

Gaussian Processes

Maurizio Filippone

EURECOM, Sophia Antipolis, France

February, 9th, 2017

Gaussian Processes for Machine Learning

Carl E. Rasmussen and Christopher K. I. Williams

Pattern Recognition and Machine Learning

C. Bishop

- Motivation – Examples
- Introduction to Gaussian Processes
 - Weight space view
 - Function space view
- Challenges
- Modern Gaussian Processes

Motivation

Motivation

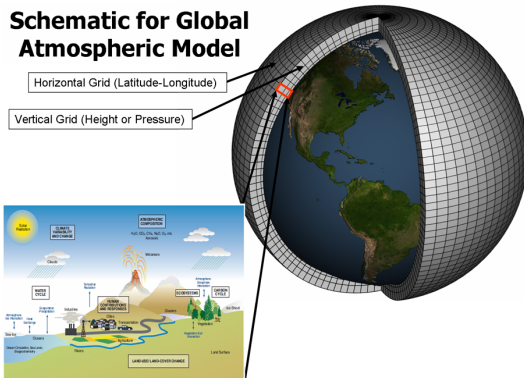
Quantification of Uncertainty with Expensive Models

- Climate modeling

Schematic for Global Atmospheric Model

Horizontal Grid (Latitude-Longitude)

Vertical Grid (Height or Pressure)

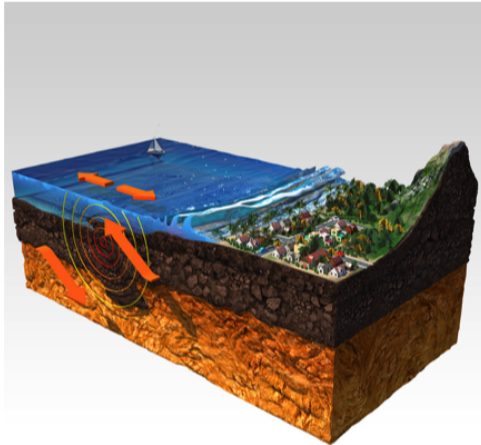


Kennedy and O'Hagan, *RSS-B*, 2001

Motivation

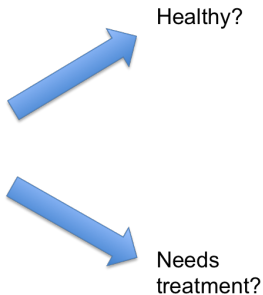
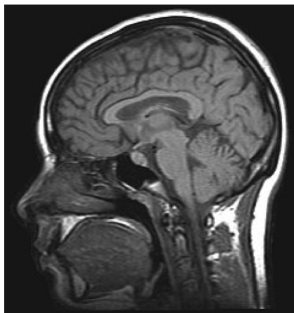
Quantification of Uncertainty with Expensive Models

- Earthquake modeling



Kennedy and O'Hagan, *RSS-B*, 2001

- Classification of neurodegenerative diseases

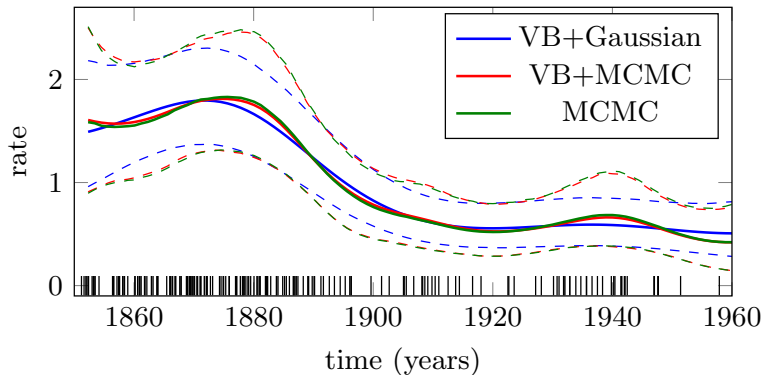


Filippone et al., *AoAS*, 2012

Motivation

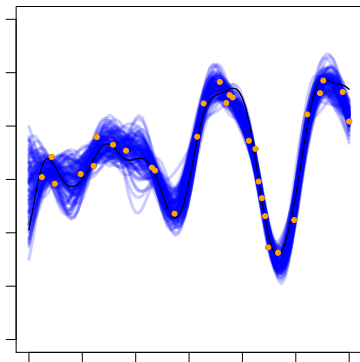
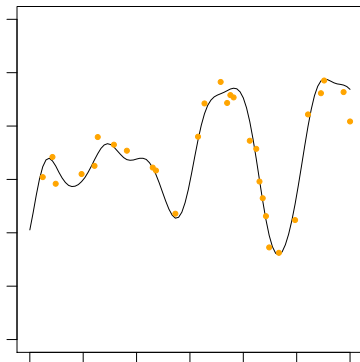
Quantification of Uncertainty with No Models

- Coal mining disaster data



Hensman, Matthews, Filippone, Ghahramani, *NIPS*, 2015

- Regression example



A Unified Framework

A model might be expensive to simulate/inaccurate

- Emulate model/discrepancy using a surrogate

A Unified Framework

A model might be expensive to simulate/inaccurate

- Emulate model/discrepancy using a surrogate

A model might not even be available

- Replace it with a flexible model

A model might be expensive to simulate/inaccurate

- Emulate model/discrepancy using a surrogate

A model might not even be available

- Replace it with a flexible model

Gaussian processes for Accurate Quantification of Uncertainty

Gaussian Processes

Gaussian Processes can be explained in two ways

- Weight Space View
 - Bayesian linear regression with infinite basis functions
- Function Space View
 - Defined as priors over functions

Gaussian Processes can be explained in two ways

- **Weight Space View**
 - **Bayesian linear regression with infinite basis functions**
- **Function Space View**
 - Defined as priors over functions

- Modeling observations as noisy realizations of a linear combination of the features:

$$p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I})$$

- Modeling observations as noisy realizations of a linear combination of the features:

$$p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I})$$

- Gaussian prior over model parameters:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{S})$$

