

Dropout as a Bayesian Approximation

Ivan Rodriguez

Main reference papers:

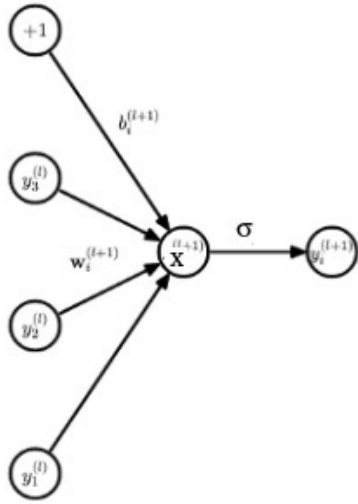
- 1 - “Dropout as a bayesian approximation: Representing model uncertainty in deep learning”, Y.Gal and Z.Ghahramani – ICML’16 ~ 6400 citations
- 2 - “Dropout as a Bayesian Approximation: Appendix”, Y.Gal and Z.Ghahramani - ICML’16

Seminar outline:

- 1 - Standard Dropout Introduction
- 2 - Gaussian Process for DNN
- 3 - Dropout from a Bayesian point of view
- 4 - Results

1 - Standard Dropout Introduction

Original dropout method

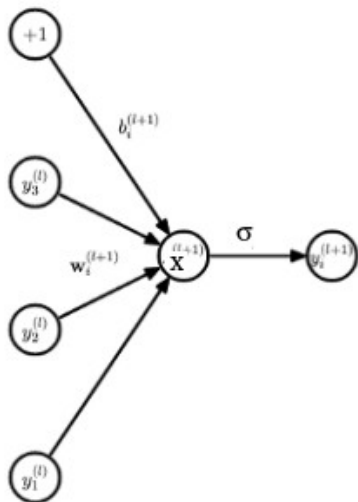


(a) Standard network

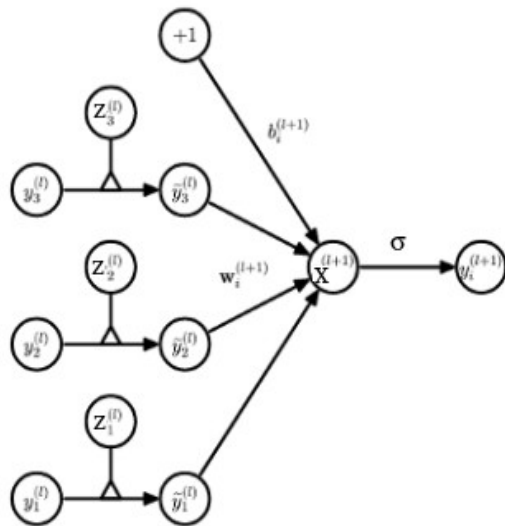
$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



(a) Standard network



(b) Dropout network

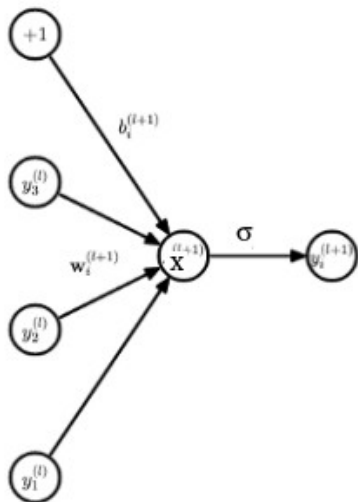
$z^{(l)} = (z_1, z_2, \dots)$ with $z_i \sim \text{Bernoulli}(p)$

$\tilde{y}^{(l)} = z^{(l)} \odot y^{(l)} = (z_1^{(l)} y_1^{(l)}, z_2^{(l)} y_2^{(l)}, \dots)$

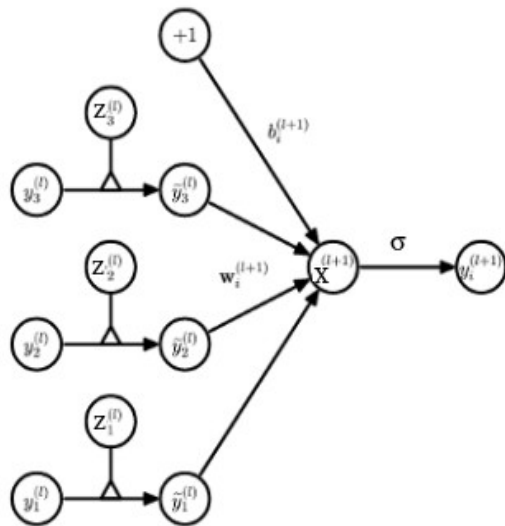
$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



(a) Standard network



(b) Dropout network

$z^{(l)} = (z_1, z_2, \dots)$ with $z_i \sim \text{Bernoulli}(p)$

$\tilde{y}^{(l)} = z^{(l)} \odot y^{(l)} = (z_1^{(l)} y_1^{(l)}, z_2^{(l)} y_2^{(l)}, \dots)$

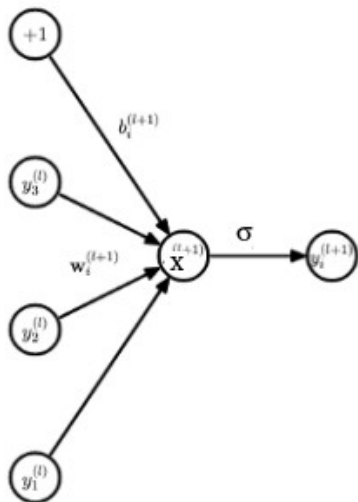
$x^{(l+1)} = W^{(l+1)} \tilde{y}^{(l)} + b^{(l+1)} = \boxed{W^{(l+1)} \text{diag}(z^{(l)})} y^{(l)} + b^{(l+1)}$

$y^{(l+1)} = \sigma(x^{(l+1)})$

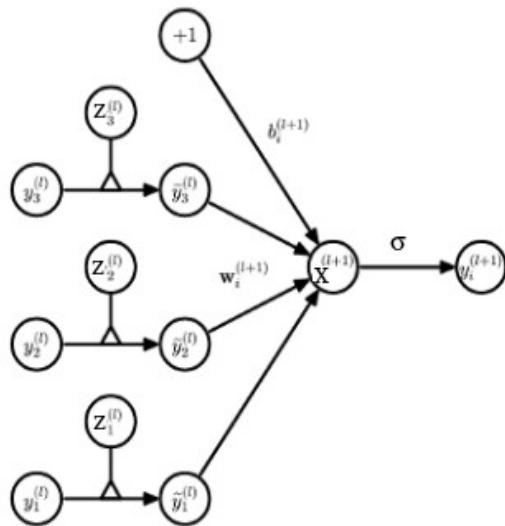
$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



(a) Standard network



(b) Dropout network

$z^{(l)} = (z_1, z_2, \dots)$ with $z_i \sim \text{Bernoulli}(p)$

$\tilde{y}^{(l)} = z^{(l)} \odot y^{(l)} = (z_1^{(l)} y_1^{(l)}, z_2^{(l)} y_2^{(l)}, \dots)$

$x^{(l+1)} = W^{(l+1)} \tilde{y}^{(l)} + b^{(l+1)} = W^{(l+1)} \text{diag}(z^{(l)}) y^{(l)} + b^{(l+1)}$

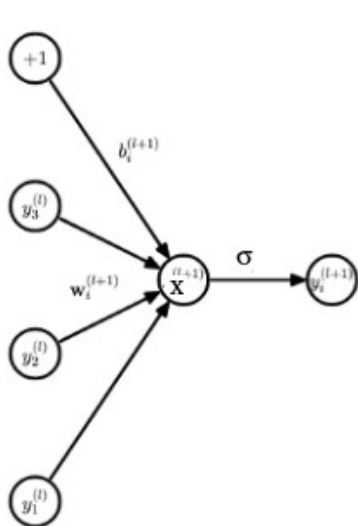
$y^{(l+1)} = \sigma(x^{(l+1)})$

$\tilde{W}_{\alpha\beta}^{(l+1)} = W_{\alpha\beta}^{(l+1)} z_{\beta}^{(l)} = 0$ if $z_{\beta} = 0$ **column β vanishes**

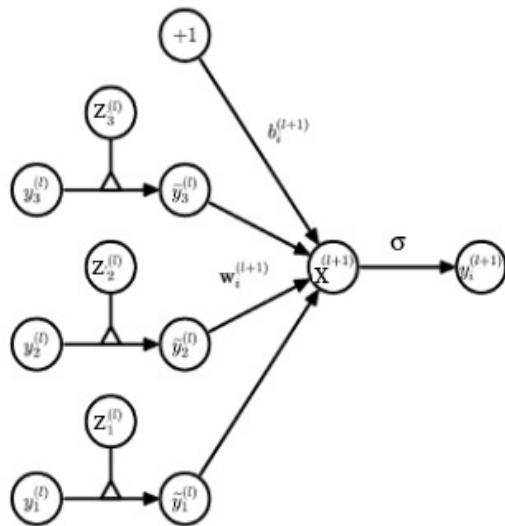
$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



(a) Standard network



(b) Dropout network

Dropout interpretation: Ensemble model

| | \tilde{y}_1 | \tilde{y}_2 | \tilde{y}_3 | |
|--|---------------|---------------|---------------|-------------|
| | 1 | 1 | 1 | Model 1 |
| | 1 | 1 | 0 | Model 2 |
| | | ⋮ | | |
| | 0 | 0 | 0 | Model 2^N |

