

Bayesian Deep Learning

Maurizio Filippone

EURECOM, Sophia Antipolis, France

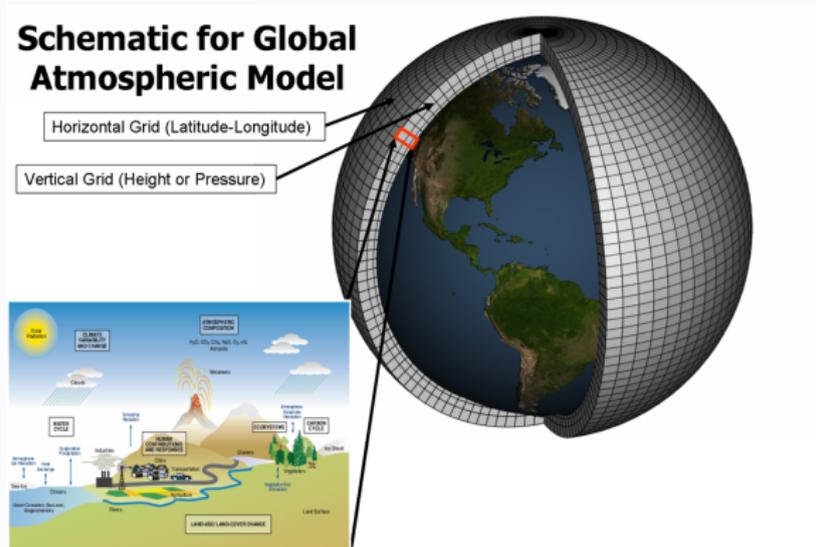
January 10th, 2019

- ① Motivation
- ② Probabilistic Deep Nets
 - Scalable Inference
 - Connections with (Deep) Gaussian Processes
- ③ Some Results
- ④ Conclusions

Motivation

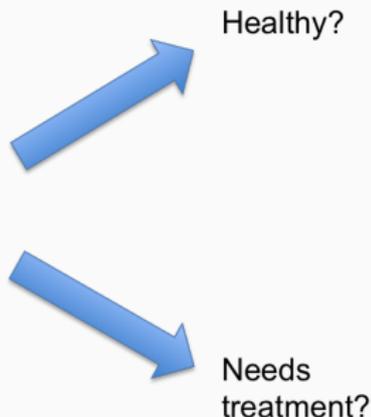
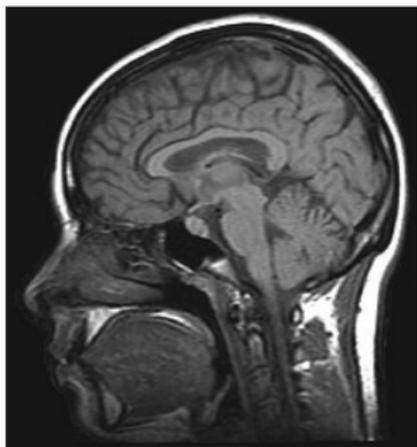
Quantification of Uncertainty with Expensive Models

- Climate modeling



Quantification of Uncertainty with No Models

- Classification and progression of neurodegenerative diseases



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 statistical model

A Unified Framework

A model might be expensive to simulate/inaccurate

- Emulate model/discrepancy using a surrogate

A model might not even be available

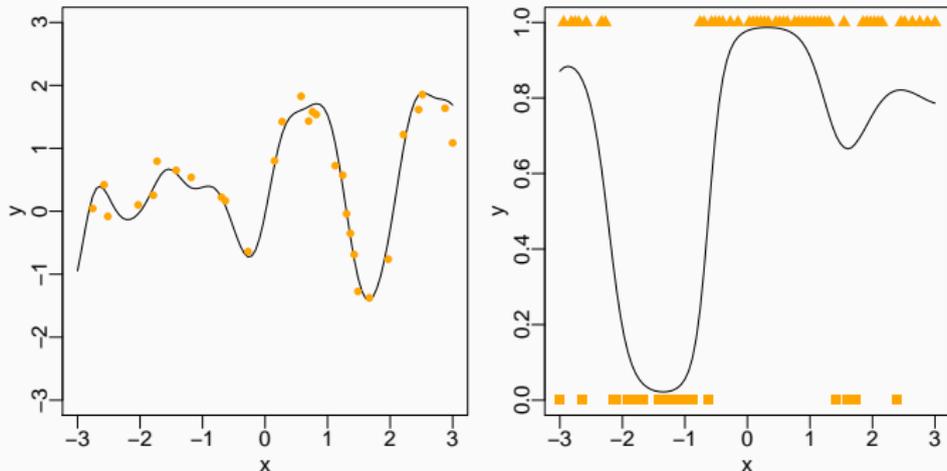
- Replace it with a flexible statistical model

**Probabilistic Deep Models for Accurate Modeling and
Quantification of Uncertainty**

Probabilistic Deep Nets

Learning from Data – Function Estimation

- Take these two examples



- We are interested in estimating a function $f(x)$ from data
- Most problems in Machine Learning can be cast this way!

