

On Weight-Space Issues in Bayesian Neural Networks

Jixiang Qing

LAI, February 11, 2026

BNN: The promise and pitfalls of BNN

- "Bayesian Neural Network **will be solved** in one or two years"
 - Some well-recognized researcher in the field, 2021(ish)

BNN: The promise and pitfalls of BNN

- "Bayesian Neural Network **will be solved** in one or two years"
 - Some well-recognized researcher in the field, 2021(ish)
- "Bayesian Deep Learning is **dead**"
 - Some well-recognized researcher in Lancaster, 2025

BNN: The promise and pitfalls of BNN

- "Bayesian Neural Network **will be solved** in one or two years"
 - Some well-recognized researcher in the field, 2021(ish)
- "Bayesian Deep Learning is **dead**"
 - Some well-recognized researcher in Lancaster, 2025
- "Is Bayesian Neural Network still a thing?"
 - Bayesian Optimisation Workshops, 2025

BNN: The promise and pitfalls of BNN

- "Bayesian Neural Network **will be solved** in one or two years"
 - Some well-recognized researcher in the field, 2021(ish)
- "Bayesian Deep Learning is **dead**"
 - Some well-recognized researcher in Lancaster, 2025
- "Is Bayesian Neural Network still a thing?"
 - Bayesian Optimisation Workshops, 2025
- "Do Bayesian Neural Networks **Actually Behave** Like **Bayesian** Models?"
 - Tom Rainforth and his colleagues, 2025

Part1: A review of BNN (Bayesian Scheme, Approximation Inference and Parameterization)

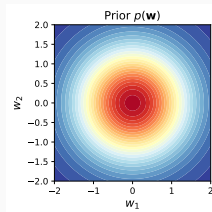
Part2: A Geometric View & Approach Enhance BNN's Performance in Weight Space

Part1: A review of BNN (Bayesian Scheme, Approximation Inference and Parameterization)

The Bayesian Scheme

Model (Generative Assumption):

- Prior: $\mathbf{w} \sim p(\mathbf{w})$ (e.g., $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$)
- Likelihood: $y_i | \mathbf{x}_i, \mathbf{w} \sim p(y | f_{\mathbf{w}}(\mathbf{x}_i))$
(e.g., $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, $y_i \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \sigma^2)$)

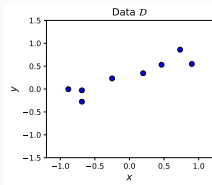
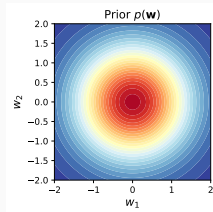


The Bayesian Scheme

Model (Generative Assumption):

- Prior: $\mathbf{w} \sim p(\mathbf{w})$ (e.g., $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$)
- Likelihood: $y_i | \mathbf{x}_i, \mathbf{w} \sim p(y | f_{\mathbf{w}}(\mathbf{x}_i))$
(e.g., $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, $y_i \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \sigma^2)$)

Setup: Observed data: $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$



The Bayesian Scheme

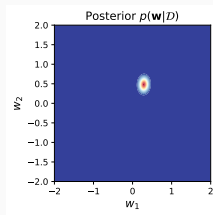
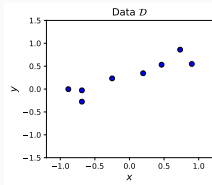
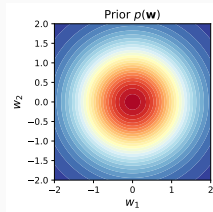
Model (Generative Assumption):

- Prior: $\mathbf{w} \sim p(\mathbf{w})$ (e.g., $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$)
- Likelihood: $y_i | \mathbf{x}_i, \mathbf{w} \sim p(y | f_{\mathbf{w}}(\mathbf{x}_i))$
(e.g., $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, $y_i \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \sigma^2)$)

Setup: Observed data: $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

Posterior: $p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$

(e.g., $\mathbf{w} | \mathcal{D} \sim \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$, closed-form)



The Bayesian Scheme

Model (Generative Assumption):

- Prior: $w \sim p(w)$ (e.g., $w \sim \mathcal{N}(0, I)$)
- Likelihood: $y_i | x_i, w \sim p(y | f_w(x_i))$
(e.g., $f_w(x) = w^\top x$, $y_i \sim \mathcal{N}(w^\top x_i, \sigma^2)$)

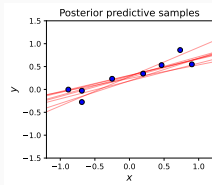
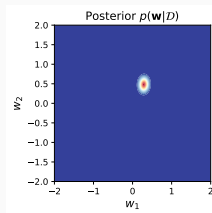
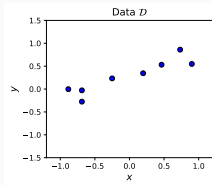
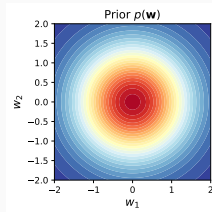
Setup: Observed data: $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Posterior: $p(w | \mathcal{D}) = \frac{p(\mathcal{D} | w)p(w)}{p(\mathcal{D})}$
(e.g., $w | \mathcal{D} \sim \mathcal{N}(\mu_N, \Sigma_N)$, closed-form)

Predictive:

$$p(y^* | x^*, \mathcal{D}) = \int p(y^* | x^*, w) p(w | \mathcal{D}) dw$$

(e.g., $y^* | x^*, \mathcal{D} \sim \mathcal{N}(\mu_N^\top x^*, \sigma^2 + x^{*\top} \Sigma_N x^*)$)



Bayesian Neural Networks

- Prior: typically $p(\boldsymbol{w}) = \mathcal{N}(0, \sigma_0^2 I)$
- Likelihood: e.g., $\mathcal{N}(y; f_{\boldsymbol{w}}(\boldsymbol{x}), \sigma^2)$ (regression)
where $f_{\boldsymbol{w}}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is a neural network: **nonlinear** w.r.t both parameters \boldsymbol{w} and input \boldsymbol{x}

Bayesian Neural Networks

- Prior: typically $p(\mathbf{w}) = \mathcal{N}(0, \sigma_0^2 I)$
- Likelihood: e.g., $\mathcal{N}(y; f_{\mathbf{w}}(\mathbf{x}), \sigma^2)$ (regression)
where $f_{\mathbf{w}}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is a neural network: **nonlinear** w.r.t both parameters \mathbf{w} and input \mathbf{x}

The Problem:

- Posterior $p(\mathbf{w}|\mathcal{D})$ has no closed form
- Predictive integral $\int p(y^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$ is intractable

Bayesian Neural Networks

- Prior: typically $p(\mathbf{w}) = \mathcal{N}(0, \sigma_0^2 I)$
- Likelihood: e.g., $\mathcal{N}(y; f_{\mathbf{w}}(\mathbf{x}), \sigma^2)$ (regression)
where $f_{\mathbf{w}}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is a neural network: **nonlinear** w.r.t both parameters \mathbf{w} and input \mathbf{x}

The Problem:

- Posterior $p(\mathbf{w}|\mathcal{D})$ has no closed form
- Predictive integral $\int p(y^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$ is intractable

Typically: \mathbf{w} is high-dimensional (millions of parameters)

Bayesian Neural Networks

- Prior: typically $p(\mathbf{w}) = \mathcal{N}(0, \sigma_0^2 I)$
- Likelihood: e.g., $\mathcal{N}(y; f_{\mathbf{w}}(\mathbf{x}), \sigma^2)$ (regression)
where $f_{\mathbf{w}}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is a neural network: **nonlinear** w.r.t both parameters \mathbf{w} and input \mathbf{x}

The Problem:

- Posterior $p(\mathbf{w}|\mathcal{D})$ has no closed form
- Predictive integral $\int p(y^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$ is intractable

Typically: \mathbf{w} is high-dimensional (millions of parameters)

\Rightarrow Need *approximate inference*!