- Bayes rule:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

Bayesian Linear Regression

- Bayes rule:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

- **Posterior density:** $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$
 - Distribution over parameters *after* observing data

- Bayes rule:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

- **Posterior density:** $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$
 - Distribution over parameters *after* observing data
- **Likelihood :** $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$
 - Measure of “fitness”

- Bayes rule:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

- **Posterior density:** $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$
 - Distribution over parameters *after* observing data
- **Likelihood :** $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$
 - Measure of “fitness”
- **Prior density:** $p(\mathbf{w})$
 - Anything we know about parameters *before* we see any data

- Bayes rule:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

- **Posterior density:** $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$
 - Distribution over parameters *after* observing data
- **Likelihood :** $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$
 - Measure of “fitness”
- **Prior density:** $p(\mathbf{w})$
 - Anything we know about parameters *before* we see any data
- **Marginal likelihood:** $p(\mathbf{y}|\mathbf{X})$
 - It is a normalization constant – ensures $\int p(\mathbf{w}|\mathbf{X}, \mathbf{y}) d\mathbf{w} = 1$.

- Ignoring normalizing constants, the posterior is:

$$\begin{aligned} p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) &\propto \exp \left\{ -\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu}) \right\} \\ &= \exp \left\{ -\frac{1}{2}(\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right\} \\ &\propto \exp \left\{ -\frac{1}{2}(\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right\} \end{aligned}$$

- Ignoring non- \mathbf{w} terms, the prior multiplied by the likelihood is:

$$\begin{aligned} & p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \sigma^2) \\ & \propto \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})\right\} \exp\left\{-\frac{1}{2}\mathbf{w}^T\mathbf{S}^{-1}\mathbf{w}\right\} \\ & \propto \exp\left\{-\frac{1}{2}\left(\mathbf{w}^T\left[\frac{1}{\sigma^2}\mathbf{X}^T\mathbf{X} + \mathbf{S}^{-1}\right]\mathbf{w} - \frac{2}{\sigma^2}\mathbf{w}^T\mathbf{X}^T\mathbf{y}\right)\right\} \end{aligned}$$

- Posterior (from previous slide):

$$\propto \exp\left\{-\frac{1}{2}(\mathbf{w}^T\mathbf{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\mathbf{\Sigma}^{-1}\boldsymbol{\mu})\right\}$$

- Equate individual terms on each side.
- Covariance:

$$\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w} = \mathbf{w}^T \left[\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{S}^{-1} \right] \mathbf{w}$$
$$\boldsymbol{\Sigma} = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{S}^{-1} \right)^{-1}$$

- Mean:

$$2\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} = \frac{2}{\sigma^2} \mathbf{w}^T \mathbf{X}^T \mathbf{y}$$
$$\boldsymbol{\mu} = \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^T \mathbf{y}$$

- Posterior **must be Gaussian**

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Covariance:

$$\boldsymbol{\Sigma} = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{S}^{-1} \right)^{-1}$$

- Mean:

$$\boldsymbol{\mu} = \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^T \mathbf{y}$$

- Predictions – same tedious exercise as before:

$$p(y_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*, \sigma^2) = \mathcal{N}(\mathbf{x}_*^T \boldsymbol{\mu}, \sigma^2 + \mathbf{x}_*^T \boldsymbol{\Sigma} \mathbf{x}_*)$$

Introducing basis functions

- Imagine transforming the inputs using a set of D functions

$$\mathbf{x} \rightarrow \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_D(\mathbf{x}))^\top$$

- The functions $\phi_1(\mathbf{x})$ are also known as *basis functions*
- Define:

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \dots & \phi_D(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_N) & \dots & \phi_D(\mathbf{x}_N) \end{bmatrix}$$

Introducing basis functions

- Applying Bayesian Linear Regression on the transformed features gives

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Covariance:

$$\boldsymbol{\Sigma} = \left(\frac{1}{\sigma^2} \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \mathbf{S}^{-1} \right)^{-1}$$

- Mean:

$$\boldsymbol{\mu} = \frac{1}{\sigma^2} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{y}$$

- Predictions:

$$p(y_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*, \sigma^2) = \mathcal{N}(\boldsymbol{\phi}_*^T \boldsymbol{\mu}, \sigma^2 + \boldsymbol{\phi}_*^T \boldsymbol{\Sigma} \boldsymbol{\phi}_*)$$

- Linear models require specifying a set of basis functions
 - Polynomials, Trigonometric, ...??

- Linear models require specifying a set of basis functions
 - Polynomials, Trigonometric, ...??
- Can we use Bayesian inference to let data tell this to us?

- Linear models require specifying a set of basis functions
 - Polynomials, Trigonometric, ...??
- Can we use Bayesian inference to let data tell this to us?
- Gaussian Processes work implicitly with an infinite set of basis functions and learn a probabilistic combination of these

Bayesian Linear Regression as a Kernel Machine

- We are going to show that predictions can be expressed exclusively in terms of scalar products as follows

$$k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\psi}(\mathbf{x})^\top \boldsymbol{\psi}(\mathbf{x}')$$

- This allows us to work with either $k(\cdot, \cdot)$ or $\boldsymbol{\psi}(\cdot)$
- Why is this useful??

Bayesian Linear Regression as a Kernel Machine

- Working with $\psi(\cdot)$ costs $O(D^2)$ storage, $O(D^3)$ time
- Working with $k(\cdot, \cdot)$ costs $O(N^2)$ storage, $O(N^3)$ time

Bayesian Linear Regression as a Kernel Machine

- Working with $\psi(\cdot)$ costs $O(D^2)$ storage, $O(D^3)$ time
- Working with $k(\cdot, \cdot)$ costs $O(N^2)$ storage, $O(N^3)$ time
- Pick the one that makes computations faster ... or

Bayesian Linear Regression as a Kernel Machine

- Working with $\psi(\cdot)$ costs $O(D^2)$ storage, $O(D^3)$ time
- Working with $k(\cdot, \cdot)$ costs $O(N^2)$ storage, $O(N^3)$ time
- Pick the one that makes computations faster ... or
- What if we could pick $k(\cdot, \cdot)$ so that $\psi(\cdot)$ is infinite dimensional?