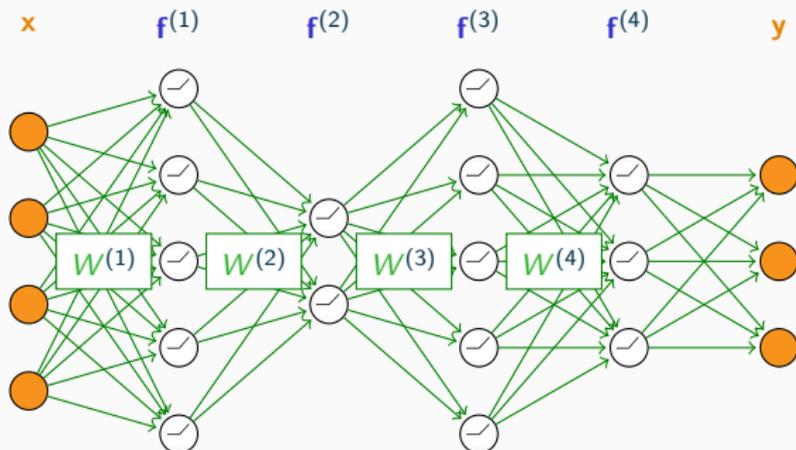
Deep Neural Networks

- Implement a composition of parametric functions

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}^{(L)} \left(\mathbf{f}^{(L-1)} \left(\dots \mathbf{f}^{(1)}(\mathbf{x}) \dots \right) \right)$$

with

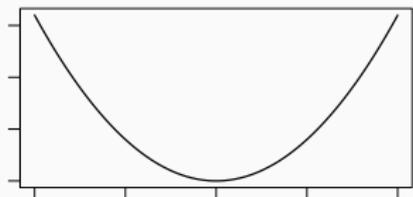
$$\mathbf{f}^{(l)}(\mathbf{h}) = \mathbf{g} \left(\mathbf{h}^\top \mathbf{W}^{(l)} \right)$$



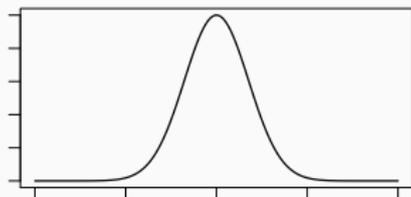
Back-propagation – Probabilistic Interpretation Loss

- Inputs : $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Labels : $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$
- Weights : $W = \{W^{(1)}, \dots, W^{(L)}\}$

Quadratic Loss



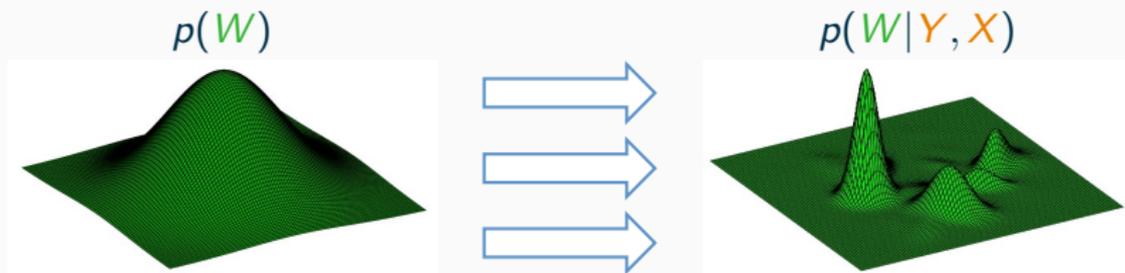
$$p(Y|X, W) \propto \exp(-\text{Loss})$$



- Back-propagation minimizes a loss function
- ... equivalent as optimizing likelihood $p(Y|X, W)$

Bayesian Inference

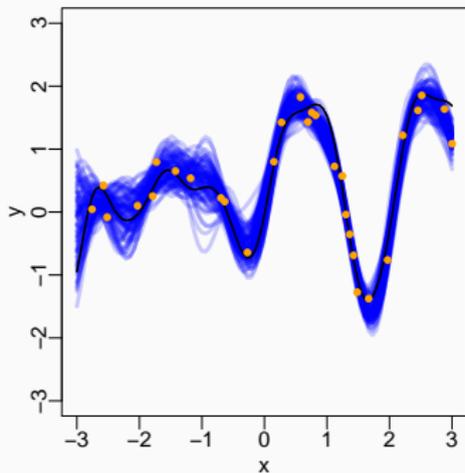
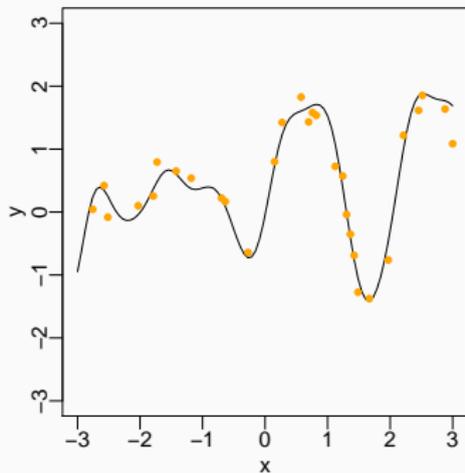
- Inputs : $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Labels : $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$
- Weights : $W = \{W^{(1)}, \dots, W^{(L)}\}$



