac{k_{2s}(\mathbf{x}, \mathbf{y})}{\|\mathbf{x}\|^m \|\mathbf{y}\|^m} = \sum_{a=0}^s \left[\sum_{i=0}^a \binom{s}{i}^2 (s+i)!(s-i)! \binom{s-i}{a-i} (-1)^{a-i} \right] \cos^{2a} \theta$$

We focus on the term between the parenthesis that we will call T_a :

$$\begin{aligned} T_a &= \sum_{i=0}^a \frac{s!^2}{(s-i)!2i!^2} (s+i)!(s-i)! \frac{(s-i)!}{(a-i)!(s-a)!} (-1)^{a-i} \\ &= \sum_{i=0}^a (s!)^2 \frac{(s+i)!}{(i!)^2 (a-i)!(s-a)!} (-1)^{a-i} \\ &= (s!)^2 \binom{s}{a} (-1)^a \sum_{i=0}^a \binom{a}{i} \binom{s+i}{i} (-1)^i \end{aligned}$$

Using the upper negation ($\binom{a}{b} = (-1)^b \binom{b-a-1}{b}$) so here $\binom{s+i}{i} = \binom{i-i-s-1}{i} (-1)^i = \binom{-s-1}{i} (-1)^i$ leads to

$$T_a = (s!)^2 \binom{s}{a} (-1)^a \sum_{i=0}^a \binom{a}{i} \binom{-s-1}{i} \quad (\text{E.1})$$

$$= (s!)^2 \binom{s}{a} (-1)^a \sum_{i=0}^a \binom{a}{i} \binom{-s-1}{-s-1-i} \quad (\text{E.2})$$

Now we can use the Vandermonde identity $\sum_{i=0}^n \binom{a}{i} \binom{b}{n-i} = \binom{a+b}{n}$, yielding:

$$T_a = (s!)^2 \binom{s}{a} (-1)^a \binom{a-s-1}{-s-1} = (s!)^2 \binom{s}{a} (-1)^a \binom{a-s-1}{a} \quad (\text{E.3})$$

$$= (s!)^2 \binom{s}{a} \binom{s}{a} = (s!)^2 \binom{s}{a}^2 \quad (\text{E.4})$$

where in the last line we used again the upper negation formula. This leads to the desired result for $m = 2s$

$$k_{2s}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x}\|^m \|\mathbf{y}\|^m \sum_{i=0}^s (s!)^2 \binom{s}{i}^2 \cos^{2i} \theta$$

Bibliography

- M. Abramowitz and I. A. Stegun. *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Dover Publications, Inc., USA, 1974.
- T. D. Ahle, M. Kapralov, J. B. T. Knudsen, R. Pagh, A. Velingker, D. P. Woodruff, and A. Zandieh. Oblivious sketching of high-degree polynomial kernels. In *Proceedings of the Thirty-First Annual ACM-SIAM Symposium on Discrete Algorithms*, page 141–160. Society for Industrial and Applied Mathematics, 2020.
- S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2425–2433, 2015.
- H. Aschard. A perspective on interaction effects in genetic association studies. *Genetic epidemiology*, 40(8):678–688, 2016.
- H. Avron, H. L. Nguyen, and D. P. Woodruff. Subspace embeddings for the polynomial kernel. In *Advances in Neural Information Processing Systems 27*, pages 2258–2266. Curran Associates, Inc., 2014.
- H. Avron, M. Kapralov, C. Musco, C. Musco, A. Velingker, and A. Zandieh. Random Fourier features for kernel ridge regression: Approximation bounds and statistical guarantees. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 253–262. PMLR, 2017.
- N. Banar, W. Daelemans, and M. Kestemont. Character-level transformer-based neural machine translation. In *Proceedings of the 4th International Conference on Natural Language Processing and Information Retrieval*, page 149–156. Association for Computing Machinery, 2020.
- M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton. Demystifying MMD GANs. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=r11U0zWCW>.
- M. Blondel, M. Ishihata, A. Fujino, and N. Ueda. Polynomial networks and factorization machines: New insights and efficient training algorithms. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, volume 48, pages 850–858. PMLR, 2016.