- It is possible to show that for

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2}\right)$$

there exists a corresponding $\psi(\cdot)$ that is infinite dimensional!!!

- There are other kernels satisfying this property

Kernels

Proof that the Gaussian kernel induces an infinite dimensional $\psi(\cdot)$

- For simplicity consider one dimensional inputs x, y
- Expand the Gaussian kernel $k(x, y)$ as

$$\exp\left(-\frac{(x-y)^2}{2}\right) = \exp\left(-\frac{x^2}{2}\right) \exp\left(-\frac{y^2}{2}\right) \exp(xy)$$

- Focusing on the last term and applying the Taylor expansion of the $\exp(\cdot)$ function

$$\exp(xy) = 1 + (xy) + \frac{(xy)^2}{2!} + \frac{(xy)^3}{3!} + \frac{(xy)^4}{4!} + \dots$$

Kernels

Proof that the Gaussian kernel induces an infinite dimensional $\psi(\cdot)$

- Define the infinite dimensional mapping

$$\psi(x) = \exp\left(-\frac{x^2}{2}\right) \left(1, x, \frac{x^2}{\sqrt{2!}}, \frac{x^3}{\sqrt{3!}}, \frac{x^4}{\sqrt{4!}}, \dots\right)^\top$$

- It is easy to verify that

$$k(x, y) = \exp\left(-\frac{(x-y)^2}{2}\right) = \psi(x)^\top \psi(y)$$

Bayesian Linear Regression as a Kernel Machine

Proof

- To show that Bayesian Linear Regression can be formulated through scalar products only, we need Woodbury identity:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

- Do not memorize this!

- Woodbury identity:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

- We can rewrite:

$$\begin{aligned}\Sigma &= \left(\frac{1}{\sigma^2} \Phi^T \Phi + \mathbf{S}^{-1} \right)^{-1} \\ &= \mathbf{S} - \mathbf{S} \Phi^T \left(\sigma^2 \mathbf{I} + \Phi \mathbf{S} \Phi^T \right)^{-1} \Phi \mathbf{S}\end{aligned}$$

- We set $A = \mathbf{S}$, $U = V^T = \Phi^T$, and $C = \frac{1}{\sigma^2} \mathbf{I}$

Bayesian Linear Regression as a Kernel Machine

Proof

- Mean and variance of the predictions:

$$p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*, \sigma^2) = \mathcal{N}(\phi_*^\top \boldsymbol{\mu}, \sigma^2 + \phi_*^\top \boldsymbol{\Sigma} \phi_*)$$

- Rewrite the variance:

$$\sigma^2 + \phi_*^\top \boldsymbol{\Sigma} \phi_* =$$

$$\sigma^2 + \phi_*^\top \mathbf{S} \phi_* - \phi_*^\top \mathbf{S} \boldsymbol{\Phi}^\top \left(\sigma^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top \right)^{-1} \boldsymbol{\Phi} \mathbf{S} \phi_*$$

... continued

Bayesian Linear Regression as a Kernel Machine

Proof

- Mean and variance of the predictions:

$$p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*, \sigma^2) = \mathcal{N}(\phi_*^\top \boldsymbol{\mu}, \sigma^2 + \phi_*^\top \boldsymbol{\Sigma} \phi_*)$$

- Rewrite the variance:

$$\begin{aligned} \sigma^2 + \phi_*^\top \mathbf{S} \phi_* - \phi_*^\top \mathbf{S} \boldsymbol{\Phi}^\top (\sigma^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top)^{-1} \boldsymbol{\Phi} \mathbf{S} \phi_* &= \\ \sigma^2 + k_{**} - \mathbf{k}_*^\top (\sigma^2 \mathbf{I} + \mathbf{K})^{-1} \mathbf{k}_* & \end{aligned}$$

- Where the mapping defining the kernel is

$$\boldsymbol{\psi}(\mathbf{x}) = \mathbf{S}^{1/2} \boldsymbol{\phi}(\mathbf{x})$$

and

$$\begin{aligned} k_{**} &= k(\mathbf{x}_*, \mathbf{x}_*) = \boldsymbol{\psi}(\mathbf{x}_*)^\top \boldsymbol{\psi}(\mathbf{x}_*) \\ (\mathbf{k}_*)_i &= k(\mathbf{x}_*, \mathbf{x}_i) = \boldsymbol{\psi}(\mathbf{x}_*)^\top \boldsymbol{\psi}(\mathbf{x}_i) \\ (\mathbf{K})_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\psi}(\mathbf{x}_i)^\top \boldsymbol{\psi}(\mathbf{x}_j) \end{aligned}$$

Bayesian Linear Regression as a Kernel Machine

Proof

- Mean and variance of the predictions:

$$p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*, \sigma^2) = \mathcal{N}(\phi_*^\top \boldsymbol{\mu}, \sigma^2 + \phi_*^\top \boldsymbol{\Sigma} \phi_*)$$

- Rewrite the mean:

$$\begin{aligned}\phi_*^\top \boldsymbol{\mu} &= \frac{1}{\sigma^2} \phi_*^\top \boldsymbol{\Sigma} \boldsymbol{\Phi}^\top \mathbf{y} \\ &= \frac{1}{\sigma^2} \phi_*^\top \left(\mathbf{S} - \mathbf{S} \boldsymbol{\Phi}^\top \left(\sigma^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top \right)^{-1} \boldsymbol{\Phi} \mathbf{S} \right) \boldsymbol{\Phi}^\top \mathbf{y} \\ &= \frac{1}{\sigma^2} \phi_*^\top \mathbf{S} \boldsymbol{\Phi}^\top \left(\mathbf{I} - \left(\sigma^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top \right)^{-1} \boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top \right) \mathbf{y} \\ &= \frac{1}{\sigma^2} \phi_*^\top \mathbf{S} \boldsymbol{\Phi}^\top \left(\mathbf{I} - \left(\mathbf{I} + \frac{\boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top}{\sigma^2} \right)^{-1} \frac{\boldsymbol{\Phi} \mathbf{S} \boldsymbol{\Phi}^\top}{\sigma^2} \right) \mathbf{y}\end{aligned}$$

