

# Dropout as a Bayesian Approximation: Insights and Applications

Yarin Gal and Zoubin Ghahramani

Discussion by: Chunyuan Li

Jan. 15, 2016

## Main idea

- ▶ In the framework of variational inference, the authors show that the **standard algorithm of SGD training with Dropout** is essentially **optimizing the stochastic lower bound of Gaussian Processes whose kernel takes the form of neural networks**.

# Outline

## Preliminaries

- Dropout

- Gaussian Processes

- Variational Inference

## A Gaussian Process Approximation

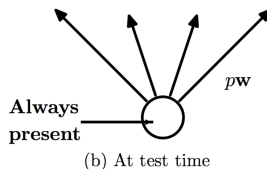
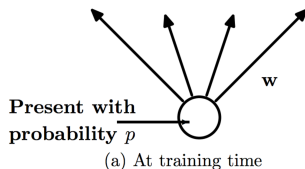
- A Gaussian Process Approximation

- Evaluating the Log Evidence Lower Bound

# Dropout

## Procedure

- ▶ Training stage: A unit is present with probability  $p$
- ▶ Testing stage: The unit is always present and the weights are multiplied by  $p$



## Intuition

- ▶ Training a neural network with dropout can be seen as training a collection of  $2^K$  thinned networks with extensive weight sharing
- ▶ A single neural network to approximate averaging output at test time

# Dropout for one-hidden-layer Neural Networks

- ▶ Dropout local units\*

$$\hat{\mathbf{y}} = (g((\mathbf{x} \odot \mathbf{b}_1)\mathbf{W}_1) \odot \mathbf{b}_2)\mathbf{W}_2 \quad (1)$$

- ▶ Input  $\mathbf{x}$  and Output  $\mathbf{y}$ ;
  - ▶  $g$ : activation function; Weights:  $\mathbf{W} \in \mathbb{R}^{K_{\ell-1} \times K_{\ell}}$
  - ▶  $\mathbf{b}_{\ell}$  is binary dropout variables
- ▶ Equivalent to multiplying the global weight matrices by the binary vectors to dropout entire rows:

$$\hat{\mathbf{y}} = g(\mathbf{x}(\text{diag}(\mathbf{b}_1)\mathbf{W}_1))(\text{diag}(\mathbf{b}_2)\mathbf{W}_2) \quad (2)$$

- ▶ Application to regression

$$\mathcal{L} = \frac{1}{2N} \sum_{n=1}^N \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|_2^2 + \lambda(\|\mathbf{W}_1\|^2 + \|\mathbf{W}_2\|^2) \quad (3)$$

---

\*bias term is ignored for simplicity

# Gaussian Processes

- ▶  $f$  is the GP function

$$\underbrace{p(f|\mathbf{X}, \mathbf{Y})}_{\text{Posterior}} \propto \underbrace{p(f)}_{\text{Prior}} \underbrace{p(\mathbf{Y}|\mathbf{X}, f)}_{\text{Likelihood}}$$

- ▶ Applications

- ▶ Regression

$$\mathbf{F}|\mathbf{X} \sim \mathcal{N}(0, \mathbf{K}(\mathbf{X}, \mathbf{X})), \quad \mathbf{Y}|\mathbf{F} \sim \mathcal{N}(\mathbf{F}, \tau^{-1}\mathbf{I}_N)$$

# Variational Inference

- ▶ The predictive distribution

$$\mathbf{K}(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{w}) \underbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{Y})}_{\approx q(\mathbf{w})} d\mathbf{w} \quad (4)$$

- ▶ Objective:  $\operatorname{argmin}_q \operatorname{KL}(q(\mathbf{w})|p(\mathbf{w}|\mathbf{X}, \mathbf{Y}))$
- ▶ Variational Prediction:  $q(\mathbf{y}^*|\mathbf{x}^*) = \int p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{w})q(\mathbf{w})d\mathbf{w}$
- ▶ Log evidence lower bound

$$\mathcal{L}_{\text{VI}} = \int q(\mathbf{w}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w})d\mathbf{w} - \operatorname{KL}(q(\mathbf{w})||p(\mathbf{w})) \quad (5)$$