$z^{(l)} = (z_1, z_2, \dots)$ with $z_i \sim \text{Bernoulli}(p)$

$\tilde{y}^{(l)} = z^{(l)} \odot y^{(l)} = (z_1^{(l)} y_1^{(l)}, z_2^{(l)} y_2^{(l)}, \dots)$

$$x^{(l+1)} = W^{(l+1)} \tilde{y}^{(l)} + b^{(l+1)} = W^{(l+1)} \text{diag}(z^{(l)}) y^{(l)} + b^{(l+1)}$$

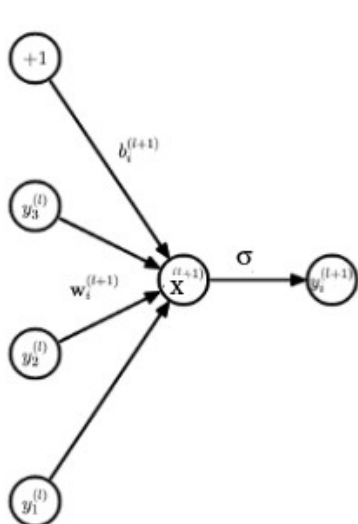
$$y^{(l+1)} = \sigma(x^{(l+1)})$$

$\tilde{W}_{\alpha\beta}^{(l+1)} = W_{\alpha\beta}^{(l+1)} z_{\beta}^{(l)} = 0$ if $z_{\beta} = 0$ **column β vanishes**

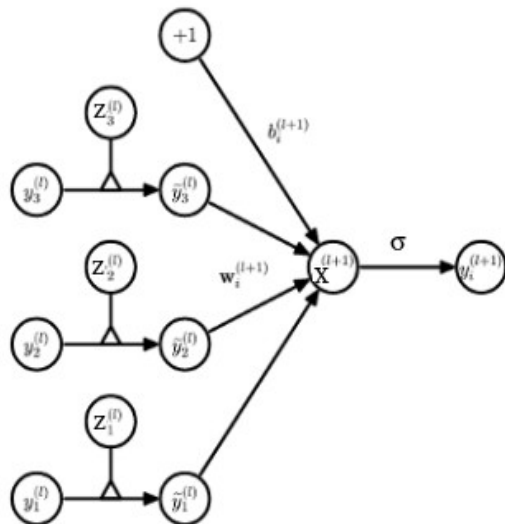
$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



(a) Standard network



(b) Dropout network

Dropout interpretation: Ensemble model

| \tilde{y}_1 | \tilde{y}_2 | \tilde{y}_3 | |
|---------------|---------------|---------------|-------------|
| 1 | 1 | 1 | Model 1 |
| 1 | 1 | 0 | Model 2 |
| | ⋮ | | |
| 0 | 0 | 0 | Model 2^N |

Inference:

$$E[\tilde{y}_i] = E[z_i y_i] = p y_i + (1-p)0 = p y_i$$

$$E[x_i] = \sum_j W_{ij} p y_j + b_i \rightarrow W \rightarrow pW$$

$$z^{(l)} = (z_1, z_2, \dots) \text{ with } z_i \sim \text{Bernoulli}(p)$$

$$\tilde{y}^{(l)} = z^{(l)} \odot y^{(l)} = (z_1^{(l)} y_1^{(l)}, z_2^{(l)} y_2^{(l)}, \dots)$$

$$x^{(l+1)} = W^{(l+1)} \tilde{y}^{(l)} + b^{(l+1)} = W^{(l+1)} \text{diag}(z^{(l)}) y^{(l)} + b^{(l+1)}$$

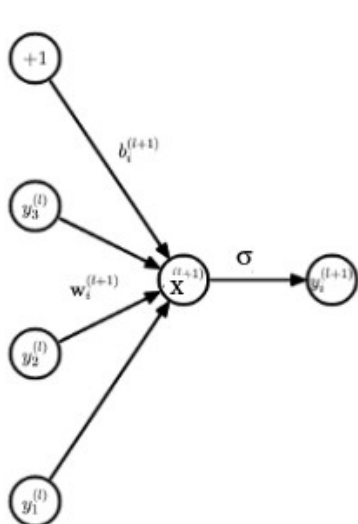
$$y^{(l+1)} = \sigma(x^{(l+1)})$$

$\tilde{W}_{\alpha\beta}^{(l+1)} = W_{\alpha\beta}^{(l+1)} z_\beta^{(l)} = 0$ if $z_\beta = 0$ **column β vanishes**

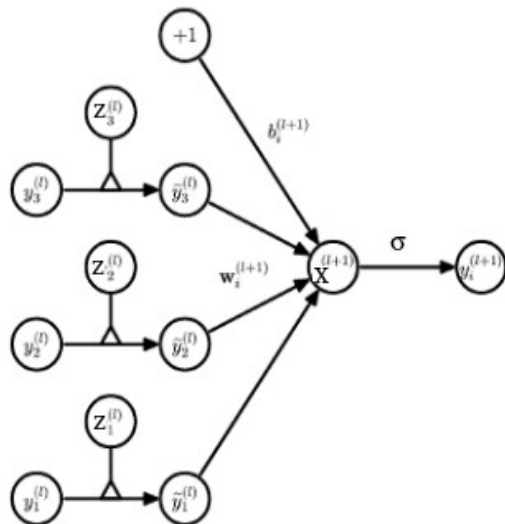
$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



(a) Standard network



(b) Dropout network

Dropout interpretation: Ensemble model

| | \tilde{y}_1 | \tilde{y}_2 | \tilde{y}_3 | |
|--------------|---------------|---------------|---------------|-------------|
| train_iter_1 | 1 | 1 | 1 | Model 1 |
| train_iter_2 | 1 | 1 | 0 | Model 2 |
| | | ⋮ | | |
| | 0 | 0 | 0 | Model 2^N |

Obviously this interpretation is an approximation as the models are not really independent !

$z^{(l)} = (z_1, z_2, \dots)$ with $z_i \sim \text{Bernoulli}(p)$

$\tilde{y}^{(l)} = z^{(l)} \odot y^{(l)} = (z_1^{(l)} y_1^{(l)}, z_2^{(l)} y_2^{(l)}, \dots)$

$$x^{(l+1)} = W^{(l+1)} \tilde{y}^{(l)} + b^{(l+1)} = W^{(l+1)} \text{diag}(z^{(l)}) y^{(l)} + b^{(l+1)}$$

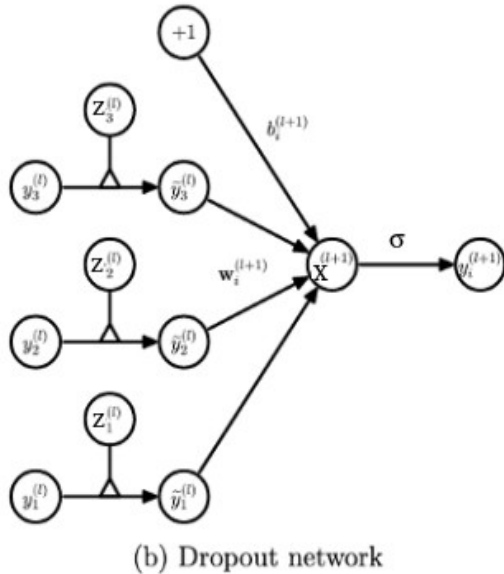
$$y^{(l+1)} = \sigma(x^{(l+1)})$$

$\tilde{W}_{\alpha\beta}^{(l+1)} = W_{\alpha\beta}^{(l+1)} z_{\beta}^{(l)} = 0$ if $z_{\beta} = 0$ **column β vanishes**

$$x^{(l+1)} = W^{(l+1)} y^{(l)} + b^{(l+1)}$$

$$y^{(l+1)} = \sigma(x^{(l+1)})$$

Original dropout method



In summary (for a 2-layer DNN): Data : $D = (X, Y) = \{(x_i, y_i) \mid i = 1, \dots, N\}$
 $x_i \sim 1 \times Q$; $y_i \sim 1 \times D$

$$\hat{y} = W^{(2)} \text{diag}(z^{(2)}) \sigma(W^{(1)} \text{diag}(z^{(1)})x + b) \quad W^{(1)} \sim Q \times K; \quad W^{(2)} \sim D \times K$$

$z^{(l)} = (z_1^{(l)}, z_2^{(l)}, \dots, z_Q^{(l)})$ with $z_j^{(l)} \sim \text{Bernoulli}(p)$

$$E = \frac{1}{2N} \sum_{n=1}^N \|y_n - \hat{y}_n\|$$

Training:

MLE with sampling $z^{(l)} \sim \text{Bern.}(p)$

Inference:

$W \rightarrow pW$

$$L_{\text{dropout}} = E + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$$

Using the Bayesian approach the ensemble picture would be more clear and we will recover the equations above and more !! .

2 - Gaussian Process for DNN

Gaussian process: short introduction



Given some observed data: $D = \{X, Y\}$

$$P(\hat{y}/D, \hat{x}) = \int df P(\hat{y}/f, \hat{x}, X) P(f/D) \rightarrow \text{Very hard to compute in general. !!}$$

Gaussian process: short introduction

Given some observed data: $D = \{X, Y\}$

$$P(\hat{y}/D, \hat{x}) = \int df P(\hat{y}/f, \hat{x}, X) \boxed{P(f/D)} \rightarrow \text{Very hard to compute in gen. !!}$$

In a Bayesian approach we use the **Bayes theorem** to get the posterior :

$$P(f/D) = \frac{\overbrace{P(Y/f, X)}^{\text{Likelihood}} \overbrace{P(f/X)}^{\text{Prior}}}{\underbrace{P(Y/X)}} \rightarrow N(y, f, \sigma^2)$$

GP with kernel K

Evidence or Normalization

$$P(Y/X) = \int df P(Y, f/X) = \int df P(Y/f, X) P(f/X)$$

Gaussian process: short introduction

Given some observed data: $D = \{X, Y\}$

$$P(\hat{y}/D, \hat{x}) = \int df P(\hat{y}/f, \hat{x}, X) \boxed{P(f/D)} \rightarrow \text{Very hard to compute in geral. !!}$$

In a Bayesian approach we use the **Bayes theorem** to get the posterior :

$$P(f/D) = \frac{\overbrace{P(Y/f, X)}^{\text{Likelihood}} \overbrace{P(f/X)}^{\text{Prior}}}{\underbrace{P(Y/X)}} \rightarrow N(y, f, \sigma^2) \quad \text{GP with kernel K}$$

Evidence or Normalization

$$P(Y/X) = \int df P(Y, f/X) = \int df P(Y/f, X) P(f/X)$$

Kernels: include info about family of functions being approximated, i.e. periodic functions, smooth functions, stochastic functions, etc. as well as information about the **confidence intervals**. See [link](#) .