... continued

Bayesian Linear Regression as a Kernel Machine

Proof

- Define $\mathbf{H} = \frac{\boldsymbol{\Phi}\mathbf{S}\boldsymbol{\Phi}^\top}{\sigma^2}$
- The term in the parenthesis

$$\left(\mathbf{I} - \left(\mathbf{I} + \frac{\boldsymbol{\Phi}\mathbf{S}\boldsymbol{\Phi}^\top}{\sigma^2} \right)^{-1} \frac{\boldsymbol{\Phi}\mathbf{S}\boldsymbol{\Phi}^\top}{\sigma^2} \right)$$

becomes

$$\left(\mathbf{I} - (\mathbf{I} + \mathbf{H})^{-1} \mathbf{H} \right) = \mathbf{I} - (\mathbf{H}^{-1} + \mathbf{I})^{-1}$$

- Using Woodbury ($A, U, V = \mathbf{I}$ and $C = \mathbf{H}^{-1}$)

$$\mathbf{I} - (\mathbf{H}^{-1} + \mathbf{I})^{-1} = (\mathbf{I} + \mathbf{H})^{-1}$$

Bayesian Linear Regression as a Kernel Machine

Proof

- Substituting into the expression of the predictive mean

$$\begin{aligned}\phi_*^\top \mu &= \frac{1}{\sigma^2} \phi_*^\top \mathbf{S} \Phi^\top \left(\mathbf{I} - \left(\mathbf{I} + \frac{\Phi \mathbf{S} \Phi^\top}{\sigma^2} \right)^{-1} \frac{\Phi \mathbf{S} \Phi^\top}{\sigma^2} \right) \mathbf{y} \\ &= \frac{1}{\sigma^2} \phi_*^\top \mathbf{S} \Phi^\top \left(\mathbf{I} + \frac{\Phi \mathbf{S} \Phi^\top}{\sigma^2} \right)^{-1} \mathbf{y} \\ &= \phi_*^\top \mathbf{S} \Phi^\top \left(\sigma^2 \mathbf{I} + \Phi \mathbf{S} \Phi^\top \right)^{-1} \mathbf{y} \\ &= \mathbf{k}_*^\top \left(\sigma^2 \mathbf{I} + \mathbf{K} \right)^{-1} \mathbf{y}\end{aligned}$$

- All definitions as in the case of the variance

$$\psi(\mathbf{x}) = \mathbf{S}^{1/2} \phi(\mathbf{x})$$

$$(\mathbf{k}_*)_i = k(\mathbf{x}_*, \mathbf{x}_i) = \psi(\mathbf{x}_*)^\top \psi(\mathbf{x}_i)$$

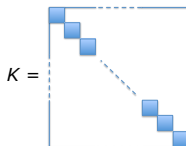
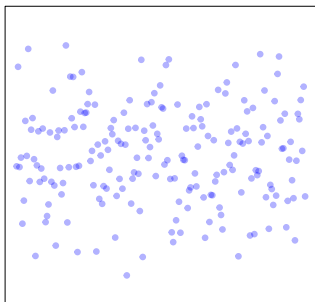
$$(\mathbf{K})_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \psi(\mathbf{x}_i)^\top \psi(\mathbf{x}_j)$$

Gaussian Processes can be explained in two ways

- Weight Space View
 - Bayesian linear regression with infinite basis functions
- **Function Space View**
 - **Defined as priors over functions**

Gaussian Processes - Priors over Functions

- Consider an infinite number of Gaussian random variables
- Think of them as indexed by the real line and as independent
- Denote them as $f(x)$



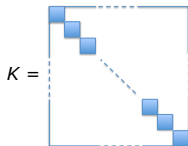
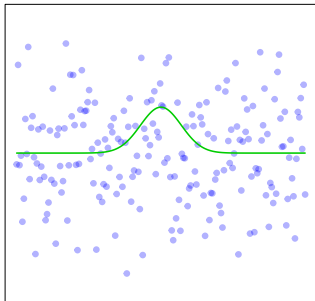
- Consider the Gaussian kernel again

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp(-\beta \|\mathbf{x} - \mathbf{x}'\|^2)$$

- We introduced some parameters for added flexibility

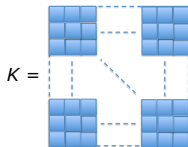
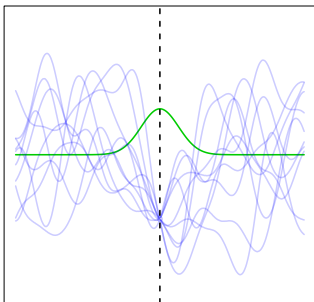
Gaussian Processes - Priors over Functions

- Impose covariance using the kernel function



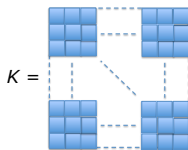
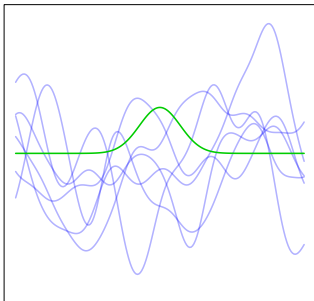
Gaussian Processes - Priors over Functions

- Draw the infinite random variables again fixing one of them (the one at $x = 0$)