- R. Boloix-Tortosa, J. J. Murillo-Fuentes, F. J. Payán-Somet, and F. Pérez-Cruz. Complex gaussian processes for regression. *IEEE Transactions on Neural Networks and Learning Systems*, 29(11):5499–5511, 2018.
- J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- T. Brown et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- Y.-W. Chang, C.-J. Hsieh, K.-W. Chang, M. Ringgaard, and C.-J. Lin. Training and testing low-degree polynomial data mappings via linear SVM. *Journal of Machine Learning Research*, 11(4), 2010.
- M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. In *Proceedings of the 29th International Colloquium on Automata, Languages and Programming*, page 693–703. Springer-Verlag, 2002.
- K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *British Machine Vision Conference*. BMVA Press, 2014.
- K. Choromanski, M. Rowland, and A. Weller. The unreasonable effectiveness of structured random orthogonal embeddings. In *Advances in Neural Information Processing Systems 31*, pages 218–227. Curran Associates Inc., 2017.
- K. Choromanski, M. Rowland, T. Sarlos, V. Sindhvani, R. Turner, and A. Weller. The geometry of random features. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84, pages 1–9. PMLR, 2018.
- K. Choromanski et al. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Ua6zuk0WRH>.
- A. Cotter, J. Keshet, and N. Srebro. Explicit approximations of the gaussian kernel. *CoRR*, abs/1109.4603, 2011.
- N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 886–893. IEEE Computer Society, 2005.

- J. Devlin, M. Chang, K. Lee, and K. Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186. Association for Computational Linguistics, 2019.
- J. Dong, S. Gigan, F. Krzakala, and G. Wainrib. Scaling up echo-state networks with multiple light scattering. In *2018 IEEE Statistical Signal Processing Workshop (SSP)*, pages 448–452. IEEE, 2018.
- J. Dong, M. Rafayelyan, F. Krzakala, and S. Gigan. Optical reservoir computing using multiple light scattering for chaotic systems prediction. *IEEE Journal of Selected Topics in Quantum Electronics*, 26(1):1–12, 2019.
- D. Dua and C. Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- B. J. Fino and V. R. Algazi. Unified matrix treatment of the fast walsh-hadamard transform. *IEEE Transactions on Computers*, 25(11):1142–1146, 1976.
- C. A. Floudas and P. M. Pardalos, editors. *Encyclopedia of Optimization, Second Edition*. Springer, 2009.
- A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Multi-modal compact bilinear pooling for visual question answering and visual grounding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 457–468. Association for Computational Linguistics, 2016.
- Y. Gao, O. Beijbom, N. Zhang, and T. Darrell. Compact bilinear pooling. *Proceedings of the 2016 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 317–326, 2016.
- D. Garreau, W. Jitkrittum, and M. Kanagawa. Large sample analysis of the median heuristic. *arXiv preprint arXiv:1707.07269*, 2017.
- M. G. Genton. Classes of kernels for machine learning: A statistics perspective. *Journal of Machine Learning Research*, 2:299–312, 2002.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773, 2012.
- S. Gupta, R. Gribonval, L. Daudet, and I. Dokmanic. Don’t take it lightly: Phasing optical random projections with unknown operators. In *Advances in Neural Information Processing Systems*, volume 32, pages 14826–14836. Curran Associates, Inc., 2019.