Gaussian process: Bayesian DNN connection

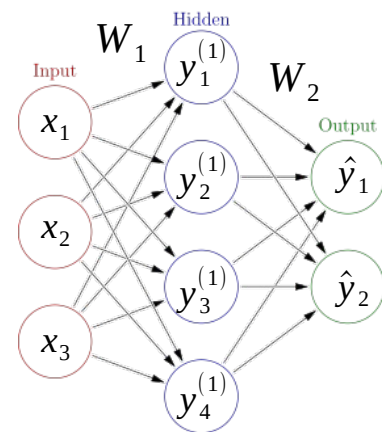


I will focus in a two layer DNN for simplicity. Suppose we have this kernel:

$$K(x, x') = \int p(w) p(b) \sigma(\mathbf{w}^T \mathbf{x} + b)^T \sigma(\mathbf{w}^T \mathbf{x}' + b) dw db$$

$w, x \sim Q \times 1 \quad b \in \mathbb{R}$

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

I will focus in a two layer DNN for simplicity. Suppose we have this kernel:

$$K(x, x') = \int p(w) p(b) \sigma(\mathbf{w}^T \mathbf{x} + b)^T \sigma(\mathbf{w}^T \mathbf{x}' + b) dw db$$

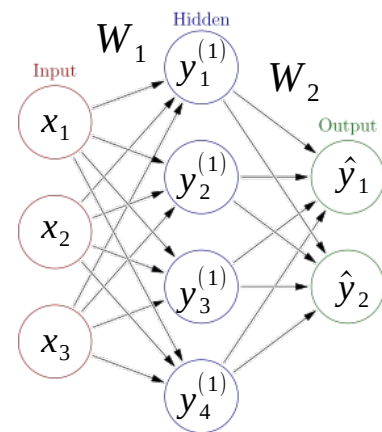
$w, \mathbf{x} \sim Q \times 1 \quad b \in \mathbb{R}$

Approximation: $\int \rightarrow \sum$

$$K(x, x') \simeq \sum_{k=1}^K \sqrt{\frac{1}{K}} \sigma(\mathbf{w}_k^T \mathbf{x} + b_k)^T \sqrt{\frac{1}{K}} \sigma(\mathbf{w}_k^T \mathbf{x}' + b_k)$$

$w_k, b_k \sim p(w), p(b)$

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

I will focus in a two layer DNN for simplicity. Suppose we have this kernel:

$$K(x, x') = \int p(w) p(b) \sigma(\mathbf{w}^T \mathbf{x} + b)^T \sigma(\mathbf{w}^T \mathbf{x}' + b) dw db$$

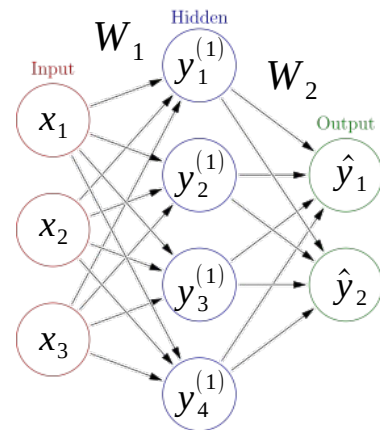
$w, \mathbf{x} \sim Q \times 1 \quad b \in \mathbb{R}$

Approximation: $\int \rightarrow \sum$

$$K(x, x') \simeq \sum_{k=1}^K \sqrt{\frac{1}{K}} \sigma(\mathbf{w}_k^T \mathbf{x} + b_k)^T \sqrt{\frac{1}{K}} \sigma(\mathbf{w}_k^T \mathbf{x}' + b_k) = \sqrt{\frac{1}{K}} \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b})^T \sqrt{\frac{1}{K}} \sigma(\mathbf{W}_1 \mathbf{x}' + \mathbf{b})$$

$w_k, b_k \sim p(w), p(b) \quad \mathbf{W}_1 \sim K \times Q$

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

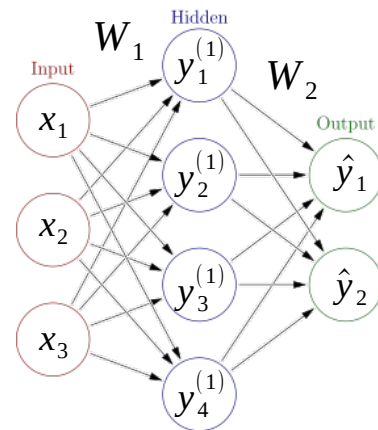
I will focus in a two layer DNN for simplicity. Suppose we have this kernel:

$$K(x, x') = \int p(w) p(b) \sigma(\mathbf{w}^T \mathbf{x} + b)^T \sigma(\mathbf{w}^T \mathbf{x}' + b) dw db$$

$w, \mathbf{x} \sim Q \times 1 \quad b \in \mathbb{R}$

Approximation: $\int \rightarrow \sum$

σ : ReLU, sigmoid, etc.



$$K(x, x') \simeq \sum_{k=1}^K \sqrt{\frac{1}{K}} \sigma(\mathbf{w}_k^T \mathbf{x} + b_k)^T \sqrt{\frac{1}{K}} \sigma(\mathbf{w}_k^T \mathbf{x}' + b_k) = \sqrt{\frac{1}{K}} \sigma(\mathbf{W}_1 \mathbf{x} + b)^T \sqrt{\frac{1}{K}} \sigma(\mathbf{W}_1 \mathbf{x}' + b)$$

$w_k, b_k \sim p(w), p(b) \quad \mathbf{W}_1 \sim K \times Q$

Gaussian process: Bayesian DNN connection

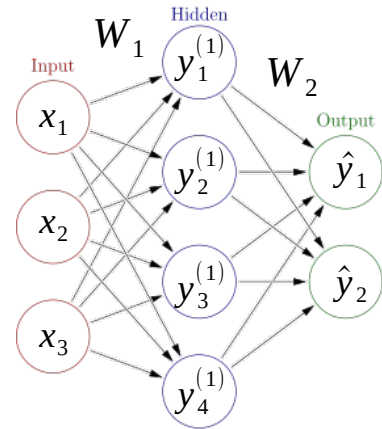


$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

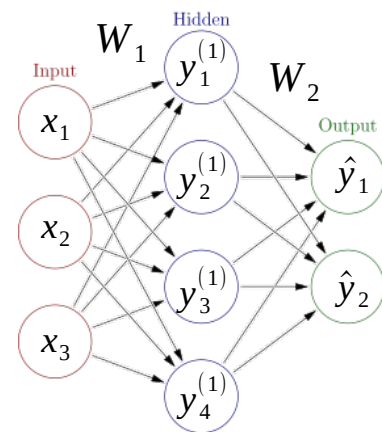
$$W_1 \sim K \times Q; \quad x \sim Q \times 1; \quad y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; \quad x \sim Q \times 1; \quad y^{(1)} \sim K \times 1$$

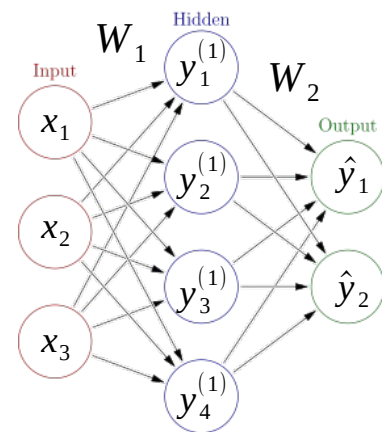
with W_1 and b random variables as usual in a Bayesian approach

To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$$P(Y/f) \sim N(Y, f, \tau^{-1} I_D) \text{ Likelihood} \quad \rightarrow \quad P(f/W_1, b, X) \sim GP(0, K)$$

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

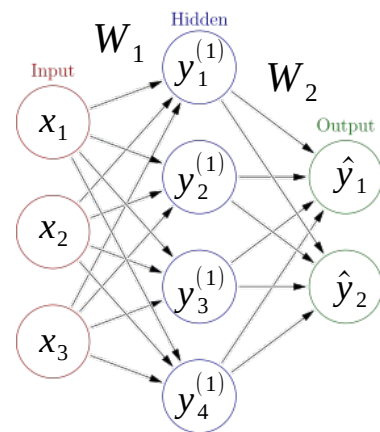
To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$$P(Y/f) \sim N(Y, f, \tau^{-1} I_D) \text{ Likelihood} \quad P(f/W_1, b, X) \sim GP(0, K)$$

$$\int df P(Y/f) P(f/W_1, b, X) = N(Y; 0, (K(X, X) + \tau^{-1}) I_D)$$

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

To find the connection between a GP and DNN's let's make use of the evidence:

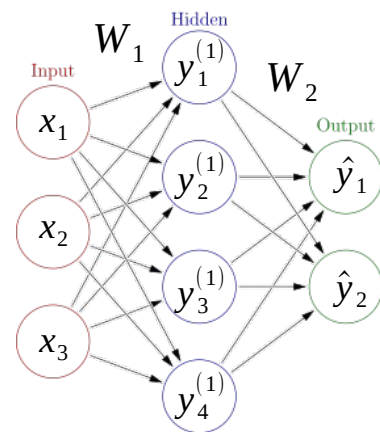
$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$$P(Y/f) \sim N(Y, f, \tau^{-1} I_D) \text{ Likelihood} \quad P(f/W_1, b, X) \sim GP(0, K)$$

$$\int df P(Y/f) P(f/W_1, b, X) = N(Y; 0, (K(X, X) + \tau^{-1} I_D)) = \int N(Y; W_2 y^{(1)}, \tau^{-1} I_D) P(W_2) dW_2$$

auxiliar matrix variables : W_2

σ : ReLU, sigmoid, etc.