$$p(W|Y, X) = \frac{p(Y|X, W)p(W)}{\int p(Y|X, W)p(W)dW}$$

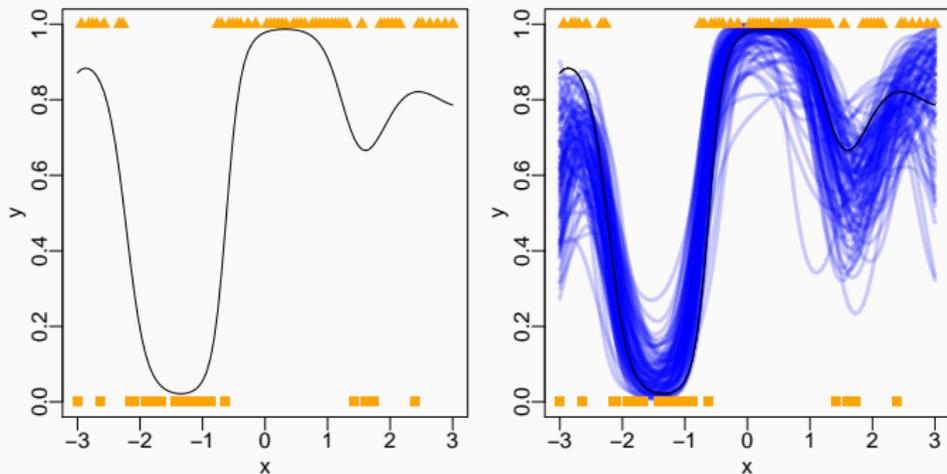
Bayesian Deep Neural Networks

- Regression example



Bayesian Deep Neural Networks

- Classification example



Stochastic Variational Inference

- Bayesian inference is intractable due to this integral

$$\log [p(Y|X)] = \log \left[\int p(Y|X, W)p(W)dW \right]$$

Stochastic Variational Inference

- Bayesian inference is intractable due to this integral

$$\log [p(Y|X)] = \log \left[\int p(Y|X, W)p(W)dW \right]$$

- Lower bound for $\log [p(Y|X)]$

$$E_{q(W)} (\log [p(Y|X, W)]) - \text{KL} [q(W)||p(W)],$$

where $q(W)$ approximates $p(W|Y, X)$.

- Kullback-Leibler divergence KL – “distance” between q and p

Stochastic Variational Inference

- Bayesian inference is intractable due to this integral

$$\log [p(Y|X)] = \log \left[\int p(Y|X, W)p(W)dW \right]$$

- Lower bound for $\log [p(Y|X)]$

$$E_{q(W)} (\log [p(Y|X, W)]) - \text{KL} [q(W)||p(W)],$$

where $q(W)$ approximates $p(W|Y, X)$.

- Kullback-Leibler divergence KL – “distance” between q and p

Optimize the lower bound wrt the parameters of $q(W)$

Stochastic Variational Inference

- Assume that the likelihood factorizes

$$p(Y|X, W) = \prod_k p(y_k|x_k, W)$$

Stochastic Variational Inference

- Assume that the likelihood factorizes

$$p(\mathbf{Y}|\mathbf{X}, W) = \prod_k p(\mathbf{y}_k|\mathbf{x}_k, W)$$

- Doubly stochastic **unbiased** estimate of the expectation term
 - Mini-batch

$$\mathbb{E}_{q(W)} (\log [p(\mathbf{Y}|\mathbf{X}, W)]) \approx \frac{n}{m} \sum_{k \in \mathcal{I}_m} \mathbb{E}_{q(W)} (\log [p(\mathbf{y}_k|\mathbf{x}_k, W)])$$

- Monte Carlo

$$\mathbb{E}_{q(W)} (\log [p(\mathbf{y}_k|\mathbf{x}_k, W)]) \approx \frac{1}{N_{\text{MC}}} \sum_{r=1}^{N_{\text{MC}}} \log [p(\mathbf{y}_k|\mathbf{x}_k, \tilde{W}_r)]$$

with $\tilde{W}_r \sim q(W)$.

- Assume a factorized Gaussian approximate posterior:

$$q(W) = \prod_{ijl} q\left(W_{ij}^{(l)}\right) = \prod_{ijl} \mathcal{N}\left(\mu_{ij}^{(l)}, (\sigma^2)_{ij}^{(l)}\right) \quad (1)$$

Stochastic Variational Inference

- Assume a factorized Gaussian approximate posterior:

$$q(\mathbf{W}) = \prod_{ijl} q\left(\mathbf{W}_{ij}^{(l)}\right) = \prod_{ijl} \mathcal{N}\left(\mu_{ij}^{(l)}, (\sigma^2)_{ij}^{(l)}\right) \quad (1)$$

- Reparameterization trick

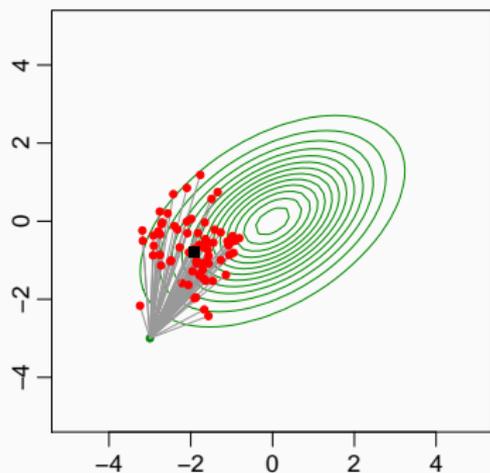
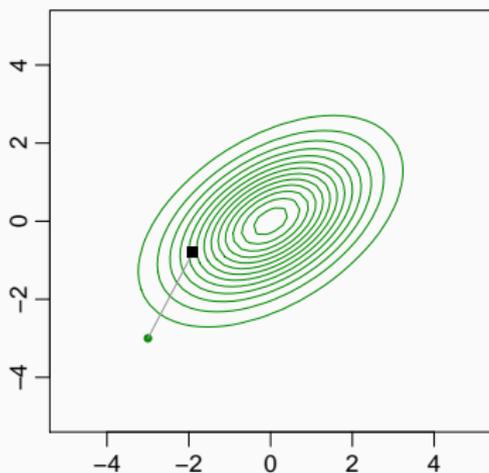
$$(\tilde{\mathbf{W}}_r^{(l)})_{ij} = \sigma_{ij}^{(l)} \varepsilon_{rij}^{(l)} + \mu_{ij}^{(l)},$$