- R. Hamid, Y. Xiao, A. Gittens, and D. DeCoste. Compact random feature maps. In *Proceedings of the 31th International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 19–27. PMLR, 2014.
- I. Han, A. Zandieh, and H. Avron. Random gegenbauer features for scalable kernel methods. *arXiv preprint arXiv:2202.03474*, 2022.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- J. Hensman, N. Durrande, and A. Solin. Variational Fourier features for Gaussian processes. *Journal of Machine Learning Research*, 18(151):1–52, 2018.
- P. Kar and H. Karnick. Random feature maps for dot product kernels. In *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*, volume 22 of *JMLR Proceedings*, pages 583–591. JMLR, 2012.
- N. Keriven, D. Garreau, and I. Poli. Newma: a new method for scalable model-free online change-point detection. *arXiv preprint arXiv:1805.08061*, 2018.
- A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- T. Lin, A. RoyChowdhury, and S. Maji. Bilinear convolutional neural networks for fine-grained visual recognition. *IEEE Trans. Pattern Anal. Mach. Intelligence*, 40(6):1309–1322, 2018.
- F. Liu, X. Huang, Y. Chen, and J. A. K. Suykens. Random features for kernel approximation: A survey in algorithms, theory, and beyond. *CoRR*, abs/2004.11154, 2020.
- Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.
- A. Liutkus, D. Martina, S. Popoff, G. Chardon, O. Katz, G. Lerosey, S. Gigan, L. Daudet, and I. Carron. Imaging with nature: Compressive imaging using a multiply scattering medium. *Scientific reports*, 4:5552, 2014.

- D. G. Lowe. Object recognition from local scale-invariant features. In *Proceedings of the International Conference on Computer Vision*, pages 1150–1157. IEEE Computer Society, 1999.
- A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150. Association for Computational Linguistics, 2011.
- D. Milios, R. Camoriano, P. Michiardi, L. Rosasco, and M. Filippone. Dirichlet-based gaussian processes for large-scale calibrated classification. In *Advances in Neural Information Processing Systems 31*, pages 6008–6018. Curran Associates, Inc., 2018.
- F. D. Neeser and J. L. Massey. Proper complex random processes with applications to information theory. *IEEE transactions on information theory*, 39(4):1293–1302, 1993.
- R. Ohana, J. Wacker, J. Dong, S. Marmin, F. Krzakala, M. Filippone, and L. Daudet. Kernel computations from large-scale random features obtained by optical processing units. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 9294–9298, 2020.
- K. I. Park. *Fundamentals of Probability and Stochastic Processes with Applications to Communications*. Springer, 1st edition, 2018.
- A. Paszke et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, pages 8026–8037. Curran Associates, Inc., 2019.
- F. Pedregosa, , et al. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- J. Pennington, F. X. X. Yu, and S. Kumar. Spherical random features for polynomial kernels. In *Advances in Neural Information Processing Systems 28*, pages 1846–1854. Curran Associates, Inc., 2015.
- N. Pham and R. Pagh. Fast and scalable polynomial kernels via explicit feature maps. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 239–247. Association for Computing Machinery, 2013.
- A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems 20*, pages 1177–1184. Curran Associates Inc., 2007.

- C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- C. E. Rasmussen and J. Quiñonero Candela. Healing the relevance vector machine through augmentation. In *Proceedings of the 22nd International Conference on Machine Learning*, Proceedings of Machine Learning Research, page 689–696. PMLR, 2005.
- A. S. Rawat, J. Chen, F. Yu, A. T. Suresh, and S. Kumar. *Sampled Softmax with Random Fourier Features*, volume 32, pages 13834–13844. Curran Associates Inc., 2019.
- S. Rendle. Factorization machines. In *Proceedings of the 2010 IEEE International Conference on Data Mining*, pages 995–1000, 2010.
- A. Rudi and L. Rosasco. Generalization properties of learning with random features. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- O. Russakovsky et al. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- A. Saade, F. Caltagirone, I. Carron, L. Daudet, A. Drémeau, S. Gigan, and F. Krzakala. Random projections through multiple optical scattering: Approximating kernels at the speed of light. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6215–6219, March 2016.
- I. J. Schoenberg. Metric spaces and completely monotone functions. *Annals of Mathematics*, 39(4):811–841, 1938.
- B. Scholkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- A. Smola, Z. Óvári, and R. C. Williamson. Regularization with dot-product kernels. In *Advances in Neural Information Processing Systems 13*, pages 308–314. Curran Associates, Inc., 2000.
- M. L. Stein. *Interpolation of spatial data*. Springer-Verlag, 1999.
- I. Steinwart and A. Christmann. *Support Vector Machines*. Information science and statistics. Springer, 2008.
- E. Strubell, A. Ganesh, and A. McCallum. Energy and policy considerations for deep learning in NLP. In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pages 3645–3650. Association for Computational Linguistics, 2019.