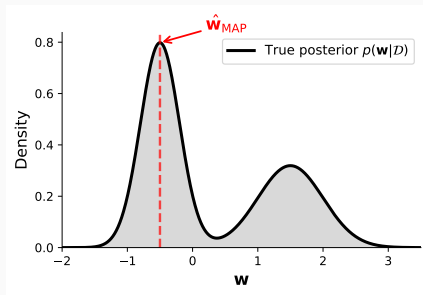
(approximate $p(\mathbf{w}|\mathcal{D})$ by sampling / simple distributions $q(\mathbf{w})$)

Laplace Approximation (MacKay, 1992)

Idea: Approximate the posterior with a Gaussian at the Maximum A Posteriori (MAP).

Step 1: Find the MAP (= training with weight decay):

$$\hat{w} = \arg \max_w [\log p(\mathcal{D}|\mathbf{w}) + \log p(\mathbf{w})]$$



Find \hat{w}_{MAP}

Laplace Approximation (MacKay, 1992)

Idea: Approximate the posterior with a Gaussian at the Maximum A Posteriori (MAP).

Step 1: Find the MAP (= training with weight decay):

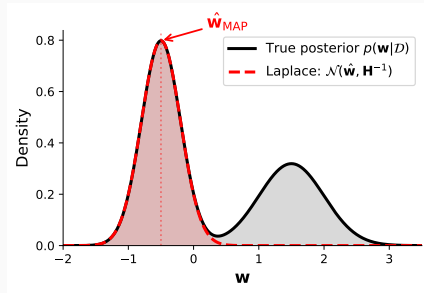
$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} [\log p(\mathcal{D}|\mathbf{w}) + \log p(\mathbf{w})]$$

Step 2: Taylor expand log-posterior around $\hat{\mathbf{w}}$:

$$\log p(\mathbf{w}|\mathcal{D}) \approx \text{const} - \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{H}(\mathbf{w} - \hat{\mathbf{w}})$$

where $\mathbf{H} = -\nabla_{\mathbf{w}}^2 \log p(\mathbf{w}|\mathcal{D})|_{\hat{\mathbf{w}}}$.

Result: $p(\mathbf{w}|\mathcal{D}) \approx \mathcal{N}(\hat{\mathbf{w}}, \mathbf{H}^{-1}) := q(\mathbf{w})$

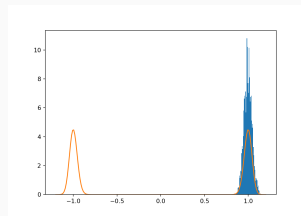


Laplace approximation of a multi-modal posterior.

Other Approximate Inference Methods

Markov Chain Monte Carlo (MCMC):

Sample $p(\boldsymbol{w}|\mathcal{D})$ asymptotically



Other Approximate Inference Methods

Markov Chain Monte Carlo (MCMC):

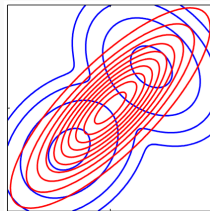
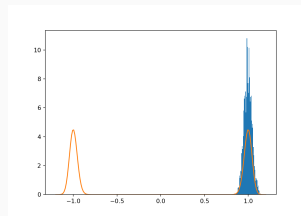
Sample $p(\mathbf{w}|\mathcal{D})$ asymptotically

Variational Inference (VI):

- Obtain a lower bound of log marginal likelihood: $\log p(\mathcal{D}) \geq$

$$\underbrace{\mathbb{E}_{q_\phi}[\log p(\mathcal{D}|\mathbf{w})] - \text{KL}(q_\phi(\mathbf{w})\|p(\mathbf{w}))}_{\text{ELBO } \mathcal{L}(\phi)}$$

- $q^* = \arg \max_{q_\phi} \mathcal{L}(q_\phi)$; typically
 $q_\phi(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$



Other Approaches to Workaround the Inference Difficulty

- **MC Dropout** (Gal & Ghahramani, 2016): Use dropout at test time as approximate VI.

$$p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{T} \sum_{t=1}^T p(y^*|\mathbf{x}^*, \hat{\mathbf{w}} \odot \mathbf{z}_t), \quad z_{t,j} \sim \text{Bernoulli}(p)$$

Other Approaches to Workaround the Inference Difficulty

- **MC Dropout** (Gal & Ghahramani, 2016): Use dropout at test time as approximate VI.

$$p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{T} \sum_{t=1}^T p(y^*|\mathbf{x}^*, \hat{\mathbf{w}} \odot \mathbf{z}_t), \quad z_{t,j} \sim \text{Bernoulli}(p)$$

- **Deep Ensembles** (Lakshminarayanan et al., 2017): Train M networks independently;
 $p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M p(y^*|\mathbf{x}^*, \hat{\mathbf{w}}_m)$

Other Approaches to Workaround the Inference Difficulty

- **MC Dropout** (Gal & Ghahramani, 2016): Use dropout at test time as approximate VI.

$$p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{T} \sum_{t=1}^T p(y^*|\mathbf{x}^*, \hat{\mathbf{w}} \odot \mathbf{z}_t), \quad z_{t,j} \sim \text{Bernoulli}(p)$$

- **Deep Ensembles** (Lakshminarayanan et al., 2017): Train M networks independently;
 $p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M p(y^*|\mathbf{x}^*, \hat{\mathbf{w}}_m)$
- **SWAG** (Maddox et al., 2019): Continue SGD past convergence, fit a Gaussian to the trajectory.

$$q(\mathbf{w}) = \mathcal{N}(\bar{\mathbf{w}}, \Sigma_{\text{low-rank+diag}}), \quad \bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$$

Other Approaches to Workaround the Inference Difficulty

- **MC Dropout** (Gal & Ghahramani, 2016): Use dropout at test time as approximate VI.

$$p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{T} \sum_{t=1}^T p(y^*|\mathbf{x}^*, \hat{\mathbf{w}} \odot \mathbf{z}_t), \quad z_{t,j} \sim \text{Bernoulli}(p)$$