with $\varepsilon_{rij}^{(l)} \sim \mathcal{N}(0, 1)$

- Optimization wrt $\mu_{ij}^{(l)}, (\sigma^2)_{ij}^{(l)}$ with automatic differentiation

Stochastic Gradient Optimization

$$\mathbb{E} \left\{ \widetilde{\nabla_{\text{par}_q} \text{LowerBound}} \right\} = \nabla_{\text{par}_q} \text{LowerBound}$$



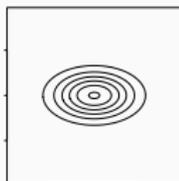
Stochastic Variational Inference - Simple Illustration

$$\text{par}_q' = \text{par}_q + \frac{\alpha_t}{2} \widetilde{\nabla}_{\text{par}_q}(\text{LowerBound}) \quad \alpha_t \rightarrow 0$$

Form of Approximating Distribution

Approximating distribution $q(W)$ can have the following forms:

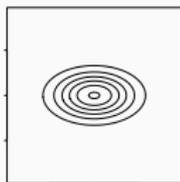
- Fully factorized



Form of Approximating Distribution

Approximating distribution $q(W)$ can have the following forms:

- Fully factorized



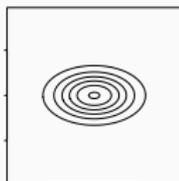
- Full covariance $\Sigma = LL^T$



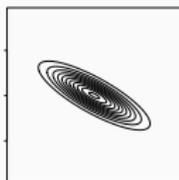
Form of Approximating Distribution

Approximating distribution $q(W)$ can have the following forms:

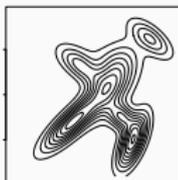
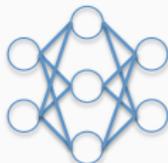
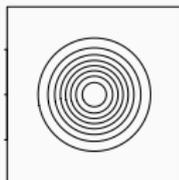
- Fully factorized



- Full covariance $\Sigma = LL^T$

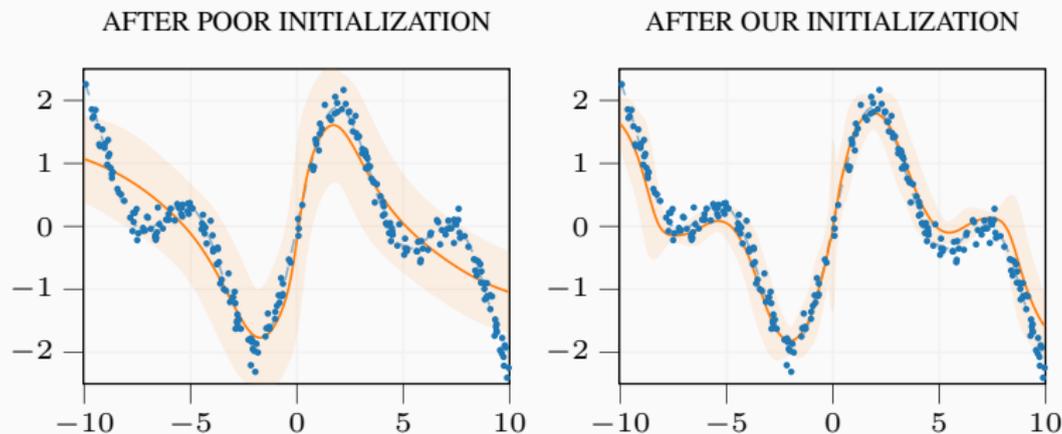


- Normalizing Flows, Real NVPs, Stein VI - change of measure determined by $\det(\text{Jacobian})$



Initialization of SVI matters

- Initialization can be an issue
- We proposed a novel way to initialize SVI well



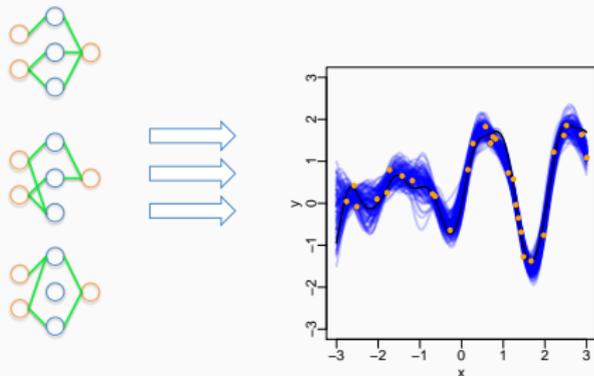
Is There an Easier Way?

- Dropout is Variational Inference with Bernoulli-like $q(W)$
- At training time, apply dropout

Iteration 1 Iteration 2 Iteration 3 ...

