

# Dropout as a Bayesian Approximation

Jaewoo Park

Department of Applied Statistics, Yonsei University

# Variational Inference (VI)

- Approximates the posterior through  $q_{\xi}(\mathbf{W}, \mathbf{b})$  (variational distribution)
- Find  $q_{\xi}(\mathbf{W}, \mathbf{b})$  that minimizes the Kullback–Leibler (KL) divergence between  $q_{\xi}(\mathbf{W}, \mathbf{b})$  and  $p(\mathbf{W}, \mathbf{b} | \mathbf{X}, \mathbf{Y})$ 
  - We need to set a class of distributions (e.g., Gaussian)
  - A practical option for BNN

# Variational Inference (contd.)

- KL divergence is defined as

$$\text{KL}(q_{\xi}(\mathbf{W}, \mathbf{b}) || p(\mathbf{W}, \mathbf{b} | \mathbf{X}, \mathbf{Y})) = \int \int q_{\xi}(\mathbf{W}, \mathbf{b}) \frac{q_{\xi}(\mathbf{W}, \mathbf{b})}{p(\mathbf{W}, \mathbf{b} | \mathbf{X}, \mathbf{Y})} d\mathbf{W} d\mathbf{b}$$

which is intractable due to  $p(\mathbf{W}, \mathbf{b} | \mathbf{X}, \mathbf{Y})$  terms

- Minimizing the KL divergence is equivalent to maximizing evidence lower-bound (ELBO)

$$\begin{aligned} \text{ELBO} = \int \int q_{\xi}(\mathbf{W}, \mathbf{b}) \log p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \mathbf{b}) d\mathbf{W} d\mathbf{b} - \\ \text{KL}(q_{\xi}(\mathbf{W}, \mathbf{b}) || p(\mathbf{W})p(\mathbf{b})) \end{aligned}$$

# Variational Inference (contd.)

- We don't need the posterior  $p(\mathbf{W}, \mathbf{b} | \mathbf{X}, \mathbf{Y})$  when we compute ELBO

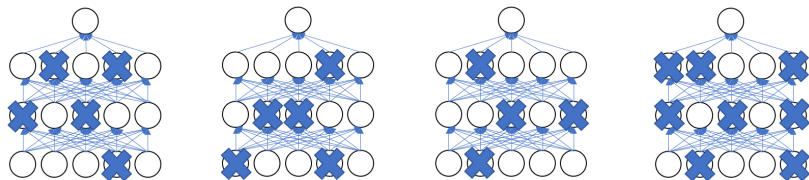
$$\text{ELBO} = \int \int q_{\xi}(\mathbf{W}, \mathbf{b}) \log p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \mathbf{b}) d\mathbf{W} d\mathbf{b} - \text{KL}(q_{\xi}(\mathbf{W}, \mathbf{b}) || p(\mathbf{W})p(\mathbf{b}))$$

- All we need is
  - Prior:  $p(\mathbf{W})p(\mathbf{b})$
  - Variational distribution:  $q_{\xi}(\mathbf{W}, \mathbf{b})$
  - Likelihood:  $p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \mathbf{b})$
- Still challenging for obtaining analytical solutions

# Monte Carlo (MC) Dropout

- MC dropout (Gal and Ghahramani, 2016) is an practical option
  - ① Fit a neural network in a usual way
  - ② Based on the fitted model, apply different dropouts to make prediction
- We are fitting the model just a single time (not multiple times)
- MC dropout can approximate ELBO (see later)

# Generating Monte Carlo Samples through Dropout



- Consider we have  $\widehat{\mathbf{W}}, \widehat{\mathbf{b}}$  from SGD
- For given  $\widehat{\mathbf{W}}, \widehat{\mathbf{b}}$ , dropout randomly removes the connections from nodes
- For each dropout, we will have different output realizations (predictive distribution)

# Dropout as a Bayesian Approximation

- Go back to VI, we have

$$\text{ELBO} = \int \int q_{\xi}(\mathbf{W}, \mathbf{b}) \log p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \mathbf{b}) d\mathbf{W} d\mathbf{b} - \text{KL}(q_{\xi}(\mathbf{W}, \mathbf{b}) || p(\mathbf{W})p(\mathbf{b}))$$

