

Summary

- We propose a function-space approach to variational inference in BNNs and derive a *tractable* function-space variational objective by approximating the BNN’s variational and prior distributions via linearization of the function mapping.
- This approach leads to competitive predictive accuracy and significantly improved predictive uncertainty estimates compared to related methods, including deep ensembles, the Laplace approximation, parameter-space variational inference, and Monte Carlo Dropout.

Background

- Consider data $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N = (\mathbf{X}_{\mathcal{D}}, \mathbf{y}_{\mathcal{D}})$ with inputs $\mathbf{x}_n \in \mathcal{X} \subseteq \mathbb{R}^D$ and targets $\mathbf{y}_n \in \mathcal{Y}$, where $\mathcal{Y} \subseteq \mathbb{R}^Q$ for regression and $\mathcal{Y} \subseteq \{0, 1\}^Q$ for classification tasks.
- Consider a function mapping defined by a neural network architecture given by $f: \mathcal{X} \times \mathbb{R}^P \rightarrow \mathbb{R}^Q$

Parameter-Space Variational Inference in BNNs

- Goal: Find posterior over parameters $p(\boldsymbol{\theta} | \mathcal{D})$.
- Find variationally via $\min_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\boldsymbol{\theta}}} \mathbb{D}_{\text{KL}}(q_{\boldsymbol{\theta}} \| p_{\boldsymbol{\theta}} | \mathcal{D})$,
 $\Leftrightarrow \max_{q(\boldsymbol{\theta}) \in \mathcal{Q}_{\boldsymbol{\theta}}} \{ \mathbb{E}_{q_{\boldsymbol{\theta}}}[\log p(\mathbf{y} | f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))] - \mathbb{D}_{\text{KL}}(q_{\boldsymbol{\theta}} \| p_{\boldsymbol{\theta}}) \}$

Induced Distributions over Functions

- $f(\cdot; \boldsymbol{\theta})$ is a random function induced by random vector $\boldsymbol{\theta}$.
- Prior distribution over functions (induced by $p_{\boldsymbol{\theta}}$):

$$p_{f(\cdot; \boldsymbol{\theta})}(f(\cdot; \boldsymbol{\theta})) = \int_{\mathbb{R}^P} p_{\boldsymbol{\theta}}(\boldsymbol{\theta}') \delta(f(\cdot; \boldsymbol{\theta}) - f(\cdot; \boldsymbol{\theta}')) d\boldsymbol{\theta}' \quad (1)$$

- Posterior distribution over functions (induced by $p_{\boldsymbol{\theta} | \mathcal{D}}$):

$$p_{f(\cdot; \boldsymbol{\theta}) | \mathcal{D}}(f(\cdot; \boldsymbol{\theta}) | \mathcal{D}) = \int_{\mathbb{R}^P} p_{\boldsymbol{\theta} | \mathcal{D}}(\boldsymbol{\theta}' | \mathcal{D}) \delta(f(\cdot; \boldsymbol{\theta}) - f(\cdot; \boldsymbol{\theta}')) d\boldsymbol{\theta}' \quad (2)$$

Function-Space Variational Inference in BNNs

Function-Space Variational Inference

- Goal: Find posterior over functions $p(f(\cdot; \boldsymbol{\theta}) | \mathcal{D})$.
- Find variationally via

$$\min_{q_{\boldsymbol{\theta}} \in \mathcal{Q}_{\boldsymbol{\theta}}} \mathbb{D}_{\text{KL}}(q_{f(\cdot; \boldsymbol{\theta})} \| p_{f(\cdot; \boldsymbol{\theta}) | \mathcal{D}}), \quad (3)$$

where

$$q_{f(\cdot; \boldsymbol{\theta})}(f(\cdot; \boldsymbol{\theta})) = \int_{\mathbb{R}^P} q_{\boldsymbol{\theta}}(\boldsymbol{\theta}') \delta(f(\cdot; \boldsymbol{\theta}) - f(\cdot; \boldsymbol{\theta}')) d\boldsymbol{\theta}' \quad (4)$$

- Data Processing Inequality (Polyanskiy and Wu, 2017):

$$\mathbb{D}_{\text{KL}}(q_{f(\cdot; \boldsymbol{\theta})} \| p_{f(\cdot; \boldsymbol{\theta})}) \leq \mathbb{D}_{\text{KL}}(q_{\boldsymbol{\theta}} \| p_{\boldsymbol{\theta}}) \quad (5)$$

- Function-space variational objective:

$$\mathcal{F}(q_{\boldsymbol{\theta}}) = \mathbb{E}_{q_{f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta})}}[\log p(\mathbf{y}_{\mathcal{D}} | f(\mathbf{X}_{\mathcal{D}}; \boldsymbol{\theta}))] - \sup_{\mathbf{X} \in \mathcal{X}_{\mathcal{N}}} \mathbb{D}_{\text{KL}}(q_{f(\mathbf{X}; \boldsymbol{\theta})} \| p_{f(\mathbf{X}; \boldsymbol{\theta})}) \quad (6)$$

where $\mathcal{X}_{\mathcal{N}} \doteq \bigcup_{n \in \mathcal{N}} \{ \mathbf{X} \in \mathcal{X}_n \mid \mathcal{X}_n \subseteq \mathbb{R}^{n \times D} \}$.