- ▶ Objective:  $\operatorname{argmax}_q \mathcal{L}_{\text{VI}}$
- ▶ A variational distribution  $q(\mathbf{w})$  that explains the data well while still being close to prior

# A single-layer neural network example

- ▶ Setup
  - ▶  $Q$ : input dimension,  
 $K$ : number of hidden units,  
 $D$ : output dimension
  - ▶ Goal: Learn  $\mathbf{W}_1 \in \mathbb{R}^{Q \times K}$  and  $\mathbf{W}_2 \in \mathbb{R}^{K \times D}$  to map  $\mathbf{X} \in \mathbb{R}^{N \times Q}$  to  $\mathbf{Y} \in \mathbb{R}^{N \times D}$
- ▶ Idea: Introduce  $\mathbf{W}_1$  and  $\mathbf{W}_2$  to GP approximation



# Introduce $\mathbf{W}_1$

- ▶ Define the kernel of GP

$$\mathbf{K}(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{w}) g(\mathbf{w}^\top \mathbf{x}) g(\mathbf{w}^\top \mathbf{x}') d\mathbf{w} \quad (6)$$

- ▶ Resort to Monte Carlo integration, with generative process:

$$\mathbf{w}_k \sim p(\mathbf{w}), \mathbf{W}_1 = [\mathbf{w}_k]_{k=1}^K, \hat{\mathbf{K}}(\mathbf{x}, \mathbf{x}') = \frac{1}{K} \sum_{k=1}^K g(\mathbf{w}_k^\top \mathbf{x}) g(\mathbf{w}_k^\top \mathbf{x}')$$

$$\mathbf{F} | \mathbf{X}, \mathbf{W}_1 \sim \mathcal{N}(\mathbf{0}, \hat{\mathbf{K}}), \mathbf{Y} | \mathbf{F} \sim \mathcal{N}(\mathbf{F}, \tau^{-1} \mathbf{I}_N) \quad (7)$$

where  $K$  is the number of hidden units

## Introduce $\mathbf{W}_2$

- Analytically integrating wrt  $\mathbf{F}$   
 The predictive distribution

$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{X}, \mathbf{W}_1)p(\mathbf{W}_1)d\mathbf{w}_1d\mathbf{f} \quad (8)$$

can be rewritten as

$$p(\mathbf{Y}|\mathbf{X}) = \int \mathcal{N}(\mathbf{Y}; \mathbf{0}, \mathbf{\Phi}\mathbf{\Phi}^\top + \tau^{-1}\mathbf{I}_N)p(\mathbf{W}_1)d\mathbf{w}_1 \quad (9)$$

where  $\mathbf{\Phi} = [\phi]_{n=1}^N$ ,  $\phi = \sqrt{\frac{1}{K}g(\mathbf{W}_1^\top \mathbf{x})}$

- For  $\mathbf{W}_2 = [\mathbf{w}_2]_{d=1}^D$ ,  $\mathbf{w}_d \sim \mathcal{N}(0, \mathbf{I}_K)$

$$\mathcal{N}(\mathbf{y}_d; \mathbf{0}, \mathbf{\Phi}\mathbf{\Phi}^\top + \tau^{-1}\mathbf{I}_N) = \int \mathcal{N}(\mathbf{y}_d; \mathbf{\Phi}\mathbf{w}_d, \tau^{-1}\mathbf{I}_N)\mathcal{N}(\mathbf{w}_d; 0, \mathbf{I}_K)d\mathbf{w}_1$$

$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}_1, \mathbf{W}_2)p(\mathbf{W}_1)p(\mathbf{W}_2)d\mathbf{w}_1d\mathbf{w}_2 \quad (10)$$