- D. J. Sutherland and J. Schneider. On the error of random fourier features. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pages 862–871. AUAI Press, 2015.
- C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- Y. Tay, M. Dehghani, S. Abnar, Y. Shen, D. Bahri, P. Pham, J. Rao, L. Yang, S. Ruder, and D. Metzler. Long range arena : A benchmark for efficient transformers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=qVyeW-grC2k>.
- M. E. Tipping. The relevance vector machine. In *Advances in Neural Information Processing Systems 12*, page 652–658. Curran Associates, Inc., 1999.
- M. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, volume 5 of *JMLR Proceedings*, pages 567–574. JMLR, 2009.
- L. N. Trefethen and D. Bau. *Numerical Linear Algebra*. SIAM, 1997.
- J. A. Tropp. Improved analysis of the subsampled randomized hadamard transform. *Advances in Adaptive Data Analysis*, 3(1-2):115–126, 2011.
- A. V. Uzilov, J. M. Keegan, and D. H. Mathews. Detection of non-coding rnas on the basis of predicted secondary structure formation free energy change. *BMC Bioinformatics*, 7:173, 2006.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pages 6000–6010. Curran Associates, Inc., 2017.
- R. Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.
- J. Wacker and M. Filippone. Local random feature approximations of the gaussian kernel. In *Proceedings of the 26th International Conference on Knowledge-Based and Intelligent Information & Engineering Systems*, Procedia Computer Science. Elsevier, 2021.
- J. Wacker, M. Ladeira, and J. E. V. Nascimento. Transfer learning for brain tumor segmentation. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries - 6th International Workshop, BrainLes 2020, Held in Conjunction with MICCAI 2020.*, volume 12658 of *Lecture Notes in Computer Science*, pages 241–251. Springer, 2020.

- J. Wacker, M. Kanagawa, and M. Filippone. Improved random features for dot product kernels. *arXiv preprint arXiv:2201.08712*, 2022a.
- J. Wacker, R. Ohana, and M. Filippone. Complex-to-real random features for polynomial kernels. *arXiv preprint arXiv:2202.02031*, 2022b.
- G. Wahba. *Spline Models for Observational Data*. SIAM, 1990.
- K. Q. Weinberger, A. Dasgupta, J. Langford, A. J. Smola, and J. Attenberg. Feature hashing for large scale multitask learning. In A. P. Danyluk, L. Bottou, and M. L. Littman, editors, *Proceedings of the 26th Annual International Conference on Machine Learning*, volume 382, pages 1113–1120. ACM, 2009.
- O. Weissbrod, D. Geiger, and S. Rosset. Multikernel linear mixed models for complex phenotype prediction. *Genome Research*, 26(7):969–979, 2016.
- P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-ucsd birds 200. Technical report, Caltech, 2010. URL <http://www.vision.caltech.edu/visipedia/CUB-200.html>.
- C. K. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems 13*, pages 682–688. Curran Associates, Inc., 2000.
- D. P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, 10(1-2):1–157, 2014.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- Y. Xiong, Z. Zeng, R. Chakraborty, M. Tan, G. Fung, Y. Li, and V. Singh. Nyströmformer: A nyström-based algorithm for approximating self-attention. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(16):14138–14148, 2021.
- H. Yamada and Y. Matsumoto. Statistical dependency analysis with support vector machines. In *Proceedings of the Eighth International Conference on Parsing Technologies*, pages 195–206, 2003.
- Z. Yang, M. Moczulski, M. Denil, N. D. Freitas, A. Smola, L. Song, and Z. Wang. Deep fried convnets. *Proceedings of the 2015 IEEE International Conference on Computer Vision*, pages 1476–1483, 2015.
- F. X. Yu, A. T. Suresh, K. Choromanski, D. Holtmann-Rice, and S. Kumar. Orthogonal random features. In *Advances in Neural Information Processing Systems*, volume 30, page 1983–1991. Curran Associates Inc., 2016.

- M. Zaheer et al. Big bird: Transformers for longer sequences. In *Advances in Neural Information Processing Systems*, volume 33, pages 17283–17297. Curran Associates, Inc., 2020.