- **Deep Ensembles** (Lakshminarayanan et al., 2017): Train M networks independently;
 $p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M p(y^*|\mathbf{x}^*, \hat{\mathbf{w}}_m)$
- **SWAG** (Maddox et al., 2019): Continue SGD past convergence, fit a Gaussian to the trajectory.

$$q(\mathbf{w}) = \mathcal{N}(\bar{\mathbf{w}}, \Sigma_{\text{low-rank+diag}}), \quad \bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$$

- **Bayesian Last Layer** (Snoek et al., 2015): Stochasticity only on last layer

Other Approaches to Workaround the Inference Difficulty

- **MC Dropout** (Gal & Ghahramani, 2016): Use dropout at test time as approximate VI.

$$p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{T} \sum_{t=1}^T p(y^*|\mathbf{x}^*, \hat{\mathbf{w}} \odot \mathbf{z}_t), \quad z_{t,j} \sim \text{Bernoulli}(p)$$

- **Deep Ensembles** (Lakshminarayanan et al., 2017): Train M networks independently;
 $p(y^*|\mathbf{x}^*, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M p(y^*|\mathbf{x}^*, \hat{\mathbf{w}}_m)$
- **SWAG** (Maddox et al., 2019): Continue SGD past convergence, fit a Gaussian to the trajectory.

$$q(\mathbf{w}) = \mathcal{N}(\bar{\mathbf{w}}, \Sigma_{\text{low-rank+diag}}), \quad \bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$$

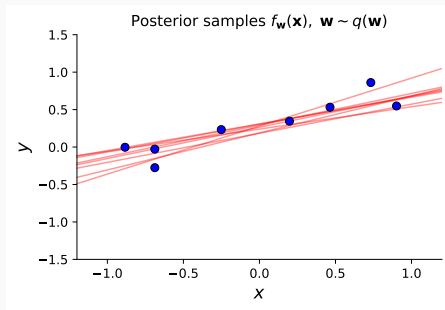
- **Bayesian Last Layer** (Snoek et al., 2015): Stochasticity only on last layer
- **Linearized Laplace** (Khan et al., 2019; Immer et al., 2021): Linearize $f_{\mathbf{w}}$ around $\hat{\mathbf{w}}$ and perform exact (Gaussian likelihood) or approximate (non-Gaussian likelihood) inference on the resulting linear model.

Making the Predictive Integral Tractable by Linearizing over Weights

Given $q(\mathbf{w}) \approx p(\mathbf{w}|\mathcal{D})$ (e.g., Laplace, SWAG, VI), the predictive integral remains intractable:

$$p(y^*|\mathbf{x}^*, \mathcal{D}) = \int p(y^*|\mathbf{x}^*, \mathbf{w}) q(\mathbf{w}) d\mathbf{w}$$

because $f_{\mathbf{w}}$ is **nonlinear** in \mathbf{w} .



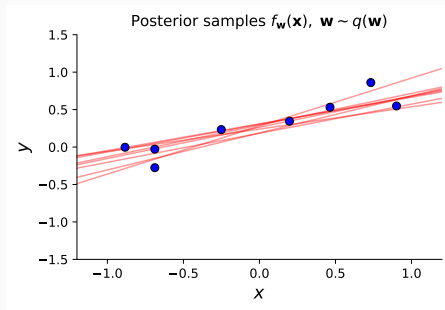
Samples from $q(\mathbf{w})$: no marginalisation.

Making the Predictive Integral Tractable by Linearizing over Weights

Given $q(\mathbf{w}) \approx p(\mathbf{w}|\mathcal{D})$ (e.g., Laplace, SWAG, VI), the predictive integral remains intractable:

$$p(y^*|\mathbf{x}^*, \mathcal{D}) = \int p(y^*|\mathbf{x}^*, \mathbf{w}) q(\mathbf{w}) d\mathbf{w}$$

because $f_{\mathbf{w}}$ is **nonlinear** in \mathbf{w} .



Samples from $q(\mathbf{w})$: no marginalisation.

Making the Predictive Integral Tractable by Linearizing over Weights

Given $q(\mathbf{w}) \approx p(\mathbf{w}|\mathcal{D})$ (e.g., Laplace, SWAG, VI), the predictive integral remains intractable:

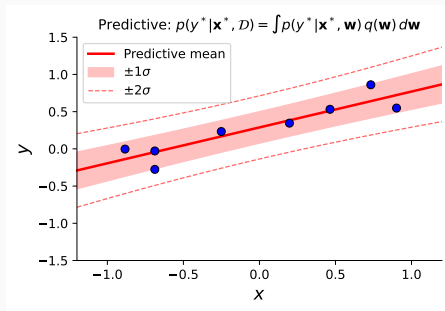
$$p(y^*|\mathbf{x}^*, \mathcal{D}) = \int p(y^*|\mathbf{x}^*, \mathbf{w}) q(\mathbf{w}) d\mathbf{w}$$

because $f_{\mathbf{w}}$ is **nonlinear** in \mathbf{w} .

Practical Approach: Linearise $f_{\mathbf{w}}$ around $\hat{\mathbf{w}}$

Khan et al. (2019); Immer et al. (2021):

$$f_{\mathbf{w}}(\mathbf{x}) \approx f_{\hat{\mathbf{w}}}(\mathbf{x}) + \mathbf{J}_{\hat{\mathbf{w}}}(\mathbf{x})(\mathbf{w} - \hat{\mathbf{w}})$$



After marginalisation: mean \pm uncertainty.

Issues in Weight-Space Inference

Another hidden problem from f_w 's parameterization

Consider a single hidden layer network with H hidden units:

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{j=1}^H v_j \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j)$$

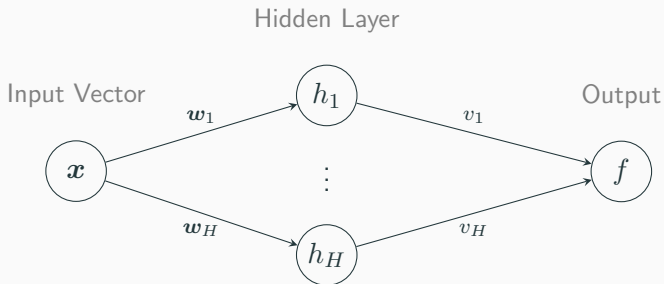
where $\mathbf{w} = (\mathbf{w}_1, b_1, v_1, \dots, \mathbf{w}_H, b_H, v_H)$ collects all parameters.

Another hidden problem from f_w 's parameterization

Consider a single hidden layer network with H hidden units:

$$f_w(\mathbf{x}) = \sum_{j=1}^H v_j \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j)$$

where $\mathbf{w} = (\mathbf{w}_1, b_1, v_1, \dots, \mathbf{w}_H, b_H, v_H)$ collects all parameters.