# Variational Inference in the Approximate Model

$$q(\mathbf{W}_1, \mathbf{W}_2) = q(\mathbf{W}_1)q(\mathbf{W}_2) \quad (11)$$

- To mimic Dropout,  $q(\mathbf{W}_1)$  is factorised over input dimension, each of them is a Gaussian mixture distribution with two components,

$$q(\mathbf{W}_1) = \prod_{q=1}^Q q(\mathbf{w}_q), \quad q(\mathbf{w}_q) = p_1 \mathcal{N}(\mathbf{m}_q, \sigma^2 \mathbf{I}_K) + (1 - p_1) \mathcal{N}(0, \sigma^2 \mathbf{I}_K) \quad (12)$$

with  $p_1 \in [0, 1]$ ,  $\mathbf{m}_q \in \mathbb{R}^K$

- The same for  $q(\mathbf{W}_2)$

$$q(\mathbf{W}_2) = \prod_{k=1}^K q(\mathbf{w}_k), \quad q(\mathbf{w}_k) = p_2 \mathcal{N}(\mathbf{m}_k, \sigma^2 \mathbf{I}_D) + (1 - p_2) \mathcal{N}(0, \sigma^2 \mathbf{I}_D) \quad (13)$$

with  $p_2 \in [0, 1]$ ,  $\mathbf{m}_k \in \mathbb{R}^D$

- Optimise over parameters, especially  $\mathbf{M}_1 = [\mathbf{m}_q]_{q=1}^Q$ ,  $\mathbf{M}_2 = [\mathbf{m}_k]_{k=1}^K$

# Evaluating the Log Evidence Lower Bound for Regression

- ▶ Log evidence lower bound

$$\mathcal{L}_{\text{GP-VI}} = \underbrace{\int q(\mathbf{W}_1, \mathbf{W}_2) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{W}_1, \mathbf{W}_2) d\mathbf{W}_1 d\mathbf{W}_2}_{\mathcal{L}_1} - \underbrace{\text{KL}(q(\mathbf{W}_1, \mathbf{W}_2) || p(\mathbf{W}_1, \mathbf{W}_2))}_{\mathcal{L}_2} \quad (14)$$

- ▶ Approximation of  $\mathcal{L}_2^\dagger$

- ▶ For large enough  $K$  we can approximate the KL divergence term as

$$\text{KL}(q(\mathbf{W}_1) || p(\mathbf{W}_1)) \approx QK(\sigma^2 - \log(\sigma^2) - 1) + \frac{p_1}{2} \sum_{q=1}^Q \mathbf{m}_q^\top \mathbf{m}_q \quad (15)$$

- ▶ Similarly for  $\text{KL}(q(\mathbf{W}_2) || p(\mathbf{W}_2))$

---

<sup>†</sup>Following Proposition 1 on KL of a Mixture of Gaussians in Appendix

# $\mathcal{L}_1$ : Monte Carlo integration

► Parameterization

$$\mathcal{L}_1 = \int q(\mathbf{b}_1, \epsilon_1, \mathbf{b}_2, \epsilon_2) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{W}_1(\mathbf{b}_1, \epsilon_1), \mathbf{W}_2(\mathbf{b}_2, \epsilon_2)) d\mathbf{b}_1 d\mathbf{b}_2 d\epsilon_1 d\epsilon_2 \quad (16)$$

$$\begin{aligned} \mathbf{W}_1 &= \text{diag}(\mathbf{b}_1)(\mathbf{M}_1 + \sigma\epsilon_1) + (1 - \text{diag}(\mathbf{b}_1))\sigma\epsilon_1 \\ \mathbf{W}_2 &= \text{diag}(\mathbf{b}_2)(\mathbf{M}_2 + \sigma\epsilon_2) + (1 - \text{diag}(\mathbf{b}_2))\sigma\epsilon_2 \end{aligned} \quad (17)$$

where

$$\begin{aligned} \epsilon_1 &\sim \mathcal{N}(\mathbf{0}, I_{Q \times K}), b_{1q} \sim \text{Bernoulli}(p_1), \\ \epsilon_2 &\sim \mathcal{N}(\mathbf{0}, I_{K \times D}), b_{2k} \sim \text{Bernoulli}(p_2), \end{aligned} \quad (18)$$

► Take a single sample

$$\mathcal{L}_{\text{GP-MC}} = \underbrace{\log p(\mathbf{Y}|\mathbf{X}, \widehat{\mathbf{W}}_1, \widehat{\mathbf{W}}_2)}_{\mathcal{L}_{1\text{-MC}}} - \text{KL}(q(\mathbf{W}_1, \mathbf{W}_2) || p(\mathbf{W}_1, \mathbf{W}_2)) \quad (19)$$

Optimising  $\mathcal{L}_{\text{GP-MC}}$  converges to the same limit as  $\mathcal{L}_{\text{GP-VI}}$ .

