

# Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning

Yeseul Jeon

January 10, 2021

# Contents

- 1 Abstract
- 2 Introduction
- 3 Dropout as a Bayesian Approximation
- 4 Obtaining Model Uncertainty
- 5 Conclusion

---

## **Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning**

---

**Yarin Gal**  
**Zoubin Ghahramani**  
University of Cambridge

YG279@CAM.AC.UK  
ZG201@CAM.AC.UK

- **Limitation:**

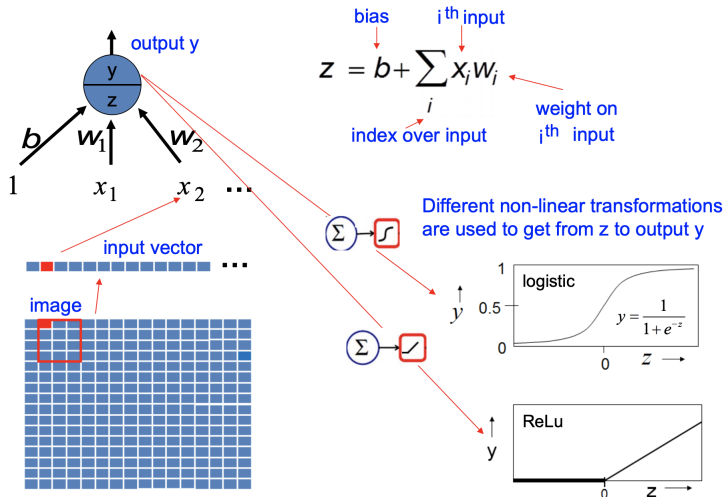
- Deep NN gives us only point estimates with no uncertainty information.
- Bayesian modeling we can get a measure of uncertainty by evaluating the posterior distribution of the NN weights.
- Bayesian model usually come with a prohibitive computational cost.

- **Idea:**

- Develop a new theoretical framework casting dropout training in deep neural networks as approximate Bayesian inference in deep Gaussian process
- This theory gives us tools to model uncertainty with dropout NNs - extracting information from existing models that has been thrown away so far.

# Introduction

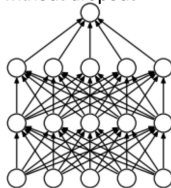
- Logistic Regression or NN with 1 neuron:



## ● Dropout:

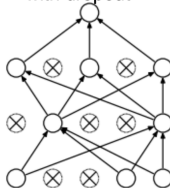
- At each training step we remove random nodes with a probability of  $p$  resulting in a sparse version of the full net and we use back-propagation to update the weights
- By averaging over these models we should be able to "reduce noise", "over-fitting".

without dropout



(a) Standard Neural Net

with dropout



(b) After applying dropout.

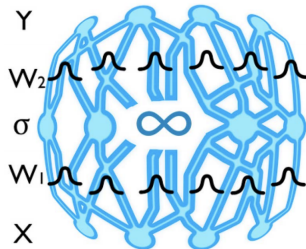
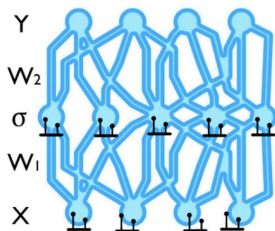
$$\mathcal{L}_{\text{dropout}} := \frac{1}{N} \sum_{i=1}^N E(\mathbf{y}_i, \hat{\mathbf{y}}_i) + \lambda \sum_{i=1}^L (\|\mathbf{W}_i\|_2^2 + \|\mathbf{b}_i\|_2^2).$$

- NN Optimisation

- A regularisation term is added.
- $E(y_i, \hat{y}_i)$ : Error function
- Sample binary variables for every input point and for every network unit in each layer.

