

# Flat Seeking Bayesian Neural Networks

Nguyen et al. 2023

Reviewer: Jihu Lee

IDEA lab  
Department of Statistics  
Seoul National University

January 8, 2024

## Objective of work

- Propose a sharpness-aware posterior
- Propose its variational approach

## Notations

- $f_\theta(x)$ : model (NN) with  $\theta \in \Theta$
- $\mathcal{S} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ : training set, where  $(x_i, y_i) \sim \mathcal{D}$
- $\mathbb{P}, p(\theta)$ : prior distribution, prior density
- $l(f_\theta(x), y)$ : loss function

## Standard Posterior

$$q(\theta|\mathcal{S}) \propto \prod_{i=1}^n p(y_i|x_i, \mathcal{S}, \theta)p(\theta) \quad (1)$$

**Likelihood** ( $\lambda \geq 0$ : regularization parameter)

$$p(y|x, \mathcal{S}, \theta) \propto \exp \left\{ -\frac{\lambda}{|\mathcal{S}|} l(f_\theta(x), y) \right\} = \exp \left\{ -\frac{\lambda}{n} l(f_\theta(x), y) \right\} \quad (2)$$

## Population & Empirical losses

$$\mathcal{L}_{\mathcal{D}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[l(f_{\theta}(x), y)] \quad (3)$$

$$\mathcal{L}_{\mathcal{S}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{S}}[l(f_{\theta}(x), y)] = \frac{1}{n} \sum_{i=1}^n l(f_{\theta}(x_i), y_i) \quad (4)$$

## Standard Posterior (rewritten)

$$q(\theta|\mathcal{S}) \propto \exp\{-\lambda \mathcal{L}_{\mathcal{S}}(\theta)\} p(\theta) \quad (5)$$

## Population & Empirical losses over $\mathbb{Q}$ (distribution over $\theta$ with density $q$ )

$$\mathcal{L}_{\mathcal{D}}(\mathbb{Q}) = \int_{\Theta} \mathcal{L}_{\mathcal{D}}(\theta) d\mathbb{Q}(\theta) = \int_{\Theta} \mathcal{L}_{\mathcal{D}}(\theta) q(\theta) d\theta \quad (6)$$

$$\mathcal{L}_{\mathcal{S}}(\mathbb{Q}) = \int_{\Theta} \mathcal{L}_{\mathcal{S}}(\theta) d\mathbb{Q}(\theta) = \int_{\Theta} \mathcal{L}_{\mathcal{S}}(\theta) q(\theta) d\theta \quad (7)$$

## Theorem (3.1)

Consider the following optimization problem:

$$\min_{\mathbb{Q} \ll \mathbb{P}} \{ \lambda \mathcal{L}_{\mathcal{S}}(\mathbb{Q}) + KL(\mathbb{Q}, \mathbb{P}) \} \quad (8)$$

This optimization has a closed-form optimal solution  $\mathbb{Q}^*$  with the density

$$q^*(\theta) \propto \exp \{ -\lambda \mathcal{L}_{\mathcal{S}}(\theta) \} p(\theta) \quad (9)$$

which is exactly the standard posterior  $\mathbb{Q}_{\mathcal{S}}$  with the density  $q(\theta|\mathcal{S})$ .

## Optimization Problem (population)

$$\min_{\mathbb{Q} \ll \mathbb{P}} \{ \lambda \mathcal{L}_{\mathcal{D}}(\mathbb{Q}) + KL(\mathbb{Q}, \mathbb{P}) \} \quad (10)$$

## Theorem (3.2)

Assume that  $\Theta$  is a compact set. Under some mild conditions given any  $\delta \in [0, 1]$ , with the probability at least  $1 - \delta$  over the choice of  $S \sim \mathcal{D}^n$ , for any distribution  $\mathbb{Q}$ , we have

$$\mathcal{L}_{\mathcal{D}}(\mathbb{Q}) \leq \mathbb{E}_{\theta \sim \mathbb{Q}} \left[ \max_{\theta': \|\theta' - \theta\| \leq \rho} \mathcal{L}_S(\theta') \right] + f \left( \max_{\theta \in \Theta} \|\theta\|^2, n \right) \quad (11)$$

where  $f$  is a non-decreasing function w.r.t. the first variable and approaches 0 when the training size  $n$  approaches  $\infty$ .