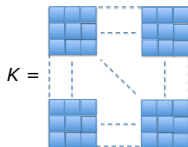
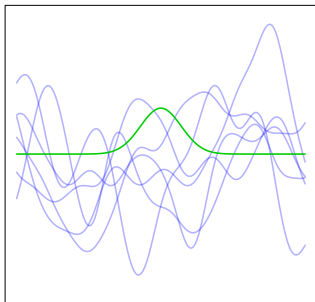
Gaussian Processes - Priors over Functions

- Draw the infinite random variables again allowing the one at $x = 0$ to be random too



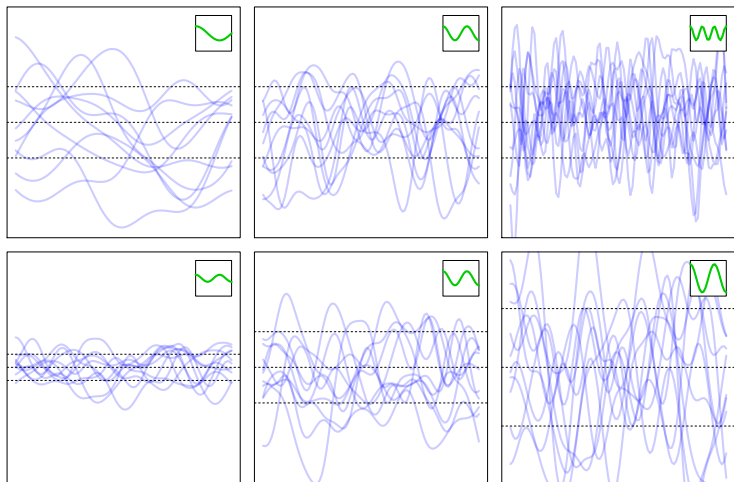
Gaussian Processes - Priors over Functions

- This can be used as a prior over functions!



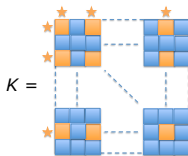
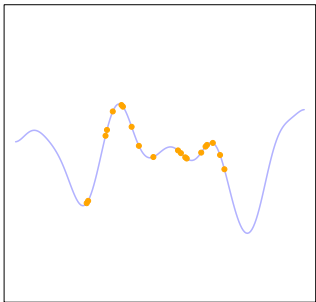
Gaussian Processes - Priors over Functions

- Infinite Gaussian random variables with parameterized and input-dependent covariance



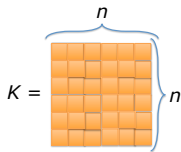
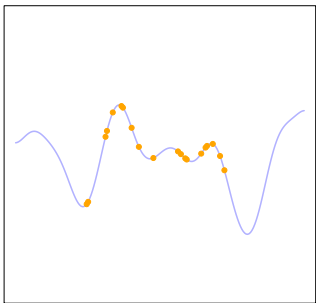
Gaussian Processes - Priors over Functions

- The distribution of N random variables $f(x_1), \dots, f(x_N)$ depends exclusively on the corresponding rows and columns of the infinite by infinite kernel matrix K



Gaussian Processes - Priors over Functions

- The distribution of N random variables $f(x_1), \dots, f(x_N)$ depends exclusively on the corresponding rows and columns of the infinite by infinite kernel matrix K



- The marginal distribution of $\mathbf{f} = (f(x_1), \dots, f(x_N))^T$ is

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

- The conditional distribution of f_* given \mathbf{f}

$$p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) = \mathcal{N}(\bar{m}, \bar{s}^2)$$

with

$$\bar{m} = \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f}$$

$$\bar{s}^2 = k_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*$$

Gaussian Processes - Priors over Functions

- Remember that when we modeled labels \mathbf{y} in the linear model we assumed noise with variance σ around $\mathbf{w}^T \mathbf{x}$
- We can do the same in Gaussian processes

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N p(y_i|f_i)$$

with

$$p(y_i|f_i) = \mathcal{N}(y_i|f_i, \sigma^2)$$

- Likelihood and prior are both Gaussian - conjugate!

Gaussian Processes - Priors over Functions

- Remember that when we modeled labels \mathbf{y} in the linear model we assumed noise with variance σ around $\mathbf{w}^T \mathbf{x}$
- We can do the same in Gaussian processes

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N p(y_i|f_i)$$

with

$$p(y_i|f_i) = \mathcal{N}(y_i|f_i, \sigma^2)$$

- Likelihood and prior are both Gaussian - conjugate!
- We can integrate out Gaussian process prior on \mathbf{f}

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X})d\mathbf{f}$$

- This gives

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2\mathbf{I})$$

- We can derive the predictive distribution of the function:

$$p(f_* | \mathbf{y}, \mathbf{x}_*, \mathbf{X}) = \int p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{y}, \mathbf{X}) d\mathbf{f} df_* = \mathcal{N}(m, s^2)$$

with

$$m = \mathbf{k}_*^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$
$$s^2 = k_{**} - \mathbf{k}_*^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_*$$

- Same expression as in the “Weight-Space View” section

- We can also make predictions as follows:

$$\begin{aligned} p(y_* | \mathbf{y}, \mathbf{x}_*, \mathbf{X}) &= \int p(t_* | f_*) p(f_* | \mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f} | \mathbf{y}, \mathbf{X}) d\mathbf{f} df_* \\ &= \mathcal{N}(m_{\mathbf{y}}, s_{\mathbf{y}}^2) \end{aligned}$$

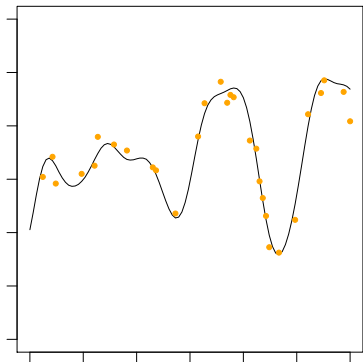
with

$$\begin{aligned} m_{\mathbf{y}} &= \mathbf{k}_*^{\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \\ s_{\mathbf{y}}^2 &= \sigma^2 + k_{**} - \mathbf{k}_*^{\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_* \end{aligned}$$

- Same expression as in the “Weight-Space View” section

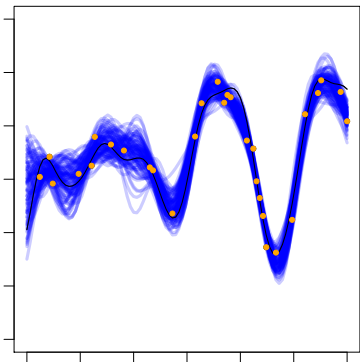
Gaussian Processes - Regression example

- Some data generated as a noisy version of some function



Gaussian Processes - Regression example

- Draws from the posterior distribution over f_* on the real line



- The kernel has parameters that have to be tuned

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp(-\beta \|\mathbf{x} - \mathbf{x}'\|^2)$$

... and there is also the noise parameter σ^2 .