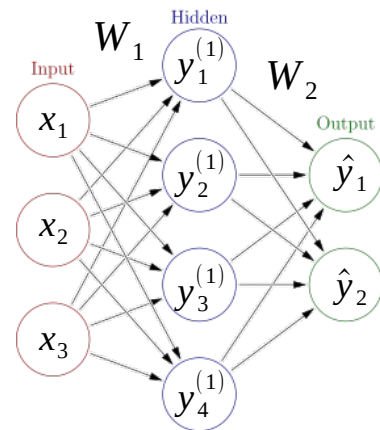
Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

σ : ReLU, sigmoid, etc.



To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$P(Y/f) \sim N(Y, f, \tau^{-1} I_D)$ Likelihood \rightarrow $P(f/W_1, b, X) \sim GP(0, K)$

$$\int df P(Y/f) P(f/W_1, b, X) = N(Y; 0, (K(X, X) + \tau^{-1} I_D)) = \int N(Y; W_2 y^{(1)}, \tau^{-1} I_D) P(W_2) dW_2$$

auxiliar matrix variables : W_2

$$f(X, W_1, W_2, b) = \hat{y} = W_2 \sigma(W_1 X + b)$$

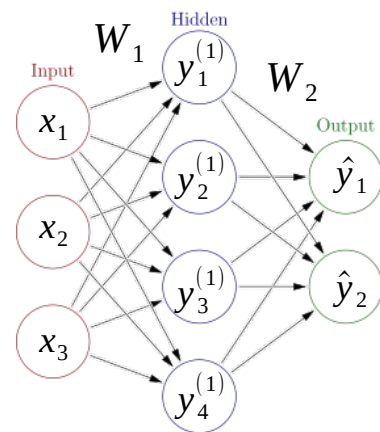
Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

σ : ReLU, sigmoid, etc.



To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$P(Y/f) \sim N(Y, f, \tau^{-1} I_D)$ Likelihood \rightarrow $P(f/W_1, b, X) \sim GP(0, K)$

$$\int df P(Y/f) P(f/W_1, b, X) = N(Y; 0, (K(X, X) + \tau^{-1} I_D)) = \int N(Y; W_2 y^{(1)}, \tau^{-1} I_D) P(W_2) dW_2$$

$P(Y/W_1, W_2, b, X)$ likelihood in parameter space

auxiliar matrix variables : W_2

$$f(X, W_1, W_2, b) = \hat{y} = W_2 \sigma(W_1 X + b)$$

Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; \quad x \sim Q \times 1; \quad y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

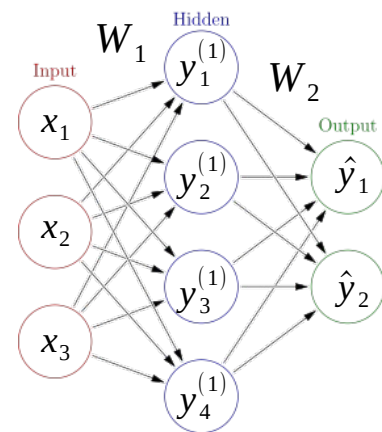
To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$$= \int P(Y/W_1, W_2, b, X) P(W_1) P(W_2) P(b) dW_1 dW_2 db$$

Bayesian parametric representation of our 2 layer DNN !!

σ : ReLU, sigmoid, etc.



Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

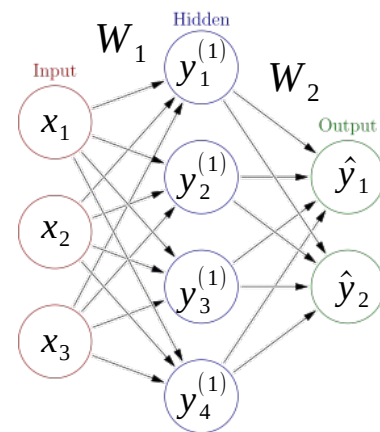
To find the connection between a GP and DNN's let's make use of the evidence:

$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$$= \int P(Y/W_1, W_2, b, X) P(W_1) P(W_2) P(b) dW_1 dW_2 db$$

Bayesian parametric representation of our 2 layer DNN !!

σ : ReLU, sigmoid, etc.



Conclusion: $K(x, x')$ right kernel to approximate DNN functions in a GP approach.

Gaussian process: Bayesian DNN connection

$$K(x, x') = \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b) \sim y^{(1)}(W_1, x, b)^T y^{(1)}(W_1, x', b)$$

$$W_1 \sim K \times Q; x \sim Q \times 1; y^{(1)} \sim K \times 1$$

with W_1 and b random variables as usual in a Bayesian approach

To find the connection between a GP and DNN's let's make use of the evidence:

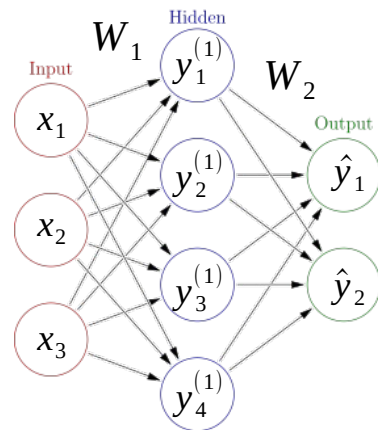
$$P(Y/X) = \int P(Y/f) P(f/W_1, b, X) P(W_1) P(b) df dW_1 db$$

$$= \int P(Y/W_1, W_2, b, X) P(W_1) P(W_2) P(b) dW_1 dW_2 db$$

Bayesian parametric representation of our 2 layer DNN !!

**Non-linearities σ non-trivial
Impact in C.I. of Bayes-DNN**

σ : ReLU, sigmoid, etc.



Conclusion: $K(x, x')$ right kernel to approximate DNN functions in a GP approach.

Note:

For a generalization to deeper DNN and classification problems see ref.2 (Appendix)

But for a more accurate connection see “Deep neural networks as Gaussian processes”, Lee et al ‘2017.

3 - Dropout from a Bayesian point of view

Variational Bayesian Inference and Dropout relationship

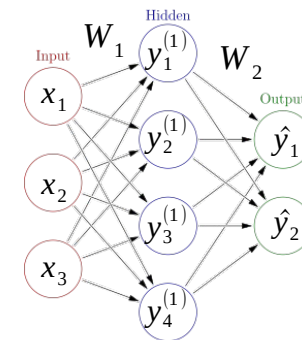


Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

σ : ReLU, sigmoid, etc.



Variational Bayesian Inference and Dropout relationship



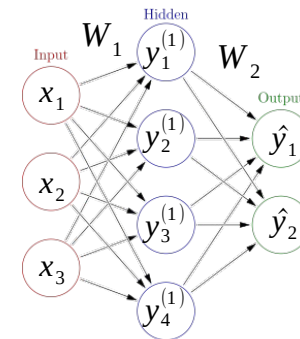
Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

Posterior: quite hard to compute !!

σ : ReLU, sigmoid, etc.



Variational Bayesian Inference and Dropout relationship



Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

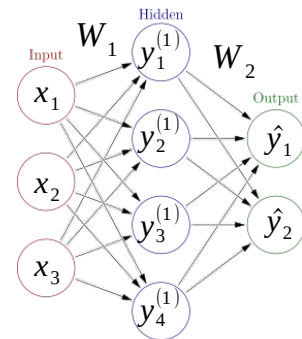
Posterior: quite hard to compute !!

$$P(W_1, W_2, b/X, Y) = \frac{P(Y/W_1, W_2, b, X) P(W_1, W_2, b/X)}{P(Y/X)}$$

$$P(Y/X) = \int P(Y/W_1, W_2, b, X) P(W_1) P(W_2) P(b) dW_1 dW_2 db$$

Intractable for large DNN !!!

σ : ReLU, sigmoid, etc.



Variational Bayesian Inference and Dropout relationship



Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

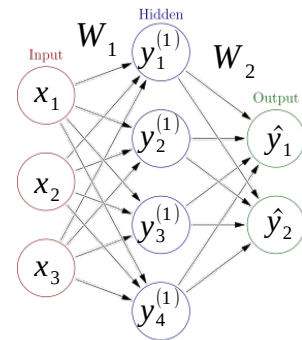
Posterior: quite hard to compute !!

$$P(W_1, W_2, b/X, Y) = \frac{P(Y/W_1, W_2, b, X) P(W_1, W_2, b/X)}{P(Y/X)}$$

$$P(Y/X) = \int P(Y/W_1, W_2, b, X) P(W_1) P(W_2) P(b) dW_1 dW_2 db$$

Intractable for large DNN !!!

σ : ReLU, sigmoid, etc.



Variational Inference: Approximation of the posterior by an ansatz distribution.