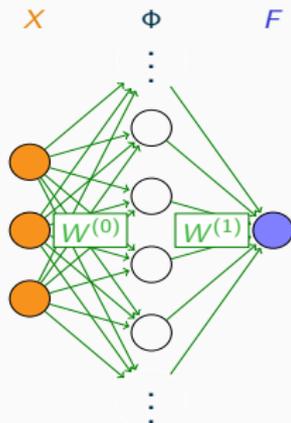


- At test time, “sample” networks with different dropout masks



Gaussian Processes as Infinitely-Wide Shallow Neural Nets

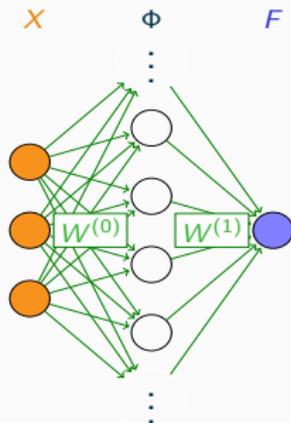
- Take $W^{(i)} \sim \mathcal{N}(\mathbf{0}, \alpha_i I)$
- Central Limit Theorem implies that F is Gaussian



- F has zero-mean
- $\text{cov}(F) = \mathbb{E}_{p(W^{(0)}, W^{(1)})} [\Phi(XW^{(0)}) W^{(1)} W^{(1)\top} \Phi(XW^{(0)})^\top]$

Gaussian Processes as Infinitely-Wide Shallow Neural Nets

- Take $W^{(i)} \sim \mathcal{N}(\mathbf{0}, \alpha_i I)$
- Central Limit Theorem implies that F is Gaussian



- F has zero-mean
- $\text{cov}(F) = \alpha_1 \mathbb{E}_{p(W^{(0)})} [\Phi(XW^{(0)})\Phi(XW^{(0)})^\top]$
- Some choices of Φ lead to analytic expression of known kernels (RBF, Matérn, arc-cosine, Brownian motion, ...)

Random Feature Expansions for DGPs - Bochner's theorem

- Continuous shift-invariant covariance function

$$k(\mathbf{x}_i - \mathbf{x}_j | \theta) = \sigma^2 \int p(\boldsymbol{\omega} | \theta) \exp\left(i(\mathbf{x}_i - \mathbf{x}_j)^\top \boldsymbol{\omega}\right) d\boldsymbol{\omega}$$

Random Feature Expansions for DGPs - Bochner's theorem

- Continuous shift-invariant covariance function

$$k(\mathbf{x}_i - \mathbf{x}_j | \theta) = \sigma^2 \int p(\boldsymbol{\omega} | \theta) \exp\left(i(\mathbf{x}_i - \mathbf{x}_j)^\top \boldsymbol{\omega}\right) d\boldsymbol{\omega}$$

- Monte Carlo estimate

$$k(\mathbf{x}_i - \mathbf{x}_j | \theta) \approx \frac{\sigma^2}{N_{\text{RF}}} \sum_{r=1}^{N_{\text{RF}}} \mathbf{z}(\mathbf{x}_i | \tilde{\boldsymbol{\omega}}_r)^\top \mathbf{z}(\mathbf{x}_j | \tilde{\boldsymbol{\omega}}_r)$$

with

$$\tilde{\boldsymbol{\omega}}_r \sim p(\boldsymbol{\omega} | \theta)$$

$$\mathbf{z}(\mathbf{x} | \boldsymbol{\omega}) = [\cos(\mathbf{x}^\top \boldsymbol{\omega}), \sin(\mathbf{x}^\top \boldsymbol{\omega})]^\top$$

Random Feature Expansions for DGPs

- Define

$$\Phi^{(l)} = \sqrt{\frac{\sigma^2}{N_{\text{RF}}^{(l)}}} \left[\cos \left(F^{(l)} \Omega^{(l)} \right), \sin \left(F^{(l)} \Omega^{(l)} \right) \right]$$

and

$$F^{(l+1)} = \Phi^{(l)} W^{(l)}$$

- We are stacking Bayesian linear models with

$$p \left(W_{\cdot i}^{(l)} \right) = \mathcal{N} \left(\mathbf{0}, I \right)$$

Random Feature Expansions for DGPs

- Define

$$\Phi^{(l)} = \sqrt{\frac{\sigma^2}{N_{\text{RF}}^{(l)}}} \left[\cos \left(F^{(l)} \Omega^{(l)} \right), \sin \left(F^{(l)} \Omega^{(l)} \right) \right]$$

and

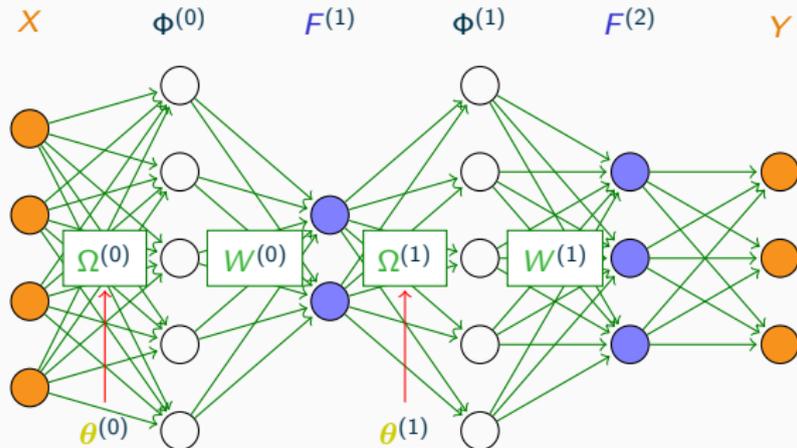
$$F^{(l+1)} = \Phi^{(l)} W^{(l)}$$

- We are stacking Bayesian linear models with

$$p \left(W_{\cdot i}^{(l)} \right) = \mathcal{N} \left(\mathbf{0}, I \right)$$

- Expansion of arc-cosine kernel yields ReLU activations!