- Define $\theta = (\alpha, \beta, \sigma^2)$
- How should we tune them?

- Define $\mathbf{K}_y = \mathbf{K} + \sigma^2 \mathbf{I}$
- Maximize the logarithm of the likelihood

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_y)$$

that is

$$-\frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} + \text{const.}$$

- Derivatives can be useful for gradient-based optimization

$$\frac{\partial \log[p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})]}{\partial \theta_i}$$

- Log-likelihood

$$-\frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} + \text{const.}$$

- Derivatives can be useful for gradient-based optimization:

$$\frac{\partial \log[p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})]}{\partial \theta_i} = -\frac{1}{2} \text{Tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbf{K}_y^{-1} \mathbf{y}$$

Challenges

- Non-Gaussian Likelihoods?
- Scalability?

- Marginal likelihood

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f}$$

can only be computed if $p(\mathbf{y}|\mathbf{f})$ is Gaussian

- What if $p(\mathbf{y}|\mathbf{f})$ is **not** Gaussian?

- Approximation options:
 - Local variational bounds (classification only)
 - Gibbs and MacKay, *IEEE TNN*, 2000

- Approximation options:
 - Local variational bounds (classification only)
 - Gibbs and MacKay, *IEEE TNN*, 2000
 - Laplace Approximation
 - Williams and Barber, *IEEE TPAMI*, 1998

- Approximation options:
 - Local variational bounds (classification only)
 - Gibbs and MacKay, *IEEE TNN*, 2000
 - Laplace Approximation
 - Williams and Barber, *IEEE TPAMI*, 1998
 - Expectation Propagation
 - Minka, *PhD thesis*, 2001

- Approximation options:
 - Local variational bounds (classification only)
 - Gibbs and MacKay, *IEEE TNN*, 2000
 - Laplace Approximation
 - Williams and Barber, *IEEE TPAMI*, 1998
 - Expectation Propagation
 - Minka, *PhD thesis*, 2001
 - Variational Bayes
 - Nickisch and Rasmussen, *JMLR*, 2008
 - Opper and Archambeau, *Neural Comp*, 2009

- Approximation options:
 - Local variational bounds (classification only)
 - Gibbs and MacKay, *IEEE TNN*, 2000
 - Laplace Approximation
 - Williams and Barber, *IEEE TPAMI*, 1998
 - Expectation Propagation
 - Minka, *PhD thesis*, 2001
 - Variational Bayes
 - Nickisch and Rasmussen, *JMLR*, 2008
 - Opper and Archambeau, *Neural Comp*, 2009
 - Markov chain Monte Carlo
 - Murray and Adams, *NIPS*, 2010
 - Filippone and Girolami, *IEEE TPAMI*, 2014

- Marginal likelihood

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f}$$

can only be computed if $p(\mathbf{y}|\mathbf{X}, \mathbf{f})$ is Gaussian

- ... even then

$$\log[p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})] = -\frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} + \text{const.}$$

where $\mathbf{K}_y = K(\mathbf{X}, \boldsymbol{\theta})$ is a $n \times n$ dense matrix!

- Complexity of exact method is $\mathcal{O}(n^3)$ time and $\mathcal{O}(n^2)$ space!

Tackling Gaussian case and $n \gg$

- Low-Rank Approximation options - $\mathcal{O}(nm^2)$
- Call P as a low rank approximation to \mathbf{K}_y
- Woodbury identity exploits low rank structure of P

$$\mathbf{K}_y = \begin{bmatrix} \square & & \\ & \square & \\ & & \ddots \\ & & & \square \end{bmatrix} + \begin{bmatrix} \square & & & \\ & \square & & \\ & & \ddots & \\ & & & \square \end{bmatrix}$$

$$P = \begin{bmatrix} \square & \square & \square & \\ \square & \square & \square & \\ \square & \square & \square & \\ & & & \square \end{bmatrix} + \begin{bmatrix} \square & & & \\ & \square & & \\ & & \ddots & \\ & & & \square \end{bmatrix}$$

$$P^{-1} = \begin{bmatrix} \square & & & \\ & \square & & \\ & & \ddots & \\ & & & \square \end{bmatrix}^{-1} - \begin{bmatrix} \square & & & \\ & \square & & \\ & & \ddots & \\ & & & \square \end{bmatrix}^{-1} \begin{bmatrix} \square & & \\ & \square & \\ & & \square \end{bmatrix} \left[\begin{bmatrix} \square & & \\ & \square & \\ & & \square \end{bmatrix}^{-1} + \begin{bmatrix} \square & \square & \square & \\ \square & \square & \square & \\ \square & \square & \square & \\ & & & \square \end{bmatrix}^{-1} \begin{bmatrix} \square & \square & \square & \\ \square & \square & \square & \\ \square & \square & \square & \\ & & & \square \end{bmatrix}^{-1} \right]$$

- Low-Rank Approximation options - $\mathcal{O}(nm^2)$
 - Subset-of-data 'sparse' methods
 - Smola and Bartlett, *NIPS*, 2001
 - Seeger and Williams, *AISTATS*, 2003