## Upper bound (rewritten)

$$\mathcal{L}_{\mathcal{D}}(\mathbb{Q}) \leq \mathcal{L}_{\mathcal{S}}(\mathbb{Q}) \tag{12}$$

$$+ \mathbb{E}_{\theta \sim \mathbb{Q}} \left[ \max_{\theta': \|\theta' - \theta\| \leq \rho} \mathcal{L}_{\mathcal{S}}(\theta') - \mathcal{L}_{\mathcal{S}}(\theta) \right] \tag{13}$$

$$+ f \left( \max_{\theta \in \Theta} \|\theta\|^2, n \right) \tag{14}$$

- First term: empirical losses over  $\mathbb{Q}$
- Second term: sharpness over  $\mathbb{Q}$
- Last term: constant



## Optimization Problem

$$\min_{\mathbb{Q} \ll \mathbb{P}} \left\{ \lambda \mathbb{E}_{\theta \sim \mathbb{Q}} \left[ \max_{\theta': \|\theta' - \theta\| \leq \rho} \mathcal{L}_{\mathcal{S}}(\theta') \right] + KL(\mathbb{Q}, \mathbb{P}) \right\} \quad (15)$$

- Considering sharpness-aware loss  $\max_{\theta': \|\theta' - \theta\| \leq \rho} \mathcal{L}_{\mathcal{S}}(\theta')$  makes generalization ability better

## Theorem (3.3)

The optimal solution of eq. (15) is the sharpness-aware posterior distribution  $\mathbb{Q}_S^{SA}$  with the density function  $q^{SA}(\theta|\mathcal{S})$ :

$$q^{SA}(\theta|\mathcal{S}) \propto \exp \left\{ -\lambda \max_{\theta': \|\theta' - \theta\| \leq \rho} \mathcal{L}_S(\theta') \right\} p(\theta) = \exp \{ -\lambda \mathcal{L}_S(s(\theta)) \} p(\theta) \quad (16)$$

where  $s(\theta) = \operatorname{argmax}_{\theta': \|\theta' - \theta\| \leq \rho} \mathcal{L}_S(\theta')$ .

- Closed form of the sharpness-aware posterior distribution  $\mathbb{Q}_S^{SA}$

## Sharpness-aware Likelihood

$$p^{SA}(y|x, \mathcal{S}, \theta) \propto \exp \left\{ -\frac{\lambda}{n} l(f_{s(\theta)}(x), y) \right\} \quad (17)$$

where  $s(\theta) = \underset{\theta': \|\theta' - \theta\| \leq \rho}{\mathcal{L}_{\mathcal{S}}(\theta')}$

## Sharpness-aware Posterior

$$q^{SA}(\theta|\mathcal{S}) \propto \prod_{i=1}^n p^{SA}(y_i|x_i, \mathcal{S}, \theta)p(\theta) \quad (18)$$

# Variational Approach

- $X = [x_1, \dots, x_n], Y = [y_1, \dots, y_n]$
- $\{q_\phi(\theta) : \phi \in \Phi\}$ : approximate posterior family

## ELBO

$$\max_{q_\phi} \left\{ \mathbb{E}_{q_\phi(\theta)} \left[ \sum_{i=1}^n \log p^{SA}(y_i | x_i, \mathcal{S}, \theta) \right] - KL(q_\phi, p) \right\} \quad (19)$$

$$\Leftrightarrow \min_{q_\phi} \left\{ \lambda \mathbb{E}_{q_\phi(\theta)} \left[ \max_{\theta' : \|\theta' - \theta\| \leq \rho} \mathcal{L}_{\mathcal{S}}(\theta') \right] + KL(q_\phi, p) \right\} \quad (20)$$

Table 1: Classification score on CIFAR-100 dataset. Each experiment is repeated three times with different random seeds and reports the mean and standard deviation.