# Dropout as a Bayesian Approximation

- Dropout vs Bayesian NN



- Dropout: Remove random nodes with a probability  $p$
- Bayesian NN: Update the posterior distribution of the weights



## • Posterior Distribution

$$p(\omega | \mathbf{X}, \mathbf{Y}) = \frac{\overset{\text{likelihood}}{p(\mathbf{Y}|\omega, \mathbf{X})} \cdot \overset{\text{prior}}{p(\omega)}}{\underset{\text{normalizer=marginal likelihood}}{p(\mathbf{Y}|\mathbf{X})}}$$

- We can approximate the posterior distribution for the model parameters via Variation Inference
- replacing the posterior distribution at the observed data  $p(w|\mathbf{X}, \mathbf{Y})$  with a member  $q(w)$  of a simpler distribution family  $Q$  that minimizes the Kullback–Leibler divergence to the posterior

## • Variational Inference

$$\text{KL}(q_{\theta}(\omega) \parallel p(\omega \mid \mathbf{X}, \mathbf{Y})) = \int q_{\theta}(\omega) \log \frac{q_{\theta}(\omega)}{p(\omega \mid \mathbf{X}, \mathbf{Y})} d\omega = E_q [\log(q_{\theta}(\omega)) - \log(p(\omega \mid \mathbf{X}, \mathbf{Y}))]$$

- Approximated  $p(w|X, Y)$  with simple distribution  $q_{\theta}(w)$
- Minimize Kullback Leibler divergence of  $q$  from the posterior w.r.t to the variational parameters  $\theta$ :

## • Variational Inference

$$\begin{aligned} L &= \log(p(Y|X)) = \log \int p(Y|X, \omega) \cdot p(\omega) d\omega = \log \int p(Y|X, \omega) \cdot p(\omega) \frac{q_\theta(\omega)}{q_\theta(\omega)} d\omega = \log \left( E_{q_\theta} \left[ \frac{p(Y|X, \omega) \cdot p(\omega)}{q_\theta(\omega)} \right] \right) \\ &\geq E_{q_\theta} \left[ \log \left( \frac{p(Y|X, \omega) \cdot p(\omega)}{q_\theta(\omega)} \right) \right] = E_{q_\theta} \left[ \log(p(Y|X, \omega)) + \log \left( \frac{p(\omega)}{q_\theta(\omega)} \right) \right] = E_{q_\theta} [\log(p(Y|X, \omega))] - E_{q_\theta} \left[ \log \left( \frac{q_\theta(\omega)}{p(\omega)} \right) \right] \\ &= \int q_\theta(\omega) \cdot \log(p(Y|X, \omega)) d\omega - KL(q_\theta(\omega) \| p(\omega)) \end{aligned}$$

- Minimizing the KL divergence of  $q$  from the posterior distribution w.r.t  $\theta$  is
- equivalent to maximizing a lower bound of the log marginal likelihood w.r.t  $\theta$

- **MC integration to approximate  $L$**

$$L_{VI}(\theta) := \int q_{\theta}(\omega) \cdot \log(p(Y|X, \omega)) d\omega - KL(q_{\theta}(\omega) \| p(\omega))$$

- Since this integral is not tractable for almost all  $q$  therefore we will MC integration to approximate this quantity.
- sample  $\hat{\omega}$  from  $q$  and each sampling step the integral is replaced by  $\log(p(Y|X, \hat{\omega}))$ .

- **Stochastic Inference**

$$L_{VI}(\theta) := \int q_{\theta}(\omega) \cdot \log(p(Y|X, \omega)) d\omega - KL(q_{\theta}(\omega) \| p(\omega))$$

$$\hat{L}(\theta) := \log(p(Y|X, \hat{\omega})) - KL(q_{\theta}(\omega) \| p(\omega))$$

- For inference repeatedly do:

- Sample  $\hat{w} \sim q_{\theta}(w)$
- Do one step of minimization w.r.t  $\theta$ :  $\hat{L}(\theta)$

