# GVI in Function Spaces

Gaussian Measures meet Bayesian Deep Learning

Veit D. Wild* , Robert Hu* and Dino Sejdinovic

Department of Statistics

25th of May 2022

UNIVERSITY OF
OXFORD

# Outline

UNIVERSITY OF
OXFORD

# Contents

UNIVERSITY OF
OXFORD

# Bayesian Deep Learning

# Bayesian Deep Learning

Let $\mathcal{D} := \big\{ (x_n, y_n) \,|\, n = 1, \ldots, N \big\} \subset \mathcal{X} \times \mathcal{Y}$ be data.

# Bayesian Deep Learning

Let $\mathcal{D} := \left\{ (x_n, y_n) \mid n = 1, \ldots, N \right\} \subset \mathcal{X} \times \mathcal{Y}$ be data.

- Supervised Deep Learning:

$$Y = f(x) + \mathcal{N}(0, \sigma^2), \tag{1}$$

where f is a neural network $f(x) = f(x; w)$ with parameters w.

# Bayesian Deep Learning

Let $\mathcal{D} := \big\{ (x_n, y_n) \,|\, n = 1, \ldots, N \big\} \subset \mathcal{X} \times \mathcal{Y}$ be data.

- Supervised Deep Learning:

$$Y = f(x) + \mathcal{N}(0, \sigma^2), \qquad (1)$$

  where f is a neural network $f(x) = f(x; w)$ with parameters w.

- Bayesian Neural Network:
  Sample $W \sim p(w)$ and obtain random function $F(x; W)$ as prior.

# Bayesian Deep Learning

Let $\mathcal{D} := \big\{ (x_n, y_n) \,|\, n = 1, \dots, N \big\} \subset \mathcal{X} \times \mathcal{Y}$ be data.

- Supervised Deep Learning:

$$Y = f(x) + \mathcal{N}(0, \sigma^2), \tag{1}$$

where f is a neural network $f(x) = f(x; w)$ with parameters w.

- Bayesian Neural Network:
Sample $W \sim p(w)$ and obtain random function $F(x; W)$ as prior.

- Predictions for arbitrary $x^* \in \mathcal{X}$ follow from Bayes rule:

$$p(y^*|\mathcal{D}) = \int p(y^*|w) p(w|\mathcal{D}) \, dw \tag{2}$$

$$= \int p\big(y^*|f(x^*; w)\big) p(w|\mathcal{D}) \, dw \tag{3}$$

UNIVERSITY OF
OXFORD

# Bayesian Deep Learning

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

- Bayesian model averaging may improve predictive performance:

$$p(y^*|\mathcal{D}) = \int p\big(y^*|f(x^*; w)\big) p(w|\mathcal{D}) \, dw \qquad (4)$$

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

- Bayesian model averaging may improve predictive performance:

$$p(y^*|\mathcal{D}) = \int p\big(y^*|f(x^*;w)\big)p(w|\mathcal{D})\,dw \qquad (4)$$

- Bayesian posterior can be used for uncertainty quantification

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

- Bayesian model averaging may improve predictive performance:

$$p(y^*|\mathcal{D}) = \int p\big(y^*|f(x^*;w)\big)p(w|\mathcal{D})\,dw \qquad (4)$$

- Bayesian posterior can be used for uncertainty quantification

Problem: $p(w|\mathcal{D})$ is intractable! Approximations required.

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

- Bayesian model averaging may improve predictive performance:

$$p(y^*|\mathcal{D}) = \int p\big(y^*|f(x^*; w)\big)p(w|\mathcal{D})\,dw \qquad (4)$$

- Bayesian posterior can be used for uncertainty quantification

Problem: $p(w|\mathcal{D})$ is intractable! Approximations required.

Sampling based approaches:

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

- Bayesian model averaging may improve predictive performance:

$$p(y^*|\mathcal{D}) = \int p(y^*|f(x^*; w))p(w|\mathcal{D})\,dw \qquad (4)$$

- Bayesian posterior can be used for uncertainty quantification

Problem: $p(w|\mathcal{D})$ is intractable! Approximations required.

Sampling based approaches:

- Hamiltonian Monte Carlo [Neal, 2012, Chen et al., 2014]

# Bayesian Deep Learning

Why Bayesian deep learning instead of standard deep learning?

- Bayesian model averaging may improve predictive performance:

$$p(y^*|\mathcal{D}) = \int p(y^*|f(x^*;w))p(w|\mathcal{D})\,dw \qquad (4)$$

- Bayesian posterior can be used for uncertainty quantification

Problem: $p(w|\mathcal{D})$ is intractable! Approximations required.