Approximations

1. Linearize mapping:

$$f(\cdot; \boldsymbol{\theta}) \approx \tilde{f}(\cdot; \boldsymbol{\theta}) \doteq f(\cdot; \mathbf{m}) + \mathcal{J}_{\mathbf{m}}(\cdot)(\boldsymbol{\theta} - \mathbf{m}) \quad (7)$$

with $\mathcal{J}_{\mathbf{m}}(\cdot) \doteq \left. \frac{\partial f(\cdot; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\mathbf{m}}$ to get

$$\tilde{q}_{\tilde{f}(\cdot; \boldsymbol{\theta})}(\tilde{f}(\cdot; \boldsymbol{\theta})) \approx q_{f(\cdot; \boldsymbol{\theta})}(f(\cdot; \boldsymbol{\theta})) \quad (8)$$

$$\tilde{p}_{\tilde{f}(\cdot; \boldsymbol{\theta})}(\tilde{f}(\cdot; \boldsymbol{\theta})) \approx p_{f(\cdot; \boldsymbol{\theta})}(f(\cdot; \boldsymbol{\theta})) \quad (9)$$

2. Estimate supremum via maximum over finite sample:

$$\max_{\mathbf{X} \in \mathcal{X}_{\mathcal{C}}^S} I(\mathbf{X}) \approx \sup_{\mathbf{X} \in \mathcal{X}_{\mathcal{N}}} I(\mathbf{X}), \text{ where } \mathcal{X}_{\mathcal{C}}^S \doteq \{ \mathbf{X}_{\mathcal{C}}^{(i)} \}_{i=1}^S \quad (10)$$

Approximate Function-Space Variational Objective

- For $q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$, maximize

$$\bar{\mathcal{F}}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{M} \sum_{j=1}^M \log p(\mathbf{y}_{\mathcal{B}} | f(\mathbf{X}_{\mathcal{B}}; \hat{\boldsymbol{\theta}}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\epsilon}^{(j)}))) - \max_{\mathbf{X} \in \mathcal{X}_{\mathcal{C}}^S} \mathbb{D}_{\text{KL}}(\tilde{q}_{\tilde{f}(\mathbf{X}; \hat{\boldsymbol{\theta}})} \| \tilde{p}_{\tilde{f}(\mathbf{X}; \hat{\boldsymbol{\theta}})}) \quad (11)$$

where $\boldsymbol{\epsilon}^{(j)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_P)$.

Empirical Evaluation

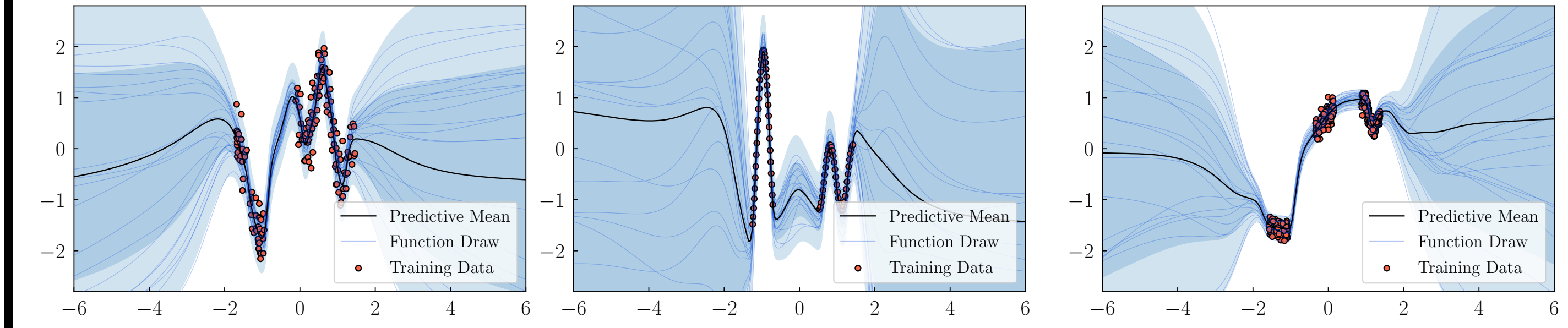


Figure 1: FSVI Posterior predictive distributions on 1D regression datasets.