Method	PreResNet-164			WideResNet28x10		
	ACC $\uparrow$	NLL $\downarrow$	ECE $\downarrow$	ACC $\uparrow$	NLL $\downarrow$	ECE $\downarrow$
<b>Variational inference</b>						
MC-Dropout	79.50 $\pm$ 0.37	0.9162 $\pm$ 0.0103	0.0993 $\pm$ 0.0033	82.30 $\pm$ 0.19	0.6500 $\pm$ 0.0049	0.0574 $\pm$ 0.0028
F-MC-Dropout	<b>81.06 <math>\pm</math> 0.44</b>	<b>0.7027 <math>\pm</math> 0.0049</b>	<b>0.0514 <math>\pm</math> 0.0047</b>	<b>83.24 <math>\pm</math> 0.11</b>	<b>0.6144 <math>\pm</math> 0.0068</b>	<b>0.0250 <math>\pm</math> 0.0027</b>
Deep-ens	82.08 $\pm$ 0.42	0.7189 $\pm$ 0.0108	0.0334 $\pm$ 0.0064	83.04 $\pm$ 0.15	0.6958 $\pm$ 0.0335	0.0483 $\pm$ 0.0017
F-Deep-ens	<b>82.54 <math>\pm</math> 0.10</b>	<b>0.6286 <math>\pm</math> 0.0022</b>	<b>0.0143 <math>\pm</math> 0.0041</b>	<b>84.52 <math>\pm</math> 0.03</b>	<b>0.5644 <math>\pm</math> 0.0106</b>	<b>0.0191 <math>\pm</math> 0.0039</b>
<b>Markov chain Monte Carlo</b>						
SGLD	80.13 $\pm$ 0.01	0.7604 $\pm$ 0.0010	0.1161 $\pm$ 0.0031	81.38 $\pm$ 0.10	0.7123 $\pm$ 0.0204	0.0958 $\pm$ 0.0004
F-SGLD	<b>80.82 <math>\pm</math> 0.02</b>	<b>0.7276 <math>\pm</math> 0.0012</b>	<b>0.1085 <math>\pm</math> 0.0008</b>	<b>82.12 <math>\pm</math> 0.16</b>	<b>0.6722 <math>\pm</math> 0.0112</b>	<b>0.0820 <math>\pm</math> 0.0021</b>
<b>Sample</b>						
SWAG-Diag	80.18 $\pm$ 0.50	0.6837 $\pm$ 0.0186	<b>0.0239 <math>\pm</math> 0.0047</b>	82.40 $\pm$ 0.09	0.6150 $\pm$ 0.0029	0.0322 $\pm$ 0.0018
F-SWAG-Diag	<b>81.01 <math>\pm</math> 0.29</b>	<b>0.6645 <math>\pm</math> 0.0050</b>	0.0242 $\pm$ 0.0039	<b>83.50 <math>\pm</math> 0.29</b>	<b>0.5763 <math>\pm</math> 0.0120</b>	<b>0.0151 <math>\pm</math> 0.0020</b>
SWAG	79.90 $\pm$ 0.50	<b>0.6595 <math>\pm</math> 0.0019</b>	0.0587 $\pm$ 0.0048	82.23 $\pm$ 0.19	0.6078 $\pm$ 0.0006	<b>0.0113 <math>\pm</math> 0.0020</b>
F-SWAG	<b>80.93 <math>\pm</math> 0.27</b>	0.6704 $\pm$ 0.0049	<b>0.0350 <math>\pm</math> 0.0025</b>	<b>83.57 <math>\pm</math> 0.26</b>	<b>0.5757 <math>\pm</math> 0.0136</b>	0.0196 $\pm$ 0.0015



Nguyen, Van-Anh et al. (2023). “Flat Seeking Bayesian Neural Networks”. In: *arXiv preprint arXiv:2302.02713*.