Sampling based approaches:

- Hamiltonian Monte Carlo [Neal, 2012, Chen et al., 2014]
- Langevin Dynamics [Welling and Teh, 2011]

UNIVERSITY OF
OXFORD

# BDL: Weight-Space inference

# BDL: Weight-Space inference

Variational approach:

# BDL: Weight-Space inference

Variational approach:
Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$.

# BDL: Weight-Space inference

Variational approach:
Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{KL}\big(q(w), p(w)\big), \qquad (5)$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.

# BDL: Weight-Space inference

Variational approach:

Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{\mathrm{KL}}\big(q(w), p(w)\big),\qquad(5)$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.

Problems:

# BDL: Weight-Space inference

Variational approach:

Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{\mathrm{KL}}\big(q(w), p(w)\big), \qquad (5)$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.

Problems:

- The parameter space for w is large and the posterior multimodal.

# BDL: Weight-Space inference

Variational approach:
Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{KL}\big(q(w), p(w)\big), \tag{5}$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.
Problems:

- The parameter space for w is large and the posterior multimodal.
  $\longrightarrow$ challenging for sampling based approaches

# BDL: Weight-Space inference

Variational approach:

Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{KL}\big(q(w), p(w)\big), \qquad (5)$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.

Problems:

- The parameter space for w is large and the posterior multimodal.
  $\longrightarrow$ challenging for sampling based approaches
- Variational approaches often introduce strong assumptions for tractability.

# BDL: Weight-Space inference

Variational approach:
Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{KL}\big(q(w), p(w)\big), \qquad (5)$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.
Problems:

- The parameter space for w is large and the posterior multimodal.
  $\longrightarrow$ challenging for sampling based approaches
- Variational approaches often introduce strong assumptions for tractability.
  $\longrightarrow$ do we still capture enough of the true posterior? [Foong et al., 2020]

UNIVERSITY OF
OXFORD

# BDL: Weight-Space inference

Variational approach:
Let $q(w) = q(w; \nu)$ be a distribution with unknown parameters $\nu$. Learn $\nu$ by maximising

$$\mathcal{L}(\nu) := \mathbb{E}_{q(w)}\big[\log p(y|w)\big] - \mathbb{D}_{\mathrm{KL}}\big(q(w), p(w)\big), \tag{5}$$

which is (often) tractable. Use $q(w; \nu) \approx p(w|\mathcal{D})$.
Problems:

- The parameter space for w is large and the posterior multimodal.
  $\longrightarrow$ challenging for sampling based approaches
- Variational approaches often introduce strong assumptions for tractability.
  $\longrightarrow$ do we still capture enough of the true posterior? [Foong et al., 2020]
- What priors on the function space are induced by $p(w)$?

UNIVERSITY OF
OXFORD

# Variational Inference in Function Spaces

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^F, \mathbb{P}^F\big), \qquad (6)$$

where $\mathbb{Q}^F, \mathbb{P}^F \in \mathcal{P}(E)$ with:

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{F}, \mathbb{P}^{F}\big), \tag{6}$$

where $\mathbb{Q}^{F}, \mathbb{P}^{F} \in \mathcal{P}(\mathrm{E})$ with:

- E an infinite dimensional (Polish) function space

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big), \tag{6}$$

where $\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}} \in \mathcal{P}(\mathrm{E})$ with:

- E an infinite dimensional (Polish) function space
- $\mathcal{P}(\mathrm{E})$ the space of Borel probability measures on E

UNIVERSITY OF
OXFORD

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big), \qquad (6)$$

where $\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}} \in \mathcal{P}(\mathrm{E})$ with:

- E an infinite dimensional (Polish) function space
- $\mathcal{P}(\mathrm{E})$ the space of Borel probability measures on E

Challenges:

UNIVERSITY OF
OXFORD

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{F}, \mathbb{P}^{F}\big), \qquad (6)$$

where $\mathbb{Q}^{F}, \mathbb{P}^{F} \in \mathcal{P}(E)$ with:

- E an infinite dimensional (Polish) function space
- $\mathcal{P}(E)$ the space of Borel probability measures on E

Challenges:

- How to specify priors on infinite dimensional function spaces?