$\mathcal{L}_1$  : Monte Carlo integration

## ► Regression

$$\begin{aligned}\mathcal{L}_{1\text{-MC}} &= \log p(\mathbf{Y}|\mathbf{X}, \widehat{\mathbf{W}}_1, \widehat{\mathbf{W}}_2) \\ &= \sum_{d=1}^D \log \mathcal{N}(\mathbf{y}_d; \phi \hat{\mathbf{w}}_d, \tau^{-1} \mathbf{I}_N) \\ &= -\frac{ND}{2} \log(2\pi) + \frac{ND}{2} \log(\tau) - \sum_{d=1}^D \frac{\tau}{2} \|\mathbf{y}_d - \phi \hat{\mathbf{w}}_d\|_2^2\end{aligned}\quad (20)$$

► Sum over the rows instead of the columns of  $\widehat{\mathbf{Y}}$ 

$$\sum_{d=1}^D \frac{\tau}{2} \|\mathbf{y}_d - \phi \hat{\mathbf{w}}_d\|_2^2 = \sum_{n=1}^N \frac{\tau}{2} \|\mathbf{y}_n - \phi \hat{\mathbf{w}}_n\|_2^2 \quad (21)$$

where  $\hat{\mathbf{y}}_n = \phi \widehat{\mathbf{W}}_2 = \sqrt{\frac{1}{K}} g(\mathbf{x}_n \widehat{\mathbf{W}}_1) \widehat{\mathbf{W}}_2$

# Recover SGD training with Dropout

- ▶ Take the approximation for  $\mathcal{L}_1$  and  $\mathcal{L}_2$  for the bound, and ignoring constant terms  $\tau$  and  $\sigma$

$$\mathcal{L}_{\text{GP-MC}} = -\frac{\tau}{2} \sum_{n=1}^N \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|_2^2 - \frac{p_1}{2} \|\mathbf{M}_1\|_2^2 - \frac{p_2}{2} \|\mathbf{M}_2\|_2^2 \quad (22)$$

- ▶ Setting  $\sigma$  tend to 0

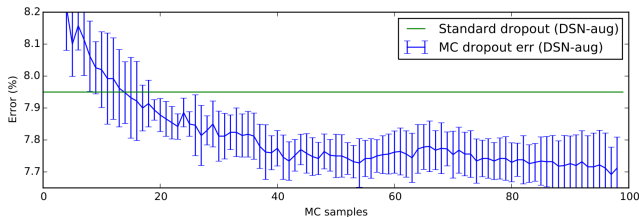
$$\begin{aligned} \mathbf{W}_1 &= \text{diag}(\mathbf{b}_1)(\mathbf{M}_1 + \sigma \mathbf{e}_1) + (1 - \text{diag}(\mathbf{b}_1))\sigma \mathbf{e}_1 \Rightarrow \\ \widehat{\mathbf{W}}_1 &\approx \text{diag}(\hat{\mathbf{b}}_1)\mathbf{M}_1, \quad \widehat{\mathbf{W}}_2 \approx \text{diag}(\hat{\mathbf{b}}_2)\mathbf{M}_2 \end{aligned} \quad (23)$$

$$\hat{\mathbf{y}}_n = \sqrt{\frac{1}{K}} g(\mathbf{x}_n \widehat{\mathbf{W}}_1) \widehat{\mathbf{W}}_2 = \sqrt{\frac{1}{K}} g(\mathbf{x}_n (\text{diag}(\hat{\mathbf{b}}_1)\mathbf{M}_1)) (\text{diag}(\hat{\mathbf{b}}_2)\mathbf{M}_2) \quad (24)$$

## More Applications

### ► Convolutional Neural Networks

*Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference*, arXiv:1506.02158, 2015



### ► Recurrent Neural Networks

*A Theoretically Grounded Application of Dropout in Recurrent Neural Networks*, arXiv:1512.05287, 2015

### ► Reinforcement Learning

*Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning*, arXiv:1506.02142, 2015