Random Feature Expansions make Deep GPs become DNNs

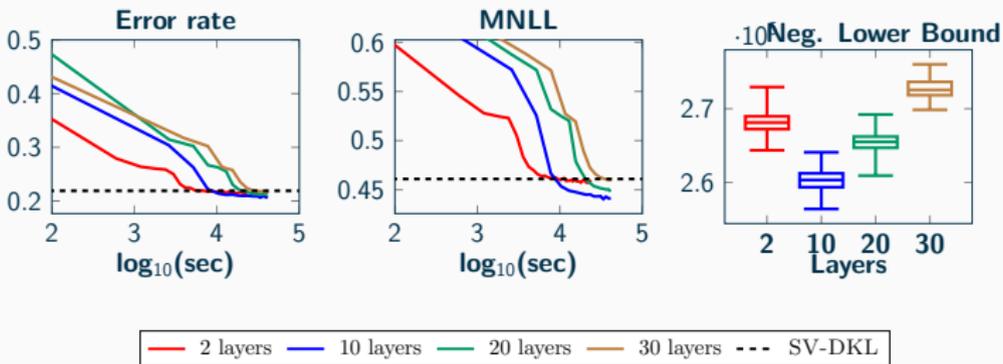


Some Results

Results - Model (Depth) Selection

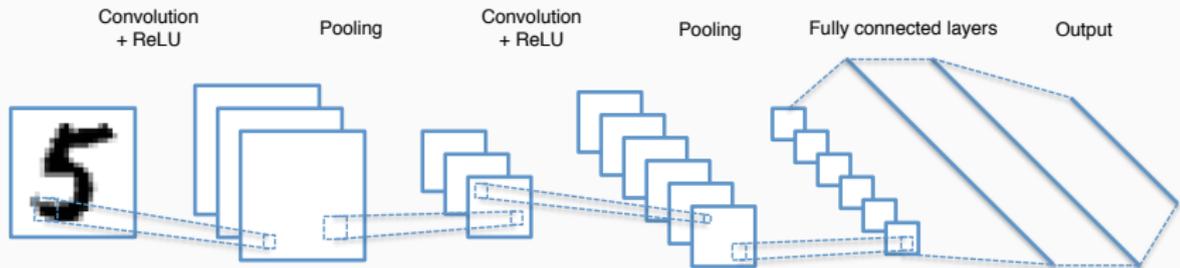
Airline dataset

($n = 5M+$, $d = 8$)



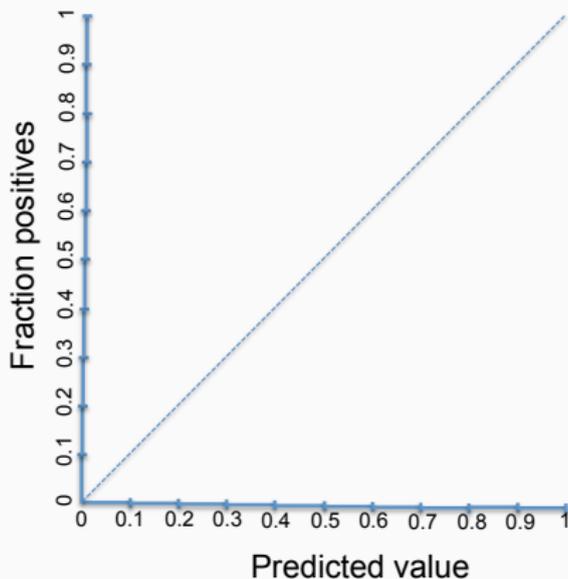
Convolutional Nets

- Convolutional nets are widely used. . .
- . . .but they are known to be overconfident!



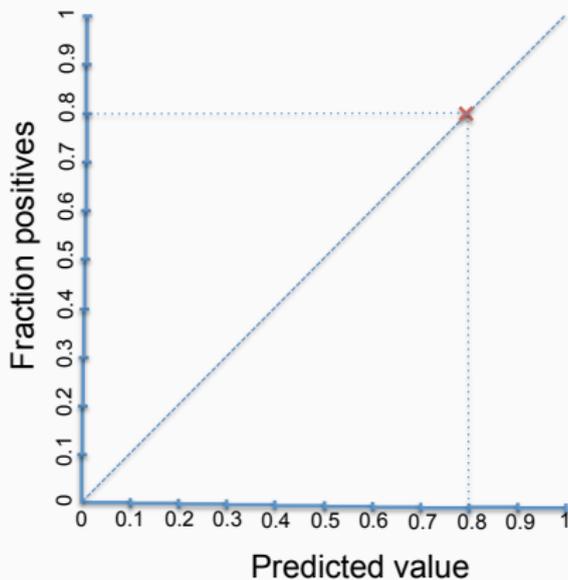
Calibration as a Measure of Quantification of Uncertainty