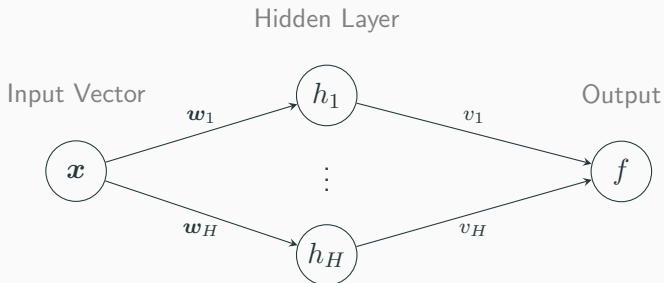


Another hidden problem from f_w 's parameterization

Consider a single hidden layer network with H hidden units:

$$f_w(\mathbf{x}) = \sum_{j=1}^H v_j \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j)$$

where $\mathbf{w} = (\mathbf{w}_1, b_1, v_1, \dots, \mathbf{w}_H, b_H, v_H)$ collects all parameters.



Key question: Can different weights produce the same function?

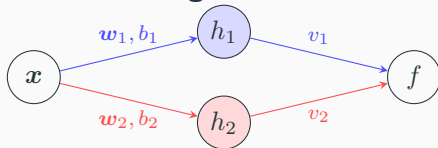
The Map $w \mapsto f_w$ Is Not Injective: Permutation Symmetry

Take $H = 2$. The two configurations below define **the same function**:

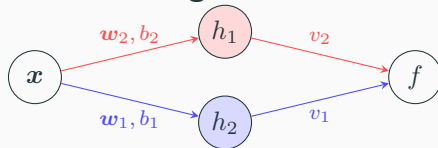
The Map $w \mapsto f_w$ Is Not Injective: Permutation Symmetry

Take $H = 2$. The two configurations below define **the same function**:

Configuration A

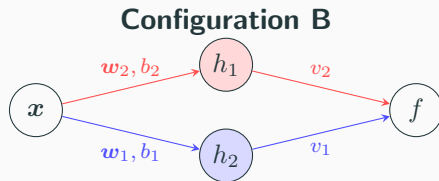
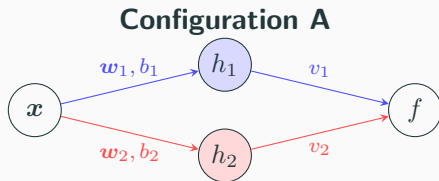


Configuration B



The Map $w \mapsto f_w$ Is Not Injective: Permutation Symmetry

Take $H = 2$. The two configurations below define **the same function**:



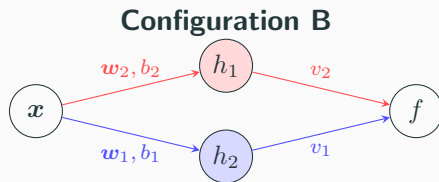
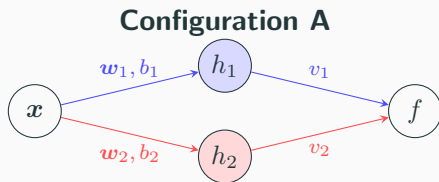
Swapping $(w_1, b_1, v_1) \leftrightarrow (w_2, b_2, v_2)$:

$$f_w(x) = v_1 \sigma(w_1^\top x + b_1) + v_2 \sigma(w_2^\top x + b_2) = f_{w'}(x)$$

Different weights $w \neq w'$, but $f_w = f_{w'}$.

The Map $w \mapsto f_w$ Is Not Injective: Permutation Symmetry

Take $H = 2$. The two configurations below define **the same function**:



Swapping $(w_1, b_1, v_1) \leftrightarrow (w_2, b_2, v_2)$:

$$f_w(x) = v_1 \sigma(w_1^\top x + b_1) + v_2 \sigma(w_2^\top x + b_2) = f_{w'}(x)$$

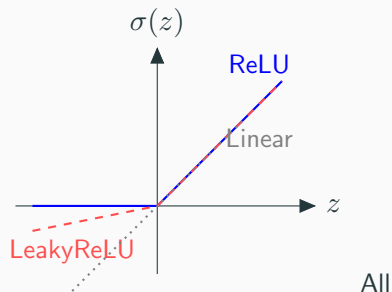
Different weights $w \neq w'$, but $f_w = f_{w'}$.

For **only one** hidden layer with H units, at least $H \times (H - 1) \times \cdots \times 1 = H!$ equivalent weight configurations.

Scaling Property of Certain Activations $\sigma(\cdot)$

Some activations **scale linearly**:

$$\sigma(\alpha z) = \alpha \sigma(z) \quad \forall \alpha > 0$$



piecewise-linear through the origin.

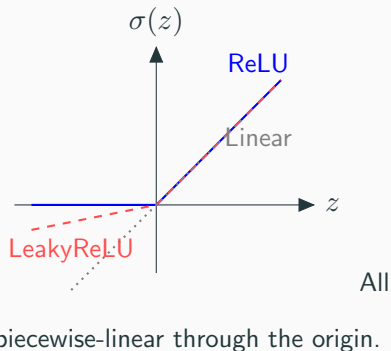
Scaling Property of Certain Activations $\sigma(\cdot)$

Some activations **scale linearly**:

$$\sigma(\alpha z) = \alpha \sigma(z) \quad \forall \alpha > 0$$

This family includes:

- ReLU: $\max(0, z)$
- LeakyReLU
- PReLU
- Linear



Scaling Symmetry

For any positively homogeneous σ , any $\alpha > 0$, and hidden unit j :

$$v_j \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j) = \underbrace{(v_j/\alpha)}_{\text{new } v_j} \sigma\left(\underbrace{(\alpha \mathbf{w}_j)^\top}_{\text{new } \mathbf{w}_j} \mathbf{x} + \underbrace{\alpha b_j}_{\text{new } b_j}\right)$$

A **continuous family** of equivalent weight configurations, parameterized by α .

Scaling Symmetry

For any positively homogeneous σ , any $\alpha > 0$, and hidden unit j :

$$v_j \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j) = \underbrace{(v_j/\alpha)}_{\text{new } v_j} \sigma\left(\underbrace{(\alpha \mathbf{w}_j)^\top}_{\text{new } \mathbf{w}_j} \mathbf{x} + \underbrace{\alpha b_j}_{\text{new } b_j}\right)$$

A **continuous family** of equivalent weight configurations, parameterized by α .

Combined consequences for the posterior:

- **Permutation** (any activation): $H!$ discrete copies of every mode per layer
- **Scaling** (homogeneous activations): continuous manifolds of equivalent weights

Scaling Symmetry

For any positively homogeneous σ , any $\alpha > 0$, and hidden unit j :

$$v_j \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j) = \underbrace{(v_j/\alpha)}_{\text{new } v_j} \sigma\left(\underbrace{(\alpha \mathbf{w}_j)^\top}_{\text{new } \mathbf{w}_j} \mathbf{x} + \underbrace{\alpha b_j}_{\text{new } b_j}\right)$$

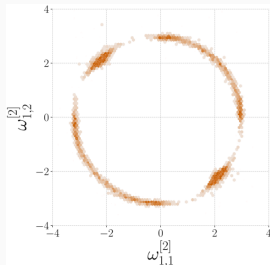
A **continuous family** of equivalent weight configurations, parameterized by α .