# What kind of q-distribution should we use?

- The deep Gaussian process

$$\mathbf{K}(\mathbf{x}, \mathbf{y}) = \int p(\mathbf{w})p(b)\sigma(\mathbf{w}^T \mathbf{x} + b)\sigma(\mathbf{w}^T \mathbf{y} + b)d\mathbf{w}db$$

$$\mathbf{w}_k \sim p(\mathbf{w}), b_k \sim p(b),$$

$$\mathbf{W}_1 = [\mathbf{w}_k]_{k=1}^K, \mathbf{b} = [b_k]_{k=1}^K$$

$$\hat{\mathbf{K}}(\mathbf{x}, \mathbf{y}) = \frac{1}{K} \sum_{k=1}^K \sigma(\mathbf{w}_k^T \mathbf{x} + b_k)\sigma(\mathbf{w}_k^T \mathbf{y} + b_k)$$

$$\mathbf{F} \mid \mathbf{X}, \mathbf{W}_1, \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \hat{\mathbf{K}}(\mathbf{X}, \mathbf{X}))$$

$$\mathbf{Y} \mid \mathbf{F} \sim \mathcal{N}(\mathbf{F}, \tau^{-1} \mathbf{I}_N),$$

- $W_i$  be a random matrix of dimensions  $K_i \times K_{i-1}$  for each layer  $i$ .
- A prior let each row of  $W_i$  distribute according to the  $p(\mathbf{w})$  above.
- Assume vectors  $m_i$  of dimensions  $K_i$  for each GP layer.

# What kind of q-distribution should we use?

- **Prediction Probability of deep GP model**

$$p(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega})p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})d\boldsymbol{\omega}$$

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}(\mathbf{x}, \boldsymbol{\omega}), \tau^{-1}\mathbf{I}_D)$$

$$\hat{\mathbf{y}}(\mathbf{x}, \boldsymbol{\omega} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\})$$

- The posterior distribution  $p(w|X, Y)$  is intractable.
- Use  $q(w)$ , a distribution over matrices whose columns are randomly set to zero
- To approximate the intractable posterior

# Define the structure of the approximate distribution $q$

- $q(\mathbf{w})$  distribution

$$\mathbf{W}_i = \mathbf{M}_i \cdot \text{diag}([\mathbf{z}_{i,j}]_{j=1}^{K_i})$$

$$\mathbf{z}_{i,j} \sim \text{Bernoulli}(p_i) \text{ for } i = 1, \dots, L, j = 1, \dots, K_{i-1}$$

$$- \int q(\boldsymbol{\omega}) \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}) d\boldsymbol{\omega} + \text{KL}(q(\boldsymbol{\omega})||p(\boldsymbol{\omega})).$$

We rewrite the first term as a sum

$$- \sum_{n=1}^N \int q(\boldsymbol{\omega}) \log p(\mathbf{y}_n|\mathbf{x}_n, \boldsymbol{\omega}) d\boldsymbol{\omega}$$

- Some probabilities  $p$
- Matrices  $M_i$  as variational parameters.
- The variational distribution  $q(\mathbf{w})$  is highly multimodal
- This corresponds to the frequencies in the sparse spectrum GP approximation



# Define the structure of the approximate distribution $q$

- KL with GP and MC

$$\begin{aligned}\mathcal{L}_{\text{GP-MC}} &\propto \frac{1}{N} \sum_{n=1}^N \frac{-\log p(\mathbf{y}_n | \mathbf{x}_n, \hat{\mathbf{w}}_n)}{\tau} \\ &\quad + \sum_{i=1}^L \left( \frac{p_i l^2}{2\tau N} \|\mathbf{M}_i\|_2^2 + \frac{l^2}{2\tau N} \|\mathbf{m}_i\|_2^2 \right). \\ \mathcal{L}_{\text{GP-MC}} &\propto -\frac{1}{2N} \sum_{n=1}^N \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|_2^2 - \frac{l^2 p_1}{2\tau N} \|\mathbf{M}_1\|_2^2 - \frac{K p_2}{2\tau N} \|\mathbf{M}_2\|_2^2 - \frac{l'^2}{2\tau N} \|\mathbf{m}\|_2^2.\end{aligned}$$

- The sampled  $\hat{\mathbf{w}}$  result in realization from the Bernoulli distribution  $z_i^n$  equivalent to the binary variables in the dropout case.