- Reliability diagrams



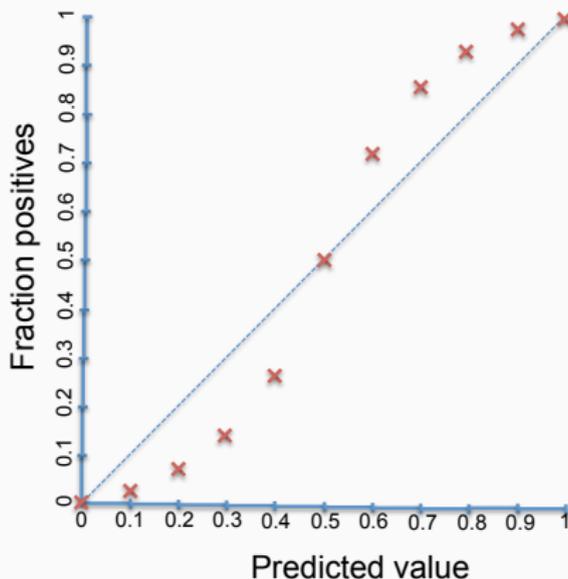
Calibration as a Measure of Quantification of Uncertainty

- Reliability diagrams



Calibration as a Measure of Quantification of Uncertainty

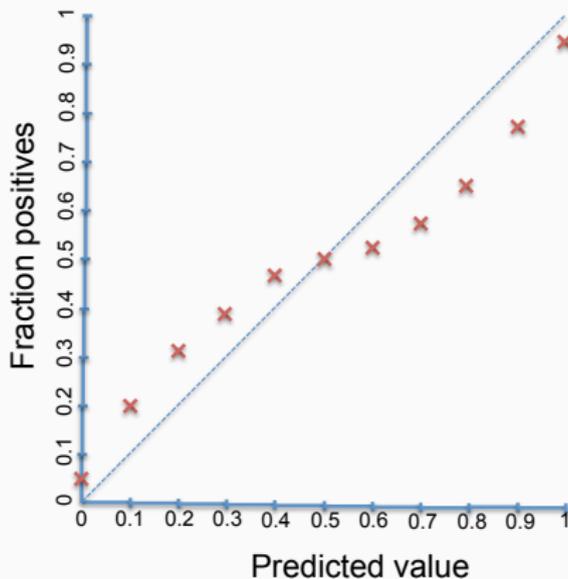
- Reliability diagrams - Under-confident predictions



- We can extract the Expected Calibration Error (ECE) score
- The BRIER score is another measure of calibration

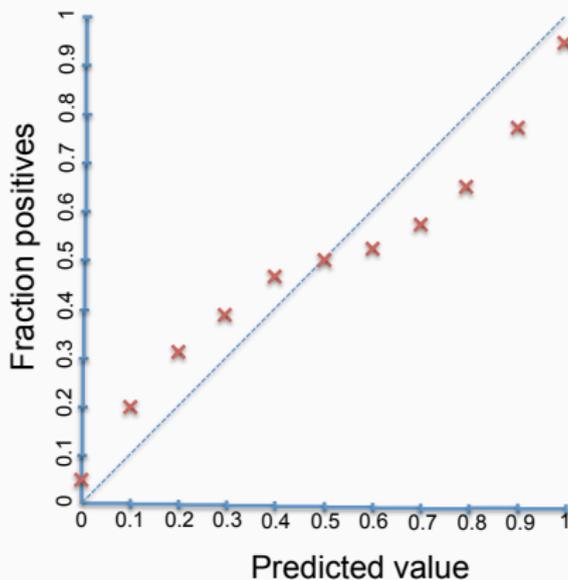
Calibration as a Measure of Quantification of Uncertainty

- Reliability diagrams - Overconfident predictions



Calibration as a Measure of Quantification of Uncertainty

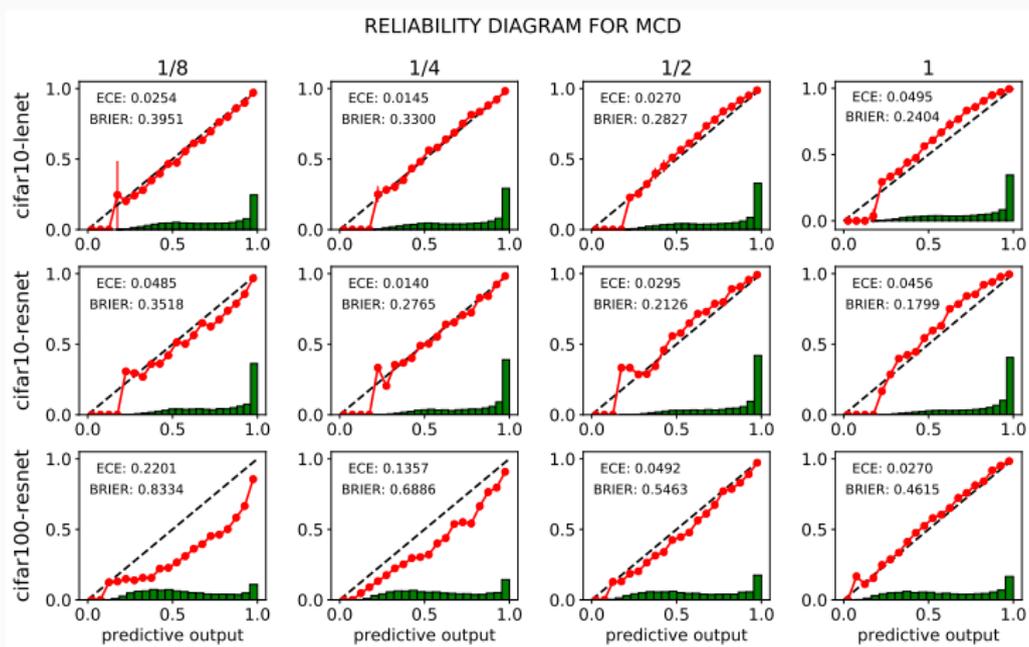
- Reliability diagrams - Overconfident predictions



Reliability diagrams of modern Deep CNNs look like this!
Bayesian treatment of filters fixes it!