Combined consequences for the posterior:

- **Permutation** (any activation): $H!$ discrete copies of every mode per layer
- **Scaling** (homogeneous activations): continuous manifolds of equivalent weights

This is a property of the **parameterization**, not of the model representation capability.

What does this mean for Approximation Inference of $p(w|\mathcal{D})$

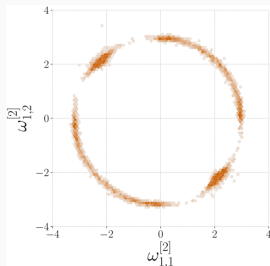


Original Posterior

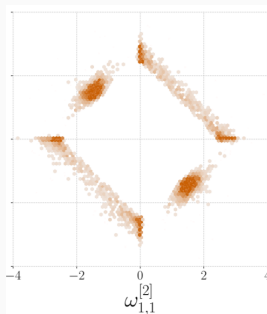
Estimated posterior over output weights v_1 vs v_2 of a 2-hidden-neuron network (Laurent et al., 2023).

1) The scaling orbit $\{(\alpha w_j, \alpha b_j, v_j/\alpha) : \alpha > 0\}$ connects equivalent weights into the ring structure; the mirror symmetry across $v_1 = v_2$ reflects the $2!$ permutation copies.

What does this mean for Approximation Inference of $p(\mathbf{w}|\mathcal{D})$



Original Posterior

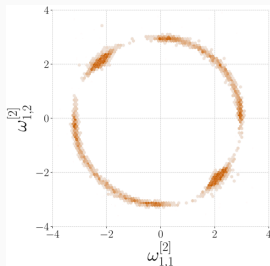


After Scaling Removal

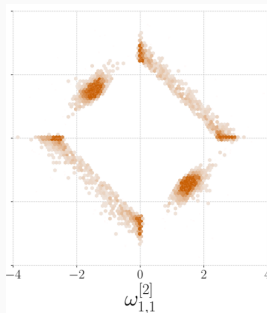
Estimated posterior over output weights v_1 vs v_2 of a 2-hidden-neuron network (Laurent et al., 2023).

1) The scaling orbit $\{(\alpha \mathbf{w}_j, \alpha b_j, v_j/\alpha) : \alpha > 0\}$ connects equivalent weights into the ring structure; the mirror symmetry across $v_1 = v_2$ reflects the $2!$ permutation copies. 2) Fixing $\|[\mathbf{w}_j; b_j]\| = 1$ determines a unique v_j , but the $2!$ copies remain.

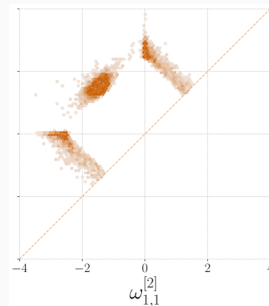
What does this mean for Approximation Inference of $p(w|\mathcal{D})$



Original Posterior



After Scaling Removal



After Scaling &
Permutation Removal

Estimated posterior over output weights v_1 vs v_2 of a 2-hidden-neuron network (Laurent et al., 2023).

- 1) The scaling orbit $\{(\alpha w_j, \alpha b_j, v_j/\alpha) : \alpha > 0\}$ connects equivalent weights into the ring structure; the mirror symmetry across $v_1 = v_2$ reflects the $2!$ permutation copies.
- 2) Fixing $\|w_j; b_j\| = 1$ determines a unique v_j , but the $2!$ copies remain.
- 3) Imposing $v_1 \geq v_2$ breaks the permutation symmetry.

Part2: A Geometric View & Approach Enhance BNN's Performance in Weight Space

A practical question

- People tried to use Laplace approximation to approximate $p(\boldsymbol{w}|\mathcal{D})$, but under the sampling scheme, the predictive uncertainty has **severe under-fit**.

A practical question

- People tried to use Laplace approximation to approximate $p(\boldsymbol{w}|\mathcal{D})$, but under the sampling scheme, the predictive uncertainty has **severe under-fit**.
- Surprisingly, when having an additional approximation over linearization of weights, under-fit has been suppressed.

A practical question

- People tried to use Laplace approximation to approximate $p(\boldsymbol{w}|\mathcal{D})$, but under the sampling scheme, the predictive uncertainty has **severe under-fit**.
- Surprisingly, when having an additional approximation over linearization of weights, under-fit has been suppressed.
- *Open Question* (Papamarkou et al., 2024): why the even crude linearization is beneficial?

Revisiting the Linearized Model (Khan et al., 2019; Immer et al., 2021)

Recall from Part 1: given a trained network with MAP estimate $\hat{\boldsymbol{w}}$, we **linearize** $f_{\boldsymbol{w}}$ w.r.t. \boldsymbol{w} around $\hat{\boldsymbol{w}}$:

$$f_{\hat{\boldsymbol{w}}}^{\text{lin}}(\boldsymbol{w}, \boldsymbol{x}) = f_{\hat{\boldsymbol{w}}}(\boldsymbol{x}) + \boldsymbol{J}_{\hat{\boldsymbol{w}}}(\boldsymbol{x}) (\boldsymbol{w} - \hat{\boldsymbol{w}})$$

where $\boldsymbol{J}_{\hat{\boldsymbol{w}}}(\boldsymbol{x}) \in \mathbb{R}^{O \times D}$ is the Jacobian (O : output dim, D : number of parameters).

Revisiting the Linearized Model (Khan et al., 2019; Immer et al., 2021)

Recall from Part 1: given a trained network with MAP estimate $\hat{\boldsymbol{w}}$, we **linearize** $f_{\boldsymbol{w}}$ w.r.t. \boldsymbol{w} around $\hat{\boldsymbol{w}}$:

$$f_{\hat{\boldsymbol{w}}}^{\text{lin}}(\boldsymbol{w}, \boldsymbol{x}) = f_{\hat{\boldsymbol{w}}}(\boldsymbol{x}) + \boldsymbol{J}_{\hat{\boldsymbol{w}}}(\boldsymbol{x}) (\boldsymbol{w} - \hat{\boldsymbol{w}})$$

where $\boldsymbol{J}_{\hat{\boldsymbol{w}}}(\boldsymbol{x}) \in \mathbb{R}^{O \times D}$ is the Jacobian (O : output dim, D : number of parameters).

In Part 1 we linearized $f_{\boldsymbol{w}}$ only in the predictive integral. Here we treat $f_{\hat{\boldsymbol{w}}}^{\text{lin}}$ as the model itself and study its structure.

Revisiting the Linearized Model (Khan et al., 2019; Immer et al., 2021)

Recall from Part 1: given a trained network with MAP estimate \hat{w} , we **linearize** f_w w.r.t. w around \hat{w} :

$$f_{\hat{w}}^{\text{lin}}(w, x) = f_{\hat{w}}(x) + J_{\hat{w}}(x) (w - \hat{w})$$

where $J_{\hat{w}}(x) \in \mathbb{R}^{O \times D}$ is the Jacobian (O : output dim, D : number of parameters).