Variational Bayesian Inference and Dropout relationship



Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

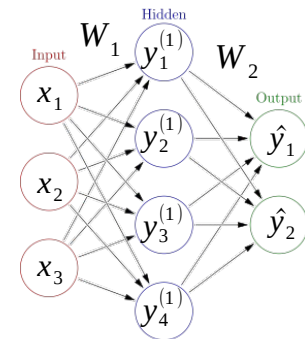
$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

Posterior: quite hard to compute !!

$$P(W_1, W_2, b/X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$$

σ : ReLU, sigmoid, etc.



Variational Bayesian Inference and Dropout relationship



Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

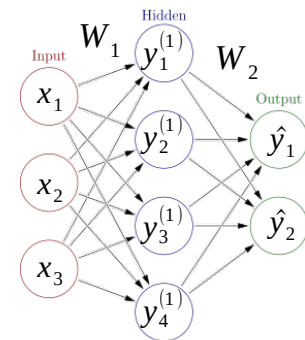
$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

Posterior: quite hard to compute !!

$$P(W_1, W_2, b/X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$$

σ : ReLU, sigmoid, etc.



$$q_M(W) = \prod_{\alpha} q_{m_{\alpha}}(w_{\alpha}) \text{ with } w_{\alpha}/m_{\alpha} \text{ the columns of } W/M$$

$$q_{m_{\alpha}}(w_{\alpha}) = p N(m_{\alpha}, \theta^2 I) + (1-p) * N(0, \theta^2 I)$$

$$q(b) = N(m, \theta^2 I)$$

Variational Bayesian Inference and Dropout relationship



Applying Bayes theorem in a parametric space the predictive probability distribution of a DNN is given by:

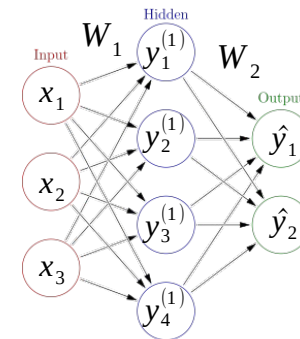
$$P(\hat{y}/\hat{x}, X, Y) = \int P(\hat{y}/\hat{x}, W_1, W_2, b) P(W_1, W_2, b/X, Y) dW_1 dW_2 db$$

$$N(\hat{y}; f = W_2 \sigma(W_1 \hat{x} + b), I)$$

Posterior: quite hard to compute !!

$$P(W_1, W_2, b/X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$$

σ : ReLU, sigmoid, etc.



$$q_M(W) = \prod_{\alpha} q_{m_{\alpha}}(w_{\alpha}) \text{ with } w_{\alpha}/m_{\alpha} \text{ the columns of } W/M$$

$$q_{m_{\alpha}}(w_{\alpha}) = p N(m_{\alpha}, \theta^2 I) + (1-p) * N(0, \theta^2 I)$$

$$q(b) = N(m, \theta^2 I)$$

Probability version of standard dropout approach !!

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

Kullback – Leibler divergence: $\omega \stackrel{\text{def}}{=} (W_1, W_2, b)$

$$KL(q_M(\omega) | P(\omega | D)) = \int q_M(\omega) \ln [P(D, \omega) / q_M(\omega)] d\omega - \ln P(D) \quad \longrightarrow \quad ELBO(q_M(\omega)) \leq \ln P(D)$$

≥ 0 $ELBO(q_M(\omega))$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

Kullback – Leibler divergence: $\omega \stackrel{\text{def}}{=} (W_1, W_2, b)$

$$KL(q_M(\omega) | P(\omega/D)) = \int q_M(\omega) \ln [P(D, \omega) / q_M(\omega)] d\omega - \ln P(D) \quad \longrightarrow \quad ELBO(q_M(\omega)) \leq \ln P(D)$$

≥ 0 $ELBO(q_M(\omega))$

$$\text{Max}_M ELBO(q_M(\omega)) \longrightarrow q_{\tilde{M}}(\omega) \rightarrow P(D/\omega)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1)q_{M_2}(W_2)q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\underbrace{\ln P(D | W_1, W_2, b)}_{\text{Likelihood}}] - \underbrace{KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))}_{\text{prior}}$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] = \int \ln P(D | W_1, W_2, b) q_M(W_1, W_2, b) dW_1 dW_2 db$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] = \int \ln P(D | W_1, W_2, b) q_M(W_1, W_2, b) dW_1 dW_2 db$$

$$\ln P(D | W_1, W_2, b) = \ln N(Y; \hat{Y} = f(X, W_1, W_2, b), \tau^{-1} I_D) \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\| \quad \text{with} \quad \hat{y}_n = f(x_n, W_1, W_2, b)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] = \int \ln P(D | W_1, W_2, b) q_M(W_1, W_2, b) dW_1 dW_2 db$$

$$\ln P(D | W_1, W_2, b) = \ln N(Y; \hat{Y} = f(X, W_1, W_2, b), \tau^{-1} I_D) \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\| \quad \text{with} \quad \hat{y}_n = f(x_n, W_1, W_2, b)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] = \frac{\mathcal{I}}{2} \sum_{n=1}^N \int \|y_n - f(x_n, W_1, W_2, b)\| \times q_M(W_1, W_2, b) dW_1 dW_2 db$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] = \frac{\mathcal{I}}{2} \sum_{n=1}^N \int \|y_n - f(x_n, W_1, W_2, b)\| \times q_M(W_1, W_2, b) dW_1 dW_2 db$$

$\int \rightarrow \Sigma$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] \sim \frac{\mathcal{I}}{2} \sum_{n=1}^N \sum_{\alpha=1}^M \|y_n - f(x_n, W_1^\alpha, W_2^\alpha, b^\alpha)\| \quad W_1^\alpha, W_2^\alpha, b^\alpha \sim q_{\tilde{M}}(W_1, W_2, b)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] \sim \frac{\mathcal{I}}{2} \sum_{n=1}^N \sum_{\alpha=1}^M \|y_n - f(x_n, W_1^\alpha, W_2^\alpha, b^\alpha)\| \quad W_1^\alpha, W_2^\alpha, b^\alpha \sim q_{\tilde{M}}(W_1, W_2, b)$$

for $N \gg 1$, $\sum_n \sum_\alpha \rightarrow \sum_n$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] \sim \frac{\mathcal{I}}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \quad W_1^n, W_2^n, b^n \sim q_{\tilde{M}}(W_1, W_2, b)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D | W_1, W_2, b)] \sim \frac{\mathcal{I}}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \quad W_1^n, W_2^n, b^n \sim q_M(W_1, W_2, b)$$

A bit of reparameterization:

$$W_1 = \text{diag}(z_1)(M_1 + \theta \epsilon_1) + (\text{Id} - \text{diag}(z_1)) \theta \epsilon_1 \quad (\text{same for } W_2)$$

$$b = m + \theta \epsilon \quad \text{with } z_i \sim \text{Bernoulli}(p_i) \text{ and } \epsilon_i \sim N(0, I)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\mathcal{I}}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \quad W_1^n, W_2^n, b^n \sim q_M(W_1, W_2, b)$$

A bit of reparameterization:

$$W_1 = \text{diag}(z_1)(M_1 + \theta \epsilon_1) + (\text{Id} - \text{diag}(z_1)) \theta \epsilon_1 \quad (\text{same for } W_2)$$

$$b = m + \theta \epsilon \quad \text{with } z_i \sim \text{Bernoulli}(p_i) \text{ and } \epsilon_i \sim N(0, I)$$

$$\theta \rightarrow 0$$

$$W_{1,2} \simeq \text{diag}(z_{1,2}) M_{1,2}$$

$$b \simeq m$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\mathcal{I}}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \quad W_1^n, W_2^n, b^n \sim q_M(W_1, W_2, b)$$

A bit of reparameterization:

$$W_1 = \text{diag}(z_1)(M_1 + \theta \epsilon_1) + (\text{Id} - \text{diag}(z_1)) \theta \epsilon_1 \quad (\text{same for } W_2)$$

$$b = m + \theta \epsilon \quad \text{with } z_i \sim \text{Bernoulli}(p_i) \text{ and } \epsilon_i \sim N(0, I)$$

$$\longrightarrow f(x_n, W_1^n, W_2^n, b^n) = \hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$$

$$\theta \rightarrow 0$$

$$W_{1,2} \simeq \text{diag}(z_{1,2}) M_{1,2}$$

$$b \simeq m$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\|$$

$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$

$\theta \rightarrow 0$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\|$$

$\theta \rightarrow 0$

First term of dropout MLE loss

$$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\|$$

$\theta \rightarrow 0$

First term of dropout MLE loss

$$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - KL(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\|$$

$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)})x + m)$
 $\theta \rightarrow 0$
First term of dropout MLE loss

$$KL(q_M(W_1, W_2, b) | P(W_1, W_2, b)) \sim -\frac{p_1}{2} \|M_1\| - \frac{p_2}{2} \|M_2\| - \frac{1}{2} \|m\|$$

$\theta \rightarrow 0$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D/W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

$$ELBO(q_M(W_1, W_2, b)) = E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] - \text{KL}(q_M(W_1, W_2, b) | P(W_1, W_2, b))$$

$$E_{W_1, W_2, b \sim q_M(W_1, W_2, b)} [\ln P(D/W_1, W_2, b)] \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - f(x_n, W_1^n, W_2^n, b^n)\| \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\|$$

$$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$$

$\theta \rightarrow 0$

First term of dropout MLE loss

$$\text{KL}(q_M(W_1, W_2, b) | P(W_1, W_2, b)) \sim -\frac{p_1}{2} \|M_1\| - \frac{p_2}{2} \|M_2\| - \frac{1}{2} \|m\|$$

$\theta \rightarrow 0$

Regularization terms of dropout MLE loss

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \longrightarrow q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

Loss function of standard dropout!!!