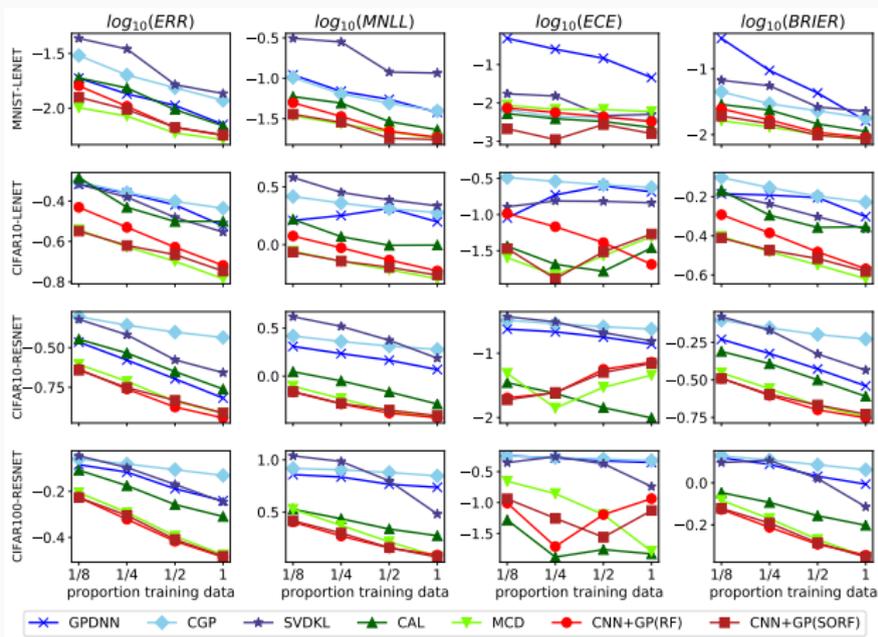
Bayesian CNNs are calibrated

- Inferring parameters of convolutional filter recovers calibration
- Example with Monte Carlo Dropout



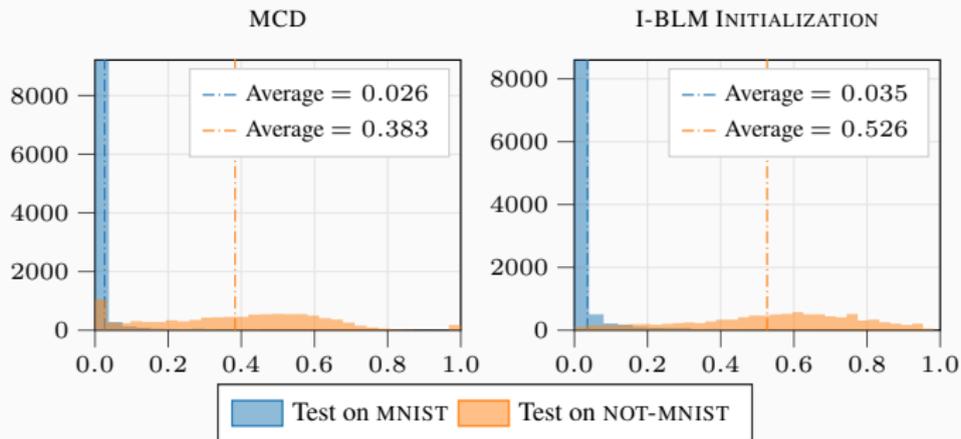
Performance Evaluation of Bayesian CNNs

- Bayesian CNNs are calibrated and achieve better performance than post calibrated CNNs



Knowing When the Model Doesn't Know

- Training on MNIST and test on not-MNIST



Conclusions

Conclusions

- Inference for Deep Nets is hard
 - Scalable stochastic-based approximate inference but...
 - ... it is difficult to assess the impact approximations on quantification of uncertainty

Conclusions

- Inference for Deep Nets is hard
 - Scalable stochastic-based approximate inference but...
 - ... it is difficult to assess the impact approximations on quantification of uncertainty
- The connection between Deep Nets and Deep Gaussian processes can have implications on
 - Understanding Deep Learning
 - Deriving sensible priors for Deep Learning
 - Improving inference borrowing algebraic/computational tricks from kernel literature

Conclusions

- Inference for Deep Nets is hard
 - Scalable stochastic-based approximate inference but...
 - ... it is difficult to assess the impact approximations on quantification of uncertainty
- The connection between Deep Nets and Deep Gaussian processes can have implications on
 - Understanding Deep Learning
 - Deriving sensible priors for Deep Learning
 - Improving inference borrowing algebraic/computational tricks from kernel literature
- Cool stuff
 - New hardware
 - Bayesian compression

We are hiring!

We are hiring PhDs, Post-docs and Assistant Professors



Thank you!



AXA
Research Fund