In Part 1 we linearized f_w only in the predictive integral. Here we treat $f_{\hat{w}}^{\text{lin}}$ as the model itself and study its structure.

Key observation: For fixed x , $f_{\hat{w}}^{\text{lin}}$ is *affine in w* . A weight perturbation $\delta = w - \hat{w}$ produces a function change of $J_{\hat{w}}\delta$.

The Jacobian $J_{\hat{w}}$ is the linear map from weight perturbations to function changes.

\Rightarrow *What does this tell us about parameterization redundancy?*

Reparameterizations of Linear Functions

Consider a linear model $f(\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{b}$. Can two different weights produce the **same** function?

That is, does there exist $\mathbf{w}' \neq \mathbf{w}$ such that $f(\mathbf{w}') = f(\mathbf{w})$?

Reparameterizations of Linear Functions

Consider a linear model $f(\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{b}$. Can two different weights produce the **same** function?

That is, does there exist $\mathbf{w}' \neq \mathbf{w}$ such that $f(\mathbf{w}') = f(\mathbf{w})$? Expanding both sides:

$$\mathbf{A}\mathbf{w}' + \mathbf{b} = \mathbf{A}\mathbf{w} + \mathbf{b} \implies \mathbf{A} \underbrace{(\mathbf{w}' - \mathbf{w})}_{\delta} = \mathbf{0}$$

Reparameterizations of Linear Functions

Consider a linear model $f(\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{b}$. Can two different weights produce the **same** function?

That is, does there exist $\mathbf{w}' \neq \mathbf{w}$ such that $f(\mathbf{w}') = f(\mathbf{w})$? Expanding both sides:

$$\mathbf{A}\mathbf{w}' + \mathbf{b} = \mathbf{A}\mathbf{w} + \mathbf{b} \implies \mathbf{A} \underbrace{(\mathbf{w}' - \mathbf{w})}_{\boldsymbol{\delta}} = \mathbf{0}$$

The set of all such $\boldsymbol{\delta}$ is known as the *null space* (kernel) of \mathbf{A} :

$$\ker(\mathbf{A}) = \{\boldsymbol{\delta} \in \mathbb{R}^D : \mathbf{A}\boldsymbol{\delta} = \mathbf{0}\}$$

These are **all and only** the directions along which you can move weights without changing the function — i.e., **reparameterizations**.

Reparameterizations of Linear Functions

Consider a linear model $f(\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{b}$. Can two different weights produce the **same** function?

That is, does there exist $\mathbf{w}' \neq \mathbf{w}$ such that $f(\mathbf{w}') = f(\mathbf{w})$? Expanding both sides:

$$\mathbf{A}\mathbf{w}' + \mathbf{b} = \mathbf{A}\mathbf{w} + \mathbf{b} \implies \mathbf{A} \underbrace{(\mathbf{w}' - \mathbf{w})}_{\boldsymbol{\delta}} = \mathbf{0}$$

The set of all such $\boldsymbol{\delta}$ is known as the *null space* (kernel) of \mathbf{A} :

$$\ker(\mathbf{A}) = \{\boldsymbol{\delta} \in \mathbb{R}^D : \mathbf{A}\boldsymbol{\delta} = \mathbf{0}\}$$

These are **all and only** the directions along which you can move weights without changing the function — i.e., **reparameterizations**.

For the linearized neural network, $\mathbf{A} = \mathbf{J}_{\hat{\mathbf{w}}}$, so:

$$\ker(\mathbf{J}_{\hat{\mathbf{w}}}) = \{\text{reparameterization directions of } f_{\hat{\mathbf{w}}}^{\text{lin}}\}$$

Two Kinds of Directions in Weight Space

We just saw: directions in $\ker(\mathbf{J}_{\hat{\mathbf{w}}})$ are reparameterizations — they don't change the function.

Applying Laplace to the linearized model, the posterior is:

$$q(\mathbf{w}|\mathcal{D}) = \mathcal{N}\left(\mathbf{w} \mid \hat{\mathbf{w}}, \left(\mathbf{J}_{\hat{\mathbf{w}}}^\top \mathbf{H} \mathbf{J}_{\hat{\mathbf{w}}} + \alpha \mathbf{I}\right)^{-1}\right),$$

where $\alpha = \sigma_0^{-2}$ is the prior precision.

The precision matrix $\mathbf{J}_{\hat{\mathbf{w}}}^\top \mathbf{H} \mathbf{J}_{\hat{\mathbf{w}}}$ is known as the Generalized Gauss-Newton matrix (*GGN*)

Two Kinds of Directions in Weight Space

We just saw: directions in $\ker(\mathbf{J}_{\hat{\mathbf{w}}})$ are reparameterizations — they don't change the function.

Applying Laplace to the linearized model, the posterior is:

$$q(\mathbf{w}|\mathcal{D}) = \mathcal{N}\left(\mathbf{w} \mid \hat{\mathbf{w}}, \left(\mathbf{J}_{\hat{\mathbf{w}}}^\top \mathbf{H} \mathbf{J}_{\hat{\mathbf{w}}} + \alpha \mathbf{I}\right)^{-1}\right),$$

where $\alpha = \sigma_0^{-2}$ is the prior precision.

The precision matrix $\mathbf{J}_{\hat{\mathbf{w}}}^\top \mathbf{H} \mathbf{J}_{\hat{\mathbf{w}}}$ is known as the Generalized Gauss-Newton matrix (GGN)

(Since symmetric), this gives a clean split of the entire weight space:

$$\mathbb{R}^D = \text{im}(\text{GGN}_{\hat{\mathbf{w}}}) \oplus \ker(\text{GGN}_{\hat{\mathbf{w}}})$$

Every perturbation δ decomposes uniquely into these two orthogonal components.

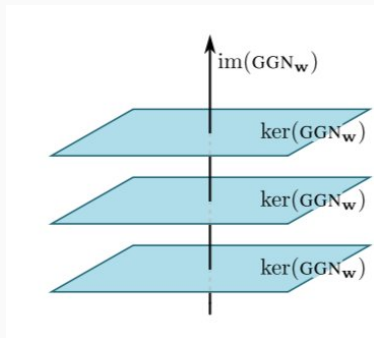
The Overparameterization Problem

Each per-input Jacobian $\mathbf{J}_{\hat{\mathbf{w}}}(\mathbf{x}_i) \in \mathbb{R}^{O \times D}$. Stacking all N data points gives an $NO \times D$ matrix, so:

$$\text{rank}(\text{GGN}_{\hat{\mathbf{w}}}) \leq NO$$

By rank-nullity, $\dim(\ker(\text{GGN}_{\hat{\mathbf{w}}})) \geq D - NO$.

Meaning: In overparameterized networks ($D \gg NO$), the **vast majority** of weight space is reparameterizations!



$\ker(\text{GGN})$: horizontal planes (vast).

$\text{im}(\text{GGN})$: vertical axis (thin).

What Does This Mean for the Posterior?