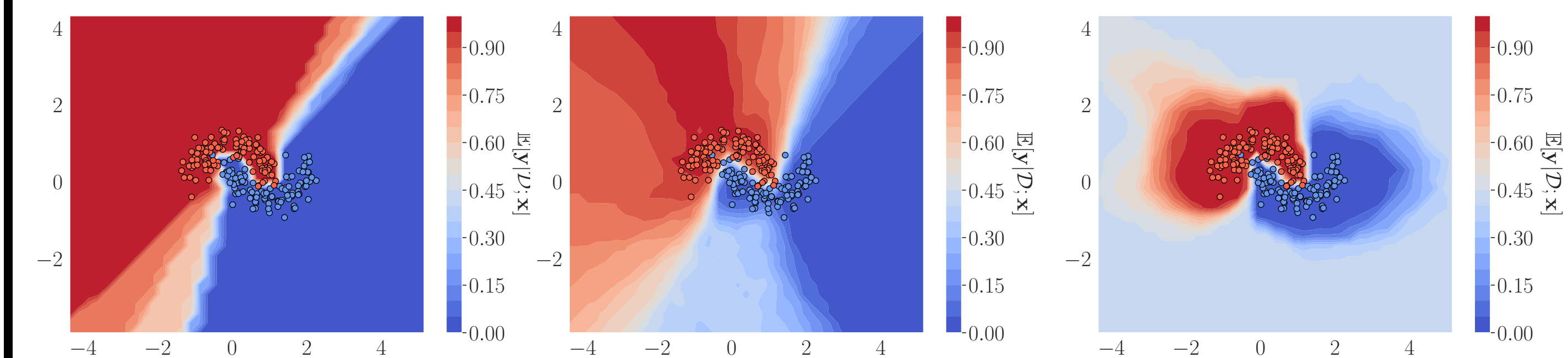


Figure 2: FSVI Posterior predictive distributions on the *Two Moons* dataset.

Table 1. Comparison of in- and out-of-distribution performance metrics (mean \pm standard error over ten random seeds).

Method	Accuracy \uparrow	ECE \downarrow	AUROC M \uparrow	AUROC NM \uparrow
MAP	91.73 \pm 0.08	0.037 \pm 0.001	87.00 \pm 0.30	74.85 \pm 1.31
MFVI	91.03 \pm 0.04	0.038 \pm 0.001	93.10 \pm 0.34	88.88 \pm 0.74
MFVI (tempered)	91.38 \pm 0.05	0.058 \pm 0.001	86.30 \pm 0.29	80.78 \pm 0.68
MFVI (radial)	90.31 \pm 0.11	0.035 \pm 0.001	84.40 \pm 0.68	82.11 \pm 1.15
MC DROPOUT	90.55 \pm 0.04	0.012 \pm 0.001	88.46 \pm 0.57	80.02 \pm 1.04
SWAG	92.56 \pm 0.05	0.043 \pm 0.001	85.18 \pm 0.35	80.31 \pm 0.30
DUQ	92.40 \pm 0.20	—	95.50 \pm 0.70	94.60 \pm 1.80
BNN-LAPLACE	92.25 \pm 0.10	0.012 \pm 0.003	95.55 \pm 0.60	—
SPG	91.60 \pm 0.14	—	95.60 \pm 6.00	—
FSVI ($p_{\mathbf{X}_{\mathcal{C}}} = \text{random monochrome}$)	92.52 \pm 0.13	0.014 \pm 0.002	96.55 \pm 0.41	95.15 \pm 0.71
FSVI ($p_{\mathbf{X}_{\mathcal{C}}} = \text{KMNIST}$)	92.67 \pm 0.15	0.012 \pm 0.002	99.81 \pm 0.19	97.44 \pm 0.24
Deep Ensemble	92.49 \pm 0.01	0.019 \pm 0.000	89.22 \pm 0.09	83.17 \pm 0.91
FSVI Ensemble	94.44 \pm 0.07	0.020 \pm 0.001	97.85 \pm 0.15	96.95 \pm 0.20

Method	Accuracy \uparrow	ECE \downarrow	OOD-AUROC \uparrow	C-CIFAR Acc \uparrow
MAP	92.19 \pm 0.15	0.046 \pm 0.001	95.17 \pm 0.40	78.55 \pm 1.01
MFVI	89.98 \pm 0.09	0.040 \pm 0.001	92.14 \pm 0.34	79.36 \pm 1.35
MFVI (tempered)	90.87 \pm 0.11	0.048 \pm 0.001	91.82 \pm 0.90	79.86 \pm 1.32
MC DROPOUT	91.32 \pm 0.06	0.041 \pm 0.001	90.32 \pm 0.57	80.19 \pm 1.44
SWAG	93.13 \pm 0.14	0.067 \pm 0.002	89.79 \pm 0.50	76.12 \pm 0.51
VOGN	84.27 \pm 0.20	0.040 \pm 0.002	87.60 \pm 0.20	—
DUQ	94.10 \pm 0.20	—	92.70 \pm 1.30	—
SPG	77.69 \pm 0.64	—	88.30 \pm 4.00	—
FSVI ($p_{\mathbf{X}_{\mathcal{C}}} = \text{random monochrome}$)	92.21 \pm 0.04	0.035 \pm 0.001	94.57 \pm 0.24	80.76 \pm 0.48
FSVI ($p_{\mathbf{X}_{\mathcal{C}}} = \text{CIFAR-100}$)	92.27 \pm 0.04	0.028 \pm 0.001	98.02 \pm 0.10	81.03 \pm 0.49
Deep Ensemble	95.13 \pm 0.06	0.019 \pm 0.001	98.04 \pm 0.07	81.22 \pm 0.37
FSVI Ensemble	95.19 \pm 0.03	0.013 \pm 0.001	99.19 \pm 0.41	81.35 \pm 0.48

Full paper: <https://timrudner.com/fsvi>