- All we need is
  - Prior:  $p(\mathbf{W})p(\mathbf{b}) \rightarrow$  standard normal prior
  - Variational distribution:  $q_{\xi}(\mathbf{W}, \mathbf{b}) = q_{\xi}(\mathbf{W})q_{\xi}(\mathbf{b}) \rightarrow$  mixture normal
  - Likelihood:  $p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \mathbf{b}) \rightarrow$  (nested) Gaussian process

# Dropout as a Bayesian Approximation: $q(\mathbf{W}, \mathbf{b})$

- We use a mixture normal distribution as

$$q_{\xi}(\mathbf{W}) = \prod_{\forall i} q_{\xi}(w_i), \quad q_{\xi}(\mathbf{b}) = \prod_{\forall i} q_{\xi}(b_i)$$

$$q_{\xi}(w_i) = pN(\mu_i^w, \sigma^2) + (1 - p)N(0, \sigma^2)$$

$$q_{\xi}(b_i) = pN(\mu_i^b, \sigma^2) + (1 - p)N(0, \sigma^2)$$

- Bayesian variable selection (dropout)
  - $p \rightarrow 0$ : likely to drop ( $w_i, b_i = 0$ ) parameters
  - $p \rightarrow 1$ : likely to include ( $w_i, b_i \neq 0$ ) parameters



# Dropout as a Bayesian Approximation: $p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \mathbf{b})$

- How can we define a likelihood function?
- Let's use a nested Gaussian process

$$\mathbf{F}_l | \mathbf{F}_{l-1} \sim N(0, \Sigma_l), \quad l = 2, \dots, L$$

$$\mathbf{Y} | \mathbf{F}_L \sim p(\mathbf{Y} | \mathbf{F}_L)$$

Here  $p(\cdot | \mathbf{F}_L)$  is a exponential family distribution

# Dropout as a Bayesian Approximation: ELBO

- Remind:

$$\text{ELBO} = \int \int q_{\xi}(\mathbf{W}) q_{\xi}(\mathbf{b}) \log p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \mathbf{b}) d\mathbf{W} d\mathbf{b} - \text{KL}(q_{\xi}(\mathbf{W}) q_{\xi}(\mathbf{b}) || p(\mathbf{W}) p(\mathbf{b}))$$

- We can approximate this as

$$\mathcal{L}_{\text{GP-MC}} := \frac{1}{M} \sum_{m=1}^M \log p(\mathbf{Y} | \mathbf{X}, \widehat{\mathbf{W}}_m, \widehat{\mathbf{b}}_m) - \text{KL}(q_{\xi}(\mathbf{W}) q_{\xi}(\mathbf{b}) || p(\mathbf{W}) p(\mathbf{b}))$$

where  $\{\widehat{\mathbf{W}}_m, \widehat{\mathbf{b}}_m\}_{m=1}^M$  are MC samples from  $q_{\xi}(\mathbf{W}) q_{\xi}(\mathbf{b})$

# Dropout as a Bayesian Approximation: Prediction

- The predictive distribution is

$$q(\mathbf{Y}^*|\mathbf{X}^*) = \int \int p(\mathbf{Y}^*|\mathbf{X}^*, \mathbf{W}, \mathbf{b}) q_{\xi}(\mathbf{W}) q_{\xi}(\mathbf{b}) d\mathbf{W} d\mathbf{b}$$

- For given unobserved  $\mathbf{X}^*, \mathbf{Y}^*$

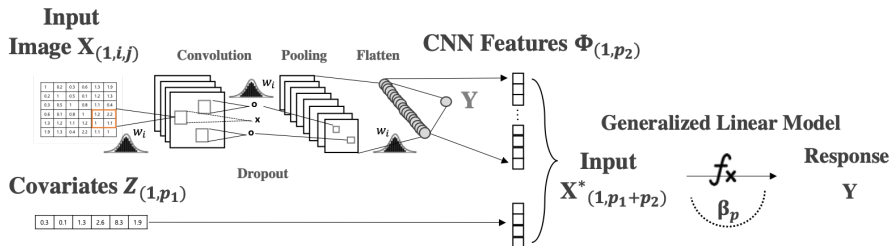
$$\mathbb{E}_q(\mathbf{Y}^*) = \int \mathbf{Y}^* q(\mathbf{Y}^*|\mathbf{X}^*) d\mathbf{Y}^* \approx \frac{1}{M} \sum_{m=1}^M \hat{\mathbf{Y}}_m^*$$