# Define the structure of the approximate distribution $q$

- $q_{\theta}(w)$

$$q_{\mathbf{M}_i}(\mathbf{W}_i) = \mathbf{M}_i \cdot \text{diag}(\lfloor \mathbf{z}_{i,j} \rfloor)$$

$$\mathbf{z}_{i,j} \sim \text{Bernoulli}(\mathbf{p}_i)$$

$$\mathbf{W}_i \sim q_{\mathbf{M}_i}(\mathbf{W}_i)$$

Approximate posterior  
of model parameters

$$\mathbf{M}_i = \text{mean}(\mathbf{W}_i)$$

$\mathbf{M}_i$  is as variational parameter of  $q$

Sampling the diagonal elements  $\mathbf{z}$  from a Bernoulli is identical to randomly setting columns of  $\mathbf{M}$  to zero which is identical to randomly setting units of the network to zero -> dropout!

- Bernoullis are computationally cheap to get multi-modality

- **Prediction distribution for uncertainty estimation**

- To get an approximation of the posterior via training

- 1) Randomly set columns of  $M_i$  to zero (do dropout)

- 2) Update the weights by doing one step

- To sample from the learned approximate posterior we just can do dropout during the test time when using the trained NN for prediction.

- From the received predictions we can estimate the predictive distribution and from this different uncertainty measures such as the variance.

- **Prediction distribution for uncertainty estimation**

- To sample from the learned approximate posterior do dropout during the test time when using the trained NN for prediction.
- From the received predictions we can estimate the predictive distribution and from this different uncertainty measures such as the variance.

- **Prediction distribution for uncertainty estimation**

$$\begin{aligned}\log p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) &= \log \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) p(\omega|\mathbf{X}, \mathbf{Y}) d\omega \\ &\approx \log \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) q(\omega) d\omega \\ &\approx \log \left( \frac{1}{T} \sum_{t=1}^T p(y^*|\mathbf{x}^*, \omega_t) \right)\end{aligned}$$

- Given a dataset  $X, Y$  and a new data point  $x^*$  we can calculate the probability of possible output values  $y^*$  using the predictive probability  $p(y^*|x^*, X, Y)$ .

- **Mean and Variance using MC-dropout**

$$\begin{aligned}\mathbb{E}_{q(\mathbf{y}^*|\mathbf{x}^*)}(\mathbf{y}^*) &\approx \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}^*(\mathbf{x}^*, \mathbf{W}_1^t, \dots, \mathbf{W}_L^t) \\ \text{Var}_{q(\mathbf{y}^*|\mathbf{x}^*)}(\mathbf{y}^*) &\approx \tau^{-1} \mathbf{I}_D \\ &+ \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{y}}^*(\mathbf{x}^*, \mathbf{W}_1^t, \dots, \mathbf{W}_L^t)^T \hat{\mathbf{y}}^*(\mathbf{x}^*, \mathbf{W}_1^t, \dots, \mathbf{W}_L^t) \\ &- \mathbb{E}_{q(\mathbf{y}^*|\mathbf{x}^*)}(\mathbf{y}^*)^T \mathbb{E}_{q(\mathbf{y}^*|\mathbf{x}^*)}(\mathbf{y}^*)\end{aligned}$$

- **Estimate Model uncertainty with MC-dropout**

- Can represent model uncertainty in deep learning, better model regularisation, computationally efficient Bayesian convolutional neural networks.
- A neural network with arbitrary depth and non-linearities and with dropout applied before every weight layer is mathematically equivalent to an approximation to the deep Gaussian process (marginalised over its covariance function parameters).