- Low-Rank Approximation options - $\mathcal{O}(nm^2)$
 - Subset-of-data 'sparse' methods
 - Smola and Bartlett, *NIPS*, 2001
 - Seeger and Williams, *AISTATS*, 2003
 - Pseudo-inputs introduced
 - Snelson and Ghahramani, *NIPS*, 2005

- Low-Rank Approximation options - $\mathcal{O}(nm^2)$
 - Subset-of-data 'sparse' methods
 - Smola and Bartlett, *NIPS*, 2001
 - Seeger and Williams, *AISTATS*, 2003
 - Pseudo-inputs introduced
 - Snelson and Ghahramani, *NIPS*, 2005
 - A unifying view brings several ideas together
 - Quiñero-Candela and Rasmussen, *JMLR*, 2005

- Low-Rank Approximation options - $\mathcal{O}(nm^2)$
 - Subset-of-data 'sparse' methods
 - Smola and Bartlett, *NIPS*, 2001
 - Seeger and Williams, *AISTATS*, 2003
 - Pseudo-inputs introduced
 - Snelson and Ghahramani, *NIPS*, 2005
 - A unifying view brings several ideas together
 - Quiñero-Candela and Rasmussen, *JMLR*, 2005
 - Variational approach for better placement of pseudo points
 - Titsias, *AISTATS*, 2009

- Low-Rank Approximation options - $\mathcal{O}(nm^2)$
 - Subset-of-data 'sparse' methods
 - Smola and Bartlett, *NIPS*, 2001
 - Seeger and Williams, *AISTATS*, 2003
 - Pseudo-inputs introduced
 - Snelson and Ghahramani, *NIPS*, 2005
 - A unifying view brings several ideas together
 - Quiñero-Candela and Rasmussen, *JMLR*, 2005
 - Variational approach for better placement of pseudo points
 - Titsias, *AISTATS*, 2009
 - Random feature expansions
 - Rahimi and Recht, *NIPS*, 2008
 - Lazaro-Gredilla et al., *JMLR*, 2010

- Approximation options:
 - Structured approximations based on Toeplitz/circulant matrices - $\mathcal{O}(dn^{\frac{d+1}{d}})$ time
 - Wilson and Nickisch, *ICML*, 2015
 - Gilboa et al., *IEEE TPAMI*, 2015

- Approximation options:
 - Structured approximations based on Toeplitz/circulant matrices - $\mathcal{O}(dn^{\frac{d+1}{d}})$ time
 - Wilson and Nickisch, *ICML*, 2015
 - Gilboa et al., *IEEE TPAMI*, 2015
 - Stochastic-gradient optimization/inference **without** model approximations - $\mathcal{O}(n^2)$ time and $\mathcal{O}(n)$ space
 - Filippone and Engler, *ICML*, 2015
 - Cutajar, Osborne, Cunningham, Filippone, *ICML*, 2016

- Marginal likelihood

$$\log[p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})] = -\frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} + \text{const.}$$

- Derivatives wrt $\boldsymbol{\theta}$

$$\frac{\partial \log[p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})]}{\partial \theta_i} = -\frac{1}{2} \text{Tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbf{K}_y^{-1} \mathbf{y}$$

- Stochastic estimate of the trace

$$\text{Tr} \left(K^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \right) = \text{Tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbb{E}[\mathbf{r}\mathbf{r}^T] \right) = \mathbb{E} \left[\mathbf{r}^T \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbf{r} \right]$$

with $\mathbb{E}[\mathbf{r}\mathbf{r}^T] = I$

- Stochastic estimate of the trace

$$\text{Tr} \left(K^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \right) = \text{Tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbb{E}[\mathbf{r}\mathbf{r}^T] \right) = \mathbb{E} \left[\mathbf{r}^T \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbf{r} \right]$$

with $\mathbb{E}[\mathbf{r}\mathbf{r}^T] = I$

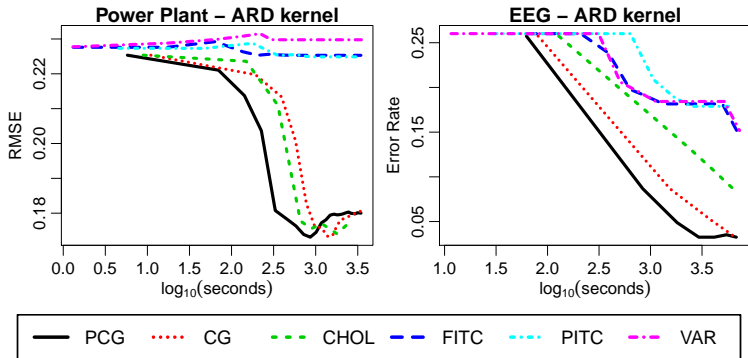
- Stochastic gradient

$$-\frac{1}{2N_r} \sum_{i=1}^{N_r} \mathbf{r}^{(i)T} \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbf{r}^{(i)} + \frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_i} \mathbf{K}_y^{-1} \mathbf{y}$$

- Linear systems only!

Teaser - Preconditioning Kernel Matrices

- Stochastic Gradient Optimization



Cutajar, Osborne, Cunningham, Filippone, *ICML*, 2016

- Non-Gaussian Likelihoods?
- Scalability?