where  $\hat{\mathbf{Y}}_m^*$  is sampled from  $q(\mathbf{Y}^*|\mathbf{X}^*)$

# A Bayesian Convolutional Neural Network-based Generalized Linear Model (Bayes CGLM)

- Goal: Analyze relationships between a response  $\mathbf{Y}$  and inputs  $\mathbf{X}, \mathbf{Z}$
- Statistical Methods:
  - Provide interpretation through regression coefficients
  - Quantify uncertainties in estimates/predictions
  - Hard to analyze high-dimensional variables (e.g., image)
- Deep Learning Methods:
  - Popular for high-dimensional variables with high accuracy
  - Hard to interpret the impact of covariates
  - Uncertainty quantification is not trivial
- Jeon et al., (2023) combine statistical and deep learning models

# Bayes CGLM



- Step 1: From  $X$ , Extract  $\{\Phi^{(m)}\}_{m=1}^M$  via Bayes CNN
- Step 2: Fit Bayes GLMs by regressing  $Y$  on  $[Z, \Phi^{(m)}]$
- Step 3: Construct an ensemble-posterior distribution

# Part 1: Extract Feature Information

- Idea: Monte Carlo dropout (Gal and Ghahramani 2015, 2016)
  - A deep Gaussian process (Damianou and Lawrence, 2013) can represent complex DNN
  - Dropout can approximate the deep Gaussian process
  - We extract MC samples  $\{\Phi^{(m)}\}_{m=1}^M$  (features) from the last layer of Bayes CNN
- Benefits:
  - Reduce the dimension of high-dimensional  $\mathbf{X}$
  - Simultaneously analyze both  $\mathbf{X}$  and  $\mathbf{Z}$  with appropriate uncertainty quantification

## Part 2: Fit Bayes GLMs

- Obtain  $m$ th feature-posterior:

- Model:

$$E[\mathbf{Y}|\mathbf{Z}, \Phi^{(m)}] = \mathbf{Z}\gamma_m + \Phi^{(m)}\delta_m = \mathbf{A}^{(m)}\beta_m$$

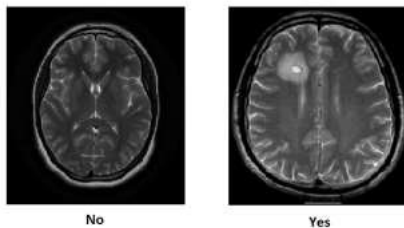
- $\mathbf{Z}$ : covariate matrix
- $\Phi^{(m)}$ : extracted feature from  $m$ th dropout sample
- Fit Bayes LM by regressing  $\mathbf{Y}$  on  $\mathbf{A}^{(m)} = [\mathbf{Z}, \Phi^{(m)}]$

## Part 3: An Ensemble Posterior

- We obtain  $M$  number of feature-posterior
- We combine them to construct the aggregated posterior
  - You can simply combine posterior samples from each  $m$
- Such aggregation can fully account for uncertainties in the feature extraction (Step 1)



# Brain Tumor Image Data



The dataset is collected from Brain Tumor Image Segmentation Challenge (BRATS)

- **Y**: Binary response (brain tumor or not)
- **X**:  $240 \times 240$  pixel gray images of brains
- **Z**: Vector covariates (first and second order features)
- $N = 2,508$  for training and  $N_{cv} = 2,007$  for validation

# Brain Tumor Image Data (contd.)

	BayesCGLM(M=500)	Bayes CNN	GLM
$\gamma_1$ (first order feature)	-5.332 (-7.049,-3.704)	0.248 -	-2.591 (-2.769,-2.412)
$\gamma_2$ (second order feature)	4.894 (3.303, 6.564)	0.160 -	2.950 (2.755, 3.144)
Accuracy	0.924	0.867	0.784
Recall	0.929	0.787	0.783
Precision	0.901	0.907	0.715
Time	293.533	103.924	0.004

**Table:** For all methods, the posterior mean of  $\gamma$ , 95% HPD interval, accuracy, recall, precision, and computing time (min) are reported in the table.

## Bayes CNN-based GLM (BayesCGLM)

- Extract the important feature of the high-dimensional variables
- Simultaneously utilize covariates with different data structures (e.g., vector and image)
- Provide accurate inference with uncertainty quantification in both estimation/prediction
- Provide interpretation of the impact of covariates
- Computationally efficient due to parallelization