The $\text{GGN}_{\hat{w}}$ is symmetric, so by spectral decomposition:

$$\text{GGN}_{\hat{w}} = \mathbf{U}_1^\top \tilde{\mathbf{\Lambda}} \mathbf{U}_1 + \mathbf{U}_2^\top \cdot \mathbf{0} \cdot \mathbf{U}_2$$

where \mathbf{U}_1 : eigenvectors with $\tilde{\lambda}_i > 0$ (span im), \mathbf{U}_2 : eigenvectors with $\tilde{\lambda}_i = 0$ (span ker).

What Does This Mean for the Posterior?

The $\text{GGN}_{\hat{w}}$ is symmetric, so by spectral decomposition:

$$\text{GGN}_{\hat{w}} = \mathbf{U}_1^\top \tilde{\Lambda} \mathbf{U}_1 + \mathbf{U}_2^\top \cdot \mathbf{0} \cdot \mathbf{U}_2$$

where \mathbf{U}_1 : eigenvectors with $\tilde{\lambda}_i > 0$ (span im), \mathbf{U}_2 : eigenvectors with $\tilde{\lambda}_i = 0$ (span ker).

Adding $\alpha \mathbf{I}$ and inverting:

$$\Sigma = (\text{GGN}_{\hat{w}} + \alpha \mathbf{I})^{-1} = \underbrace{\mathbf{U}_1^\top (\tilde{\Lambda} + \alpha \mathbf{I})^{-1} \mathbf{U}_1}_{\text{image: constrained by data}} + \underbrace{\alpha^{-1} \mathbf{U}_2^\top \mathbf{U}_2}_{\text{kernel: prior only}}$$

What Does This Mean for the Posterior?

The $\text{GGN}_{\hat{w}}$ is symmetric, so by spectral decomposition:

$$\text{GGN}_{\hat{w}} = U_1^\top \tilde{\Lambda} U_1 + U_2^\top \cdot \mathbf{0} \cdot U_2$$

where U_1 : eigenvectors with $\tilde{\lambda}_i > 0$ (span im), U_2 : eigenvectors with $\tilde{\lambda}_i = 0$ (span ker).

Adding αI and inverting:

$$\Sigma = (\text{GGN}_{\hat{w}} + \alpha I)^{-1} = \underbrace{U_1^\top (\tilde{\Lambda} + \alpha I)^{-1} U_1}_{\text{image: constrained by data}} + \underbrace{\alpha^{-1} U_2^\top U_2}_{\text{kernel: prior only}}$$

With a weak prior (large $\sigma_0 \rightarrow$ small α), the kernel variance α^{-1} dominates. The posterior assigns large variance to directions that **do not change the function** (i.e., $\ker(\text{GGN}_{\hat{w}})$).

Two Questions

The posterior assigns large variance along reparameterization directions that locally do not change the function.

Q1: Why does the linearized Laplace approximation not suffer from this?

Q2: Why does this cause underfit when sampling from the nonlinear network?

Why Linearized Laplace Does Not Underfit

Key fact: $\ker(\text{GGN}_{\hat{w}}) = \ker(\mathbf{J}_{\hat{w}})$: the directions with prior-only variance are exactly the directions where $\mathbf{J}_{\hat{w}} \mathbf{w}_{\ker} = \mathbf{0}$.

Why Linearized Laplace Does Not Underfit

Key fact: $\ker(\text{GGN}_{\hat{w}}) = \ker(\mathbf{J}_{\hat{w}})$: the directions with prior-only variance are exactly the directions where $\mathbf{J}_{\hat{w}} \mathbf{w}_{\ker} = \mathbf{0}$.

Any sample decomposes as $\mathbf{w} = \hat{\mathbf{w}} + \mathbf{w}_{\ker} + \mathbf{w}_{\text{im}}$.

Linearized Laplace (LLA):

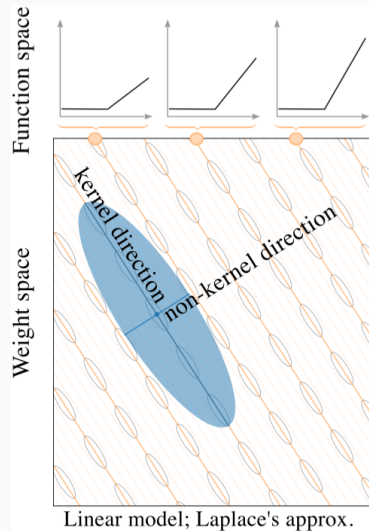
$$f^{\text{lin}}(\mathbf{w}, \mathbf{x}) = f_{\hat{\mathbf{w}}}(\mathbf{x}) + \underbrace{\mathbf{J}_{\hat{\mathbf{w}}} \mathbf{w}_{\ker}}_{=0} + \mathbf{J}_{\hat{\mathbf{w}}} \mathbf{w}_{\text{im}}$$

Kernel projected out by $\mathbf{J}_{\hat{\mathbf{w}}}$. The predictive only reflects variance in $\text{im}(\text{GGN})$.

Why Sampled Laplace Underfits

For the nonlinear network, w_{ker} is *not* filtered out:

$$f(\hat{w} + w_{\text{ker}} + w_{\text{im}}, x) \neq f(\hat{w} + w_{\text{im}}, x)$$

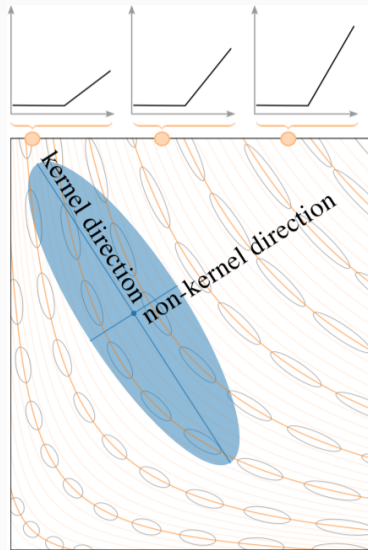


Why Sampled Laplace Underfits

For the nonlinear network, w_{ker} is *not* filtered out:

$$f(\hat{w} + w_{\text{ker}} + w_{\text{im}}, x) \neq f(\hat{w} + w_{\text{im}}, x)$$

The $\text{GGN}_{\hat{w}}$ only captures the tangent space of the reparameterization manifold locally at \hat{w} .



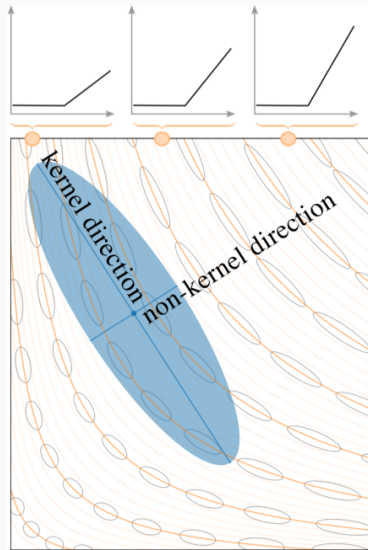
Nonlinear model; Laplace's approx.