Modern GP works tackle both

Modern Gaussian Processes

- Mini-batch-based learning - $\mathcal{O}(1)$ time for each iteration!
- Exploit GPU and distributed computing
- Automatic differentiation
- Application-specific representations (e.g., convolutional)

Stochastic Gradient Optimization

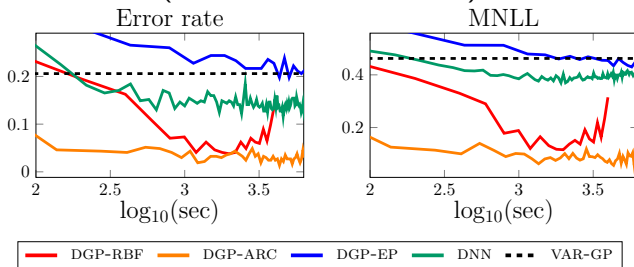
Robbins and Monro, *AoMS*, 1951

- Approximation options:
 - Scalable Expectation Propagation
 - Bui et al., *ICML*, 2016

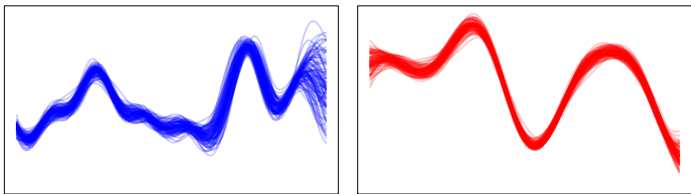
- Approximation options:
 - Scalable Expectation Propagation
 - Bui et al., *ICML*, 2016
 - Inducing points methods
 - Hensman et al., *AISTATS*, 2013
 - Hensman, Matthews, Ghahramani, Filippone, *NIPS*, 2015

- Approximation options:
 - Scalable Expectation Propagation
 - Bui et al., *ICML*, 2016
 - Inducing points methods
 - Hensman et al., *AISTATS*, 2013
 - Hensman, Matthews, Ghahramani, Filippone, *NIPS*, 2015
 - Random feature expansions
 - Gal, Ghahramani, *ICML*, 2016
 - Cutajar, Bonilla, Michiardi, Filippone, *ICML*, 2017

EEG dataset ($n = 14979, d = 14$)



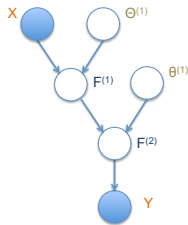
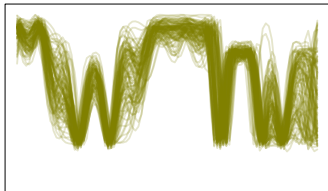
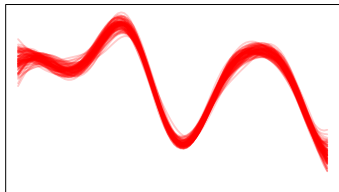
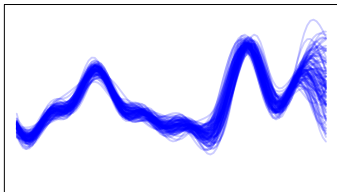
- Composition of processes - Deep Gaussian Processes



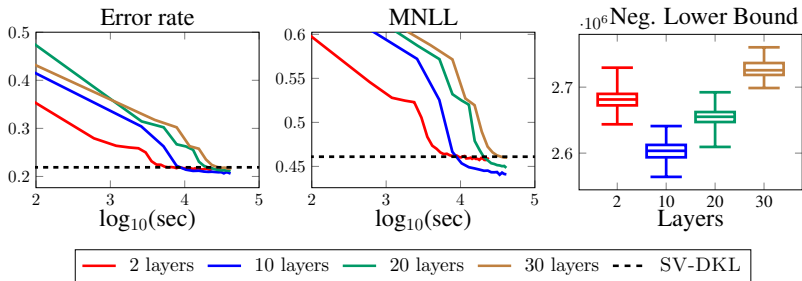
$$(f \circ g)(x)??$$

Teaser - Modern GPs - Any likelihood and $n \gg$

- Composition of processes



Airline dataset ($n = 5\text{M}+$, $d = 8$)



Other interesting topics

- Bayesian Optimization
 - Jones et al., *JoGO*, 1998

Other interesting topics

- Bayesian Optimization
 - Jones et al., *JoGO*, 1998
- Unsupervised Learning
 - Lawrence, *NIPS*, 2004

Other interesting topics

- Bayesian Optimization
 - Jones et al., *JoGO*, 1998
- Unsupervised Learning
 - Lawrence, *NIPS*, 2004
- Deep Gaussian Processes
 - Damianou and Lawrence, *AISTATS*, 2013
 - Cutajar, Bonilla, Michiardi, and Filippone, *ICML*, 2017

- Bayesian Optimization
 - Jones et al., *JoGO*, 1998
- Unsupervised Learning
 - Lawrence, *NIPS*, 2004
- Deep Gaussian Processes
 - Damianou and Lawrence, *AISTATS*, 2013
 - Cutajar, Bonilla, Michiardi, and Filippone, *ICML*, 2017
- Convolutional Gaussian Processes
 - Wilson et al., *AISTATS*, 2015
 - Wilson et al., *NIPS*, 2016
 - van der Wilk et al., *NIPS*, 2017

- Bayesian Optimization
 - Jones et al., *JoGO*, 1998
- Unsupervised Learning
 - Lawrence, *NIPS*, 2004
- Deep Gaussian Processes
 - Damianou and Lawrence, *AISTATS*, 2013
 - Cutajar, Bonilla, Michiardi, and Filippone, *ICML*, 2017
- Convolutional Gaussian Processes
 - Wilson et al., *AISTATS*, 2015
 - Wilson et al., *NIPS*, 2016
 - van der Wilk et al., *NIPS*, 2017
- Structured output
 - Galliani et al., *AISTATS*, 2017

- Bayesian Optimization
 - Jones et al., *JoGO*, 1998
- Unsupervised Learning
 - Lawrence, *NIPS*, 2004
- Deep Gaussian Processes
 - Damianou and Lawrence, *AISTATS*, 2013
 - Cutajar, Bonilla, Michiardi, and Filippone, *ICML*, 2017
- Convolutional Gaussian Processes
 - Wilson et al., *AISTATS*, 2015
 - Wilson et al., *NIPS*, 2016
 - van der Wilk et al., *NIPS*, 2017
- Structured output
 - Galliani et al., *AISTATS*, 2017
- Probabilistic Numerics
 - Fitzsimons, Cutajar, Osborne, Roberts, Filippone, *UAI*, 2017

Thank you!



AXA
Research Fund