$$ELBO(q_M(W_1, W_2, b)) \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\| - \frac{p_1}{2} \|M_1\| - \frac{p_2}{2} \|M_2\| - \frac{1}{2} \|m\|$$

$$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$$

Variational Bayesian Inference and Dropout relationship



Let's find the M optimal parameters: $P(W_1, W_2, b | X, Y) \sim q_M(W_1, W_2, b) = q_{M_1}(W_1) q_{M_2}(W_2) q_m(b)$

$$\text{Max}_M ELBO(q_M(W_1, W_2, b)) \quad \longrightarrow \quad q_{\tilde{M}}(W_1, W_2, b) \rightarrow P(D | W_1, W_2, b)$$

From ELBO definition: (spoiler: $ELBO \sim L_{dropout} = 1/(2N) \sum_{n=1}^N \|y_n - \hat{y}_n\| + \lambda_1 \|W_1\| + \lambda_2 \|W_2\| + \lambda_3 \|b\|$)

Loss function of standard dropout!!!

$$ELBO(q_M(W_1, W_2, b)) \sim \frac{\tau}{2} \sum_{n=1}^N \|y_n - \hat{y}_n\| - \frac{p_1}{2} \|M_1\| - \frac{p_2}{2} \|M_2\| - \frac{1}{2} \|m\|$$

$$\hat{y}_n = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)}) x + m)$$

We will maximize the ELBO through standard MLE methods (Gradient descent, etc)

Inference: let's compute the predictions and variances



$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

Inference: let's compute the predictions and variances



$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

Inference: let's compute the predictions and variances



$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

$$E(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \int y^* P(y^*/x^*, W_{1,t}, W_{2,t}, b_t) * d y^*$$

$$W_1 = \text{diag}(z_1) M_1$$

$$W_2 = \text{diag}(z_2) M_2$$

$$b = m$$

Inference: let's compute the predictions and variances



$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

$$E(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \int y^* P(y^*/x^*, W_{1,t}, W_{2,t}, b_t) * dy^* = \frac{1}{T} \sum_{t=1}^T \int y^* N(y^*, \hat{y}^*(x^*, z_{1,t}, z_{2,t})) * dy^*$$

$\hat{y}^* = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)})x^* + m)$

$$W_1 = \text{diag}(z_1) M_1$$

$$W_2 = \text{diag}(z_2) M_2$$

$$b = m$$

Inference: let's compute the predictions and variances

$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

$$E(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \int y^* P(y^*/x^*, W_{1,t}, W_{2,t}, b_t) * dy^* = \frac{1}{T} \sum_{t=1}^T \int y^* N(y^*, \hat{y}^*(x^*, z_{1,t}, z_{2,t})) * dy^*$$

$\hat{y}^* = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)})x^* + m)$

$= \hat{y}^*(x^*, z_{1,t}, z_{2,t})$

$$W_1 = \text{diag}(z_1) M_1$$

$$W_2 = \text{diag}(z_2) M_2$$

$$b = m$$

Inference: let's compute the predictions and variances



$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

$$E(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \int y^* P(y^*/x^*, W_{1,t}, W_{2,t}, b_t) * dy^* = \frac{1}{T} \sum_{t=1}^T \int y^* N(y^*, \hat{y}^*(x^*, z_{1,t}, z_{2,t})) * dy^* = \frac{1}{T} \sum_{t=1}^T \hat{y}^*(x^*, z_{1,t}, z_{2,t})$$

$= \hat{y}^*(x^*, z_{1,t}, z_{2,t})$

$\hat{y}^* = M^{(2)} \text{diag}(z^{(2)}) \sigma(M^{(1)} \text{diag}(z^{(1)})x^* + m)$

$$W_1 = \text{diag}(z_1) M_1$$

$$W_2 = \text{diag}(z_2) M_2$$

$$b = m$$

Inference: let's compute the predictions and variances



$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

$$E(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}^*(x^*, z_{1,t}, z_{2,t})$$

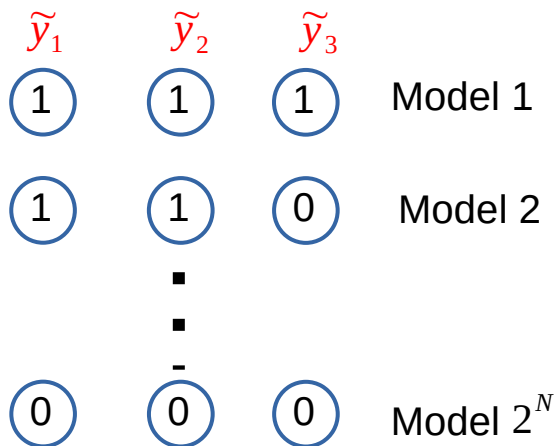
Inference: let's compute the predictions and variances

$$E(y^*) = \int y^* P(y^*/x^*, X, Y) dy^* = \int y^* P(y^*/x^*, W_1, W_2, b) * P(W_1, W_2, b/X, Y) dW_1 dW_2 db dy^*$$

$$P(W_1, W_2, b/X, Y) \rightarrow q_M(W_1, W_2, b)$$

$$E(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}^*(x^*, z_{1,t}, z_{2,t})$$

Dropout interpretation: Ensemble model



In agreement with our initial interpretation of ensemble model or model averaging !!

Inference: let's compute the predictions and variances



MC dropout method

$$E_{q_M(y^*/x^*)}(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}^*(x^*, z_1^t, z_2^t, \dots) \quad \text{Mean}$$

$$\text{Var}_{q_M(y^*/x^*)}(y^*) \simeq \tau^{-1} I_D + \frac{1}{T} \sum_{t=1}^T \hat{y}^*(x^*, z_1^t, z_2^t, \dots)^T \hat{y}^*(x^*, z_1^t, z_2^t, \dots) - E_{q_M(y^*/x^*)}(y^*)^T E_{q_M(y^*/x^*)}(y^*) \quad \text{Variance}$$

$$\hat{y}^*(x^*, z_1, z_2, \dots) = (M_L \text{diag}(z_L)) \sigma(\dots (M_2 \text{diag}(z_2)) \sigma((M_1 \text{diag}(z_1)) x^* + m_1)) \quad z_1, z_2 \sim \text{Bern}(p_1), \text{Bern}(p_2)$$

4 - Results

Regression



(a),(c) and (d): DNN with 4 layers and 1024 hidden units – $p \sim 0.2$

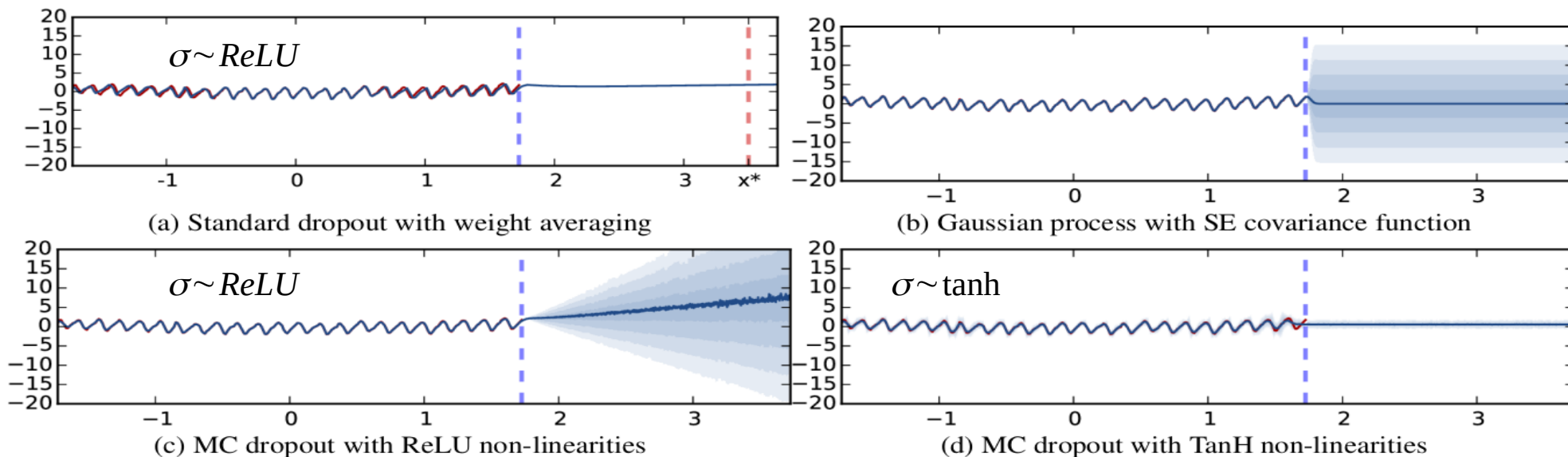


Figure 2. Predictive mean and uncertainties on the Mauna Loa CO₂ concentrations dataset, for various models. In red is the observed function (left of the dashed blue line); in blue is the predictive mean plus/minus two standard deviations (8 for fig. 2d). Different shades of blue represent half a standard deviation. Marked with a dashed red line is a point far away from the data: standard dropout confidently predicts an insensible value for the point; the other models predict insensible values as well but with the additional information that the models are uncertain about their predictions.