Why Sampled Laplace Underfits

For the nonlinear network, w_{ker} is *not* filtered out:

$$f(\hat{w} + w_{\text{ker}} + w_{\text{im}}, x) \neq f(\hat{w} + w_{\text{im}}, x)$$

The $\text{GGN}_{\hat{w}}$ only captures the tangent space of the reparameterization manifold locally at \hat{w} . Large posterior variance along kernel directions pushes samples beyond the local tangent approximation and off the true reparameterization manifold.



Nonlinear model; Laplace's approx.

Why Sampled Laplace Underfits

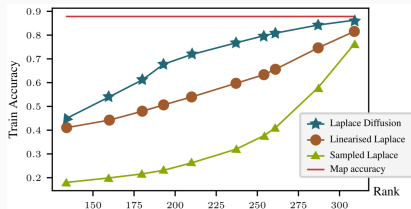
For the nonlinear network, w_{ker} is *not* filtered out:

$$f(\hat{w} + w_{\text{ker}} + w_{\text{im}}, x) \neq f(\hat{w} + w_{\text{im}}, x)$$

The $\text{GGN}_{\hat{w}}$ only captures the tangent space of the reparameterization manifold locally at \hat{w} .

Large posterior variance along kernel directions pushes samples beyond the local tangent approximation and off the true reparameterization manifold.

Sampled functions no longer reflect the training data hence underfit!



Higher GGN rank \rightarrow smaller kernel \rightarrow less underfitting. Roy et al. (2024)

Beyond the Linear Case (Roy et al., 2024)

Problem: For nonlinear networks, reparameterization sets are **curved manifolds**, not flat subspaces. A fixed Gaussian posterior cannot adapt to this curvature.

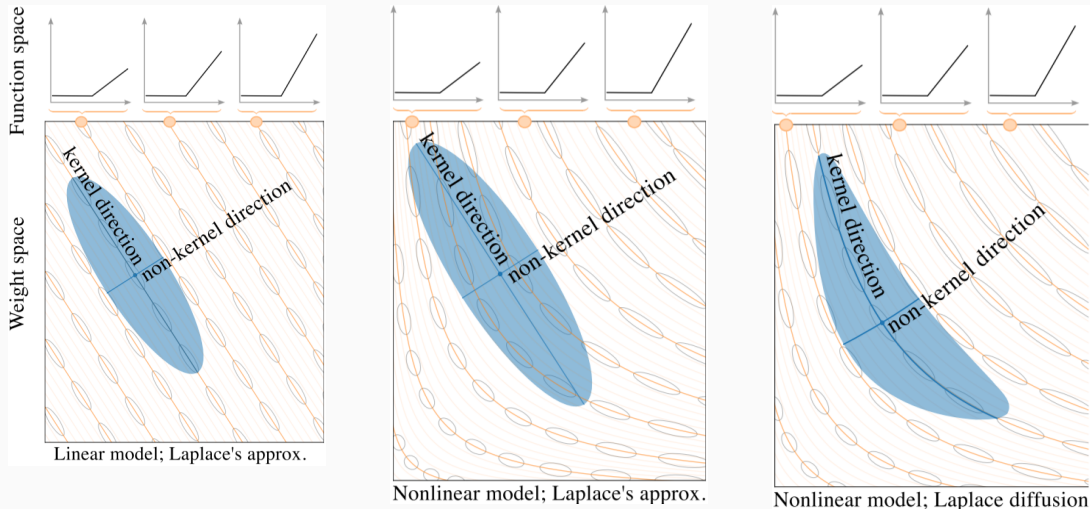
Beyond the Linear Case (Roy et al., 2024)

Problem: For nonlinear networks, reparameterization sets are **curved manifolds**, not flat subspaces. A fixed Gaussian posterior cannot adapt to this curvature.

Solution: Use the GGN as a *pseudo-Riemannian metric* it assigns zero distance to reparameterizations and adapts locally. A diffusion process under this metric gives a **reparameterization invariant posterior**.

The linearized Laplace approximation is a single-step special case of this diffusion.

Reparameterization Non-Invariance (Roy et al., 2024)



Takeaways

BNN: inference is similar to standard Bayesian statistics, however, due to its specific parameterization redundancy, this caused reparameterization challenges for the inference problem.

Takeaways

BNN: inference is similar to standard Bayesian statistics, however, due to its specific parameterization redundancy, this caused reparameterization challenges for the inference problem.

The core problem: In overparameterized networks, the vast majority of weight space consists of reparameterizations. Standard approximate posteriors waste density there.

Takeaways

BNN: inference is similar to standard Bayesian statistics, however, due to its specific parameterization redundancy, this caused reparameterization challenges for the inference problem.

The core problem: In overparameterized networks, the vast majority of weight space consists of reparameterizations. Standard approximate posteriors waste density there.

Why LLA helps: Although the LLA posterior still assign density along reparameterization direction, the predictive variance project them out due to they lie in $\ker(\mathbf{J}_{\hat{\mathbf{w}}})$. Hence does not cause any underfit issue.

Takeaways

BNN: inference is similar to standard Bayesian statistics, however, due to its specific parameterization redundancy, this caused reparameterization challenges for the inference problem.

The core problem: In overparameterized networks, the vast majority of weight space consists of reparameterizations. Standard approximate posteriors waste density there.

Why LLA helps: Although the LLA posterior still assign density along reparameterization direction, the predictive variance project them out due to they lie in $\ker(\mathbf{J}_{\hat{\mathbf{w}}})$. Hence does not cause any underfit issue.

A solution exists: Using the GGN as a pseudo-Riemannian metric yields a posterior that is reparameterization invariant by construction Roy et al. (2024).

References

- A. Immer, M. Korzepa, and M. Bauer. Improving predictions of bayesian neural nets via local linearization. In *International conference on artificial intelligence and statistics*, pages 703–711. PMLR, 2021.
- M. E. Khan, A. Immer, E. Abedi, and M. Korzepa. Approximate inference turns deep networks into gaussian processes. *Advances in neural information processing systems*, 32, 2019.
- O. Laurent, E. Aldea, and G. Franchi. A symmetry-aware exploration of bayesian neural network posteriors. *arXiv preprint arXiv:2310.08287*, 2023.
- D. J. MacKay. A practical bayesian framework for backpropagation networks. *Neural Computation*, 4 (3):448–472, 1992.

- T. Papamarkou, M. Skoularidou, K. Palla, L. Aitchison, J. Arbel, D. Dunson, M. Filippone, V. Fortuin, P. Hennig, J. M. Hernández-Lobato, et al. Position: Bayesian deep learning is needed in the age of large-scale ai. In *International Conference on Machine Learning*, 2024.
- H. Roy, M. Miani, C. H. Ek, P. Hennig, M. Pförtner, L. Tatzel, and S. Hauberg. Reparameterization invariance in approximate bayesian inference. *Advances in Neural Information Processing Systems*, 37:8132–8164, 2024.