Regression



(a),(c) and (d): DNN with 4 layers and 1024 hidden units – $p \sim 0.2$

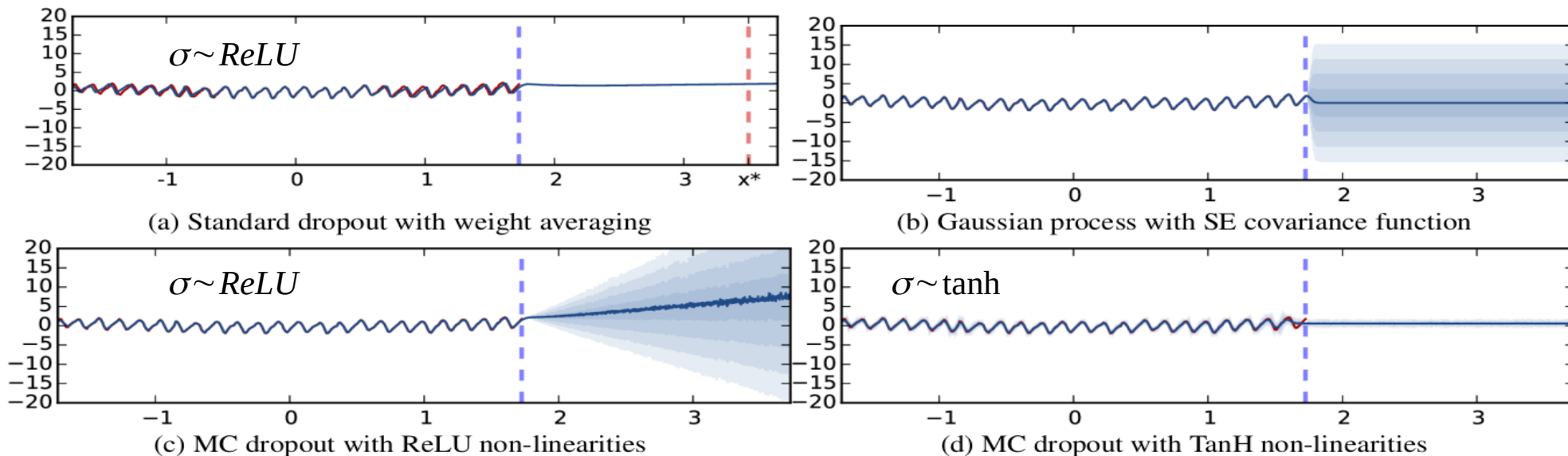
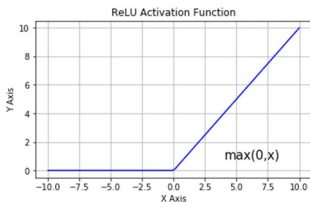
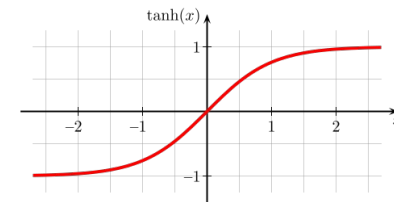


Figure 2. Predictive mean and uncertainties on the Mauna Loa CO₂ concentrations dataset, for various models. In red is the observed function (left of the dashed blue line); in blue is the predictive mean plus/minus two standard deviations (8 for fig. 2d). Different shades of blue represent half a standard deviation. Marked with a dashed red line is a point far away from the data: standard dropout confidently predicts an insensible value for the point; the other models predict insensible values as well but with the additional information that the models are uncertain about their predictions.



Variance $\rightarrow K(x, x') \simeq \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b)$



(a),(c) and (d): DNN with 4 layers and 1024 hidden units – $p \sim 0.2$

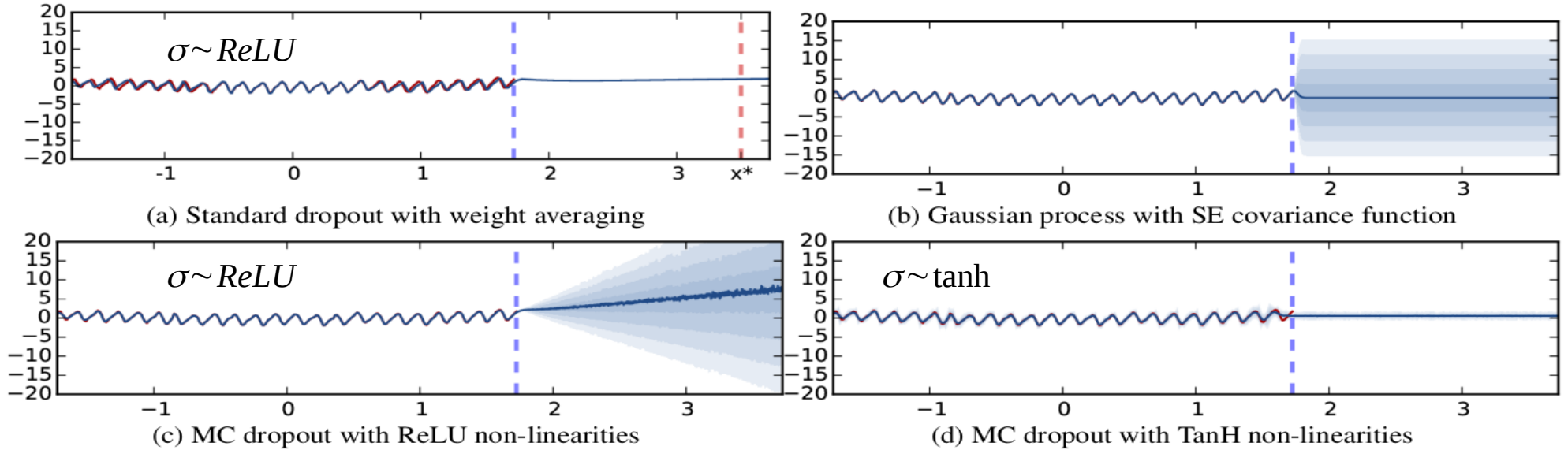
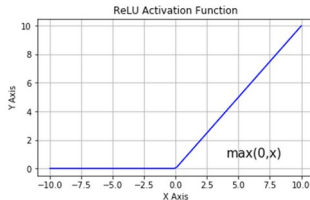
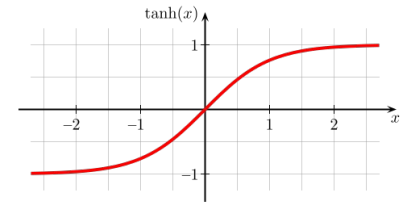


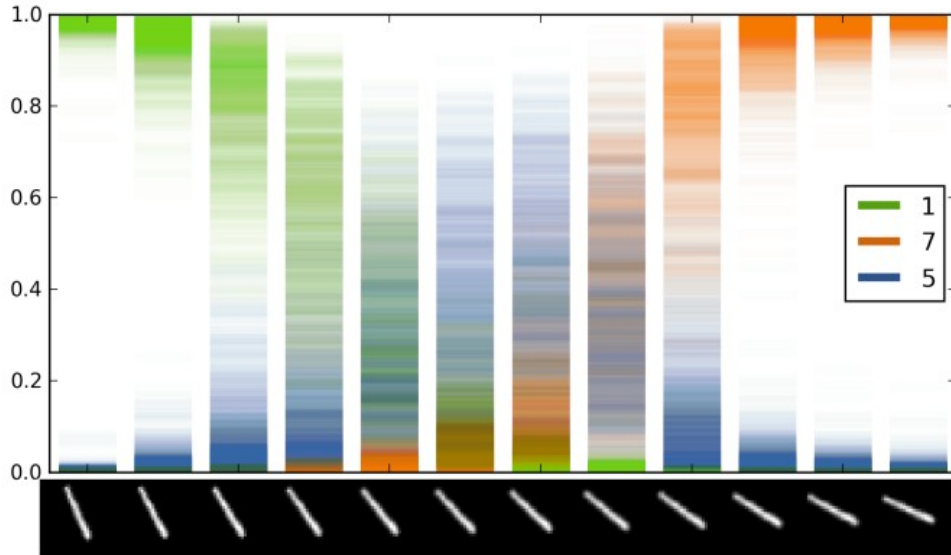
Figure 2. Predictive mean and uncertainties on the Mauna Loa CO₂ concentrations dataset, for various models. In red is the observed function (left of the dashed blue line); in blue is the predictive mean plus/minus two standard deviations (8 for fig. 2d). Different shades of blue represent half a standard deviation. Marked with a dashed red line is a point far away from the data: standard dropout confidently predicts an insensible value for the point; the other models predict insensible values as well but with the additional information that the models are uncertain about their predictions.



Variance $\rightarrow K(x, x') \simeq \sqrt{\frac{1}{K}} \sigma(W_1 x + b)^T \sqrt{\frac{1}{K}} \sigma(W_1 x' + b)$



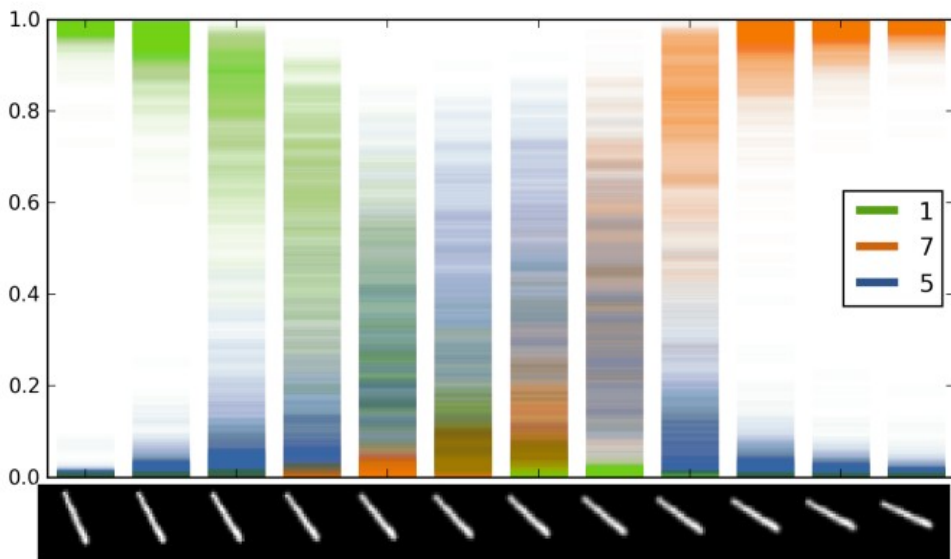
LeNet CNN on MNIST with dropout before last fully connected layer (p~0.5)



$$E_{q_M(y^*/x^*)}(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}_t^*(x^*, z_1^t, z_2^t, \dots) \quad T = 100$$

$$\hat{y}_t^* = (y_1, y_2, y_3, \dots, y_{10})_t$$

LeNet CNN on MNIST with dropout before last fully connected layer ($p \sim 0.5$)



$$E_{q_M(y^*/x^*)}(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}_t^*(x^*, z_1^t, z_2^t, \dots) \quad T = 100$$

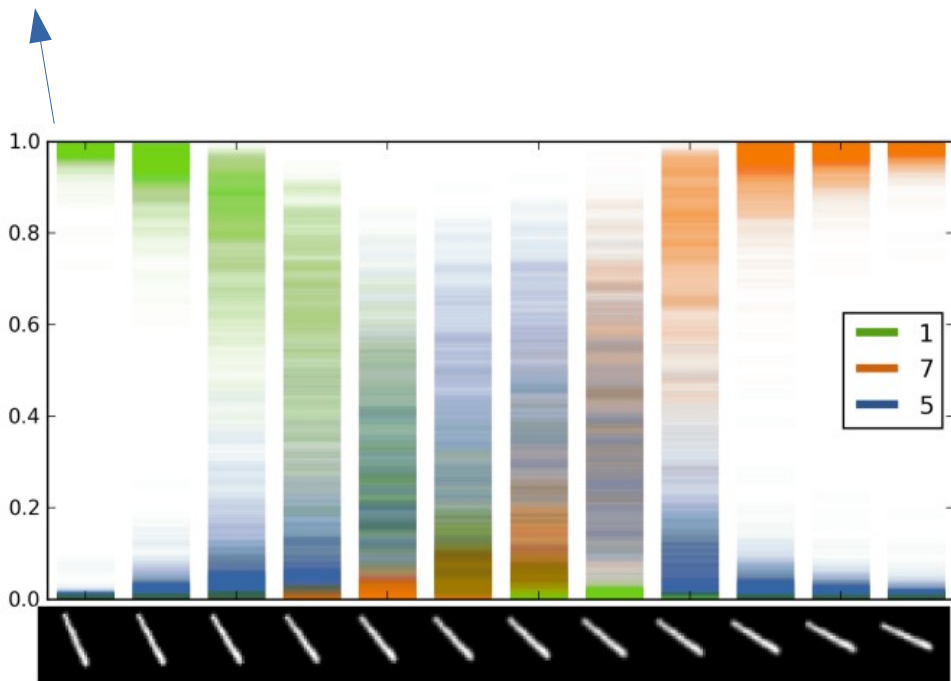
$$\hat{y}_t^* = (y_1, y_2, y_3, \dots, y_{10})_t$$

For every t plot largest three probs. of $(y_1, y_2, y_3, \dots, y_{10})_t$

LeNet CNN on MNIST with dropout before last fully connected layer ($p \sim 0.5$)

$$E(\hat{y}_1) \sim 1, E(\hat{y}_5) \sim 0, E(\hat{y}_7) \sim 0$$

$$\text{Var}(\hat{y}_1) \sim 0, \text{Var}(\hat{y}_5) \sim 0, \text{Var}(\hat{y}_7) \sim 0$$



$$E_{q_M(y^*/x^*)}(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}_t^*(x^*, z_1^t, z_2^t, \dots) \quad T = 100$$

$$\hat{y}_t^* = (y_1, y_2, y_3, \dots, y_{10})_t$$

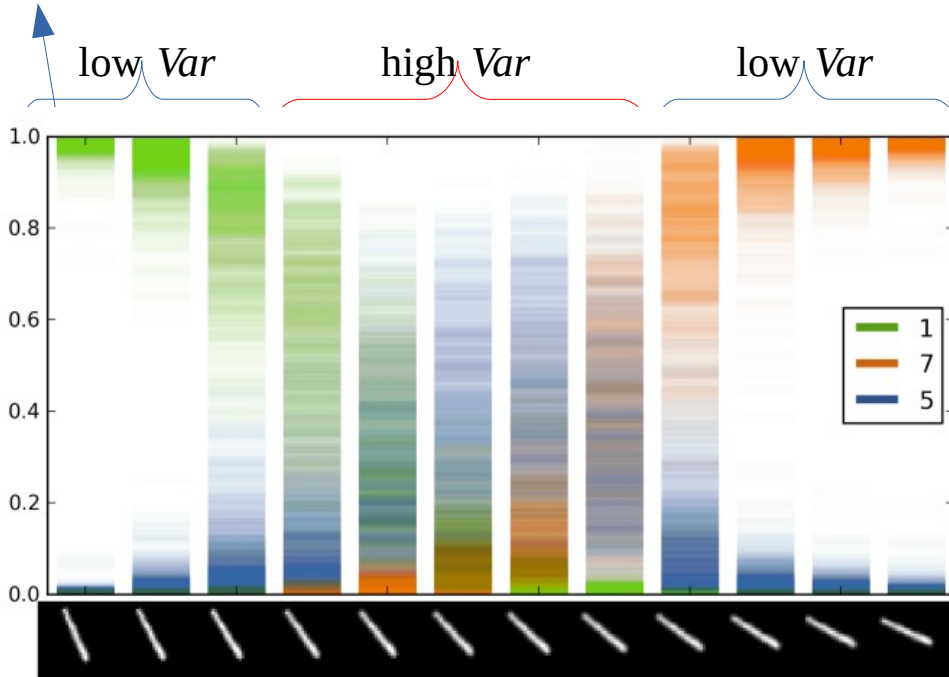
For every t plot largest three probs. of $(y_1, y_2, y_3, \dots, y_{10})_t$

Classification



LeNet CNN on MNIST with dropout before last fully connected layer ($p \sim 0.5$)

$$E(\hat{y}_1) \sim 1, E(\hat{y}_5) \sim 0, E(\hat{y}_7) \sim 0$$
$$\text{Var}(\hat{y}_1) \sim 0, \text{Var}(\hat{y}_5) \sim 0, \text{Var}(\hat{y}_7) \sim 0$$



$$E_{q_M(y^*/x^*)}(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}_t^*(x^*, z_1^t, z_2^t, \dots) \quad T = 100$$
$$\hat{y}_t^* = (y_1, y_2, y_3, \dots, y_{10})_t$$

For every t plot largest three probs. of $(y_1, y_2, y_3, \dots, y_{10})_t$

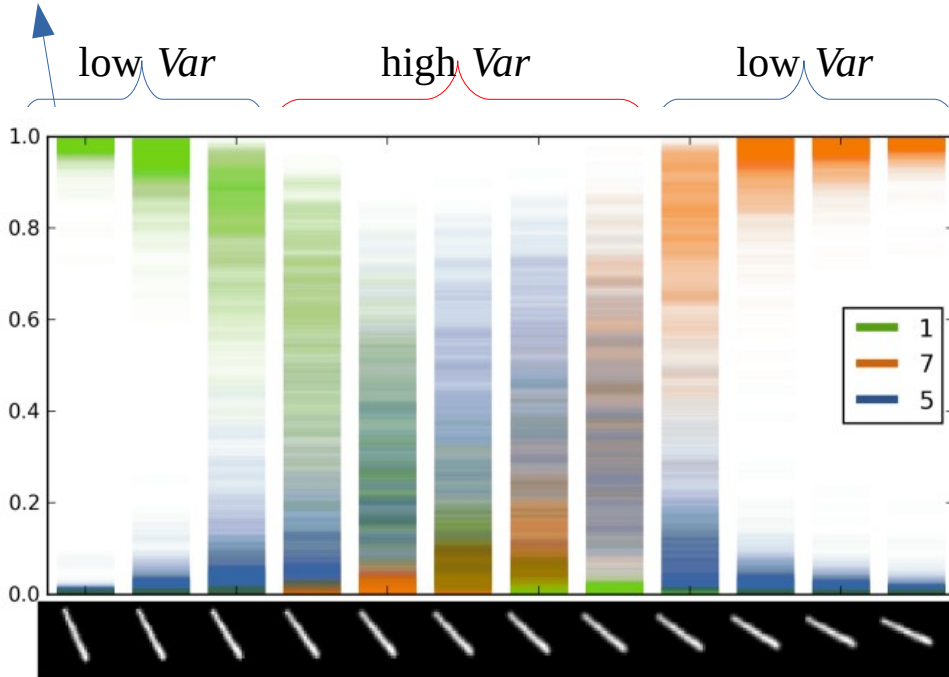
Classification



LeNet CNN on MNIST with dropout before last fully connected layer ($p \sim 0.5$)

$$E(\hat{y}_1) \sim 1, E(\hat{y}_5) \sim 0, E(\hat{y}_7) \sim 0$$
$$\text{Var}(\hat{y}_1) \sim 0, \text{Var}(\hat{y}_5) \sim 0, \text{Var}(\hat{y}_7) \sim 0$$

Q. Is the model better calibrated in this way?



$$E_{q_M(y^*/x^*)}(y^*) \simeq \frac{1}{T} \sum_{t=1}^T \hat{y}_t^*(x^*, z_1^t, z_2^t, \dots) \quad T = 100$$
$$\hat{y}_t^* = (y_1, y_2, y_3, \dots, y_{10})_t$$

For every t plot largest three probs. of $(y_1, y_2, y_3, \dots, y_{10})_t$

THANKS !!!