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big), \tag{6}$$

where $\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}} \in \mathcal{P}(\mathrm{E})$ with:

- E an infinite dimensional (Polish) function space
- $\mathcal{P}(\mathrm{E})$ the space of Borel probability measures on E

Challenges:

- How to specify priors on infinite dimensional function spaces?
  $\rightarrow$ Gaussian measures on Hilbert spaces

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big), \tag{6}$$

where $\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}} \in \mathcal{P}(\mathrm{E})$ with:

- E an infinite dimensional (Polish) function space
- $\mathcal{P}(\mathrm{E})$ the space of Borel probability measures on E

Challenges:

- How to specify priors on infinite dimensional function spaces?
  $\rightarrow$ Gaussian measures on Hilbert spaces
- The KL-divergence is (in general) intractable in infinite dimensions and may even be infinite [Burt et al., 2020].

UNIVERSITY OF
OXFORD

# Variational Inference in Function Spaces

Idea: perform inference in function space [Ma et al., 2019, Sun et al., 2019, Rudner et al., 2020, Ma and Hernández-Lobato, 2021]

$$\mathcal{L} = \mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] - \mathbb{D}_{\mathrm{KL}}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big), \tag{6}$$

where $\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}} \in \mathcal{P}(\mathrm{E})$ with:

- E an infinite dimensional (Polish) function space
- $\mathcal{P}(\mathrm{E})$ the space of Borel probability measures on E

Challenges:

- How to specify priors on infinite dimensional function spaces?
  $\rightarrow$ Gaussian measures on Hilbert spaces
- The KL-divergence is (in general) intractable in infinite dimensions and may even be infinite [Burt et al., 2020].
  $\rightarrow$ use generalised variational inference in infinite dimensions!

# The Rule of Three [Knoblauch et al., 2019]

# The Rule of Three [Knoblauch et al., 2019]

Find posterior as

# The Rule of Three [Knoblauch et al., 2019]

Find posterior as

$$q^*(w) := \underset{q \in \mathcal{Q}}{\arg\min} \left\{ \mathbb{E}_{q(w)} \left[ \sum_{n=1}^{N} \ell(y_n, w) \right] + D\big(q(w), p(w)\big) \right\}, \qquad (7)$$

UNIVERSITY OF
OXFORD

# The Rule of Three [Knoblauch et al., 2019]

Find posterior as

$$q^*(w) := \operatorname*{argmin}_{q \in \mathcal{Q}} \left\{ \mathbb{E}_{q(w)} \left[ \sum_{n=1}^{N} \ell(y_n, w) \right] + D\big(q(w), p(w)\big) \right\}, \qquad (7)$$

where:

# The Rule of Three [Knoblauch et al., 2019]

Find posterior as

$$q^*(w) := \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \left\{ \mathbb{E}_{q(w)}\left[ \sum_{n=1}^{N} \ell(y_n, w) \right] + D\big(q(w), p(w)\big) \right\}, \qquad (7)$$

where:

- $\mathcal{Q}$ is a set of tractable pdfs

# The Rule of Three [Knoblauch et al., 2019]

Find posterior as

$$q^*(w) := \underset{q \in \mathcal{Q}}{\arg\min} \left\{ \mathbb{E}_{q(w)} \left[ \sum_{n=1}^{N} \ell(y_n, w) \right] + D\big(q(w), p(w)\big) \right\}, \qquad (7)$$

where:

- $\mathcal{Q}$ is a set of tractable pdfs
- $\ell$ is a loss function

# The Rule of Three [Knoblauch et al., 2019]

Find posterior as

$$q^*(w) := \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \left\{ \mathbb{E}_{q(w)} \left[ \sum_{n=1}^{N} \ell(y_n, w) \right] + D\big(q(w), p(w)\big) \right\}, \qquad (7)$$

where:

- $\mathcal{Q}$ is a set of tractable pdfs
- $\ell$ is a loss function
- $D$ an arbitrary divergence

UNIVERSITY OF
OXFORD

# GVI in Function Spaces

# GVI in Function Spaces

- Idea: Use rule of three in infinite dimensional function spaces

# GVI in Function Spaces

- Idea: Use rule of three in infinite dimensional function spaces
- Theorem 1 in Knoblauch et al. [2019] holds for infinite dimensional parameter spaces

# GVI in Function Spaces

- Idea: Use rule of three in infinite dimensional function spaces
- Theorem 1 in Knoblauch et al. [2019] holds for infinite dimensional parameter spaces
- We can target

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log \mathrm{p}(\mathrm{y}|\mathrm{F})\big] + \mathbb{D}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big), \qquad (8)$$

for inference where $\mathbb{D}$ is an arbitrary divergence.

# GVI in Function Spaces

- Idea: Use rule of three in infinite dimensional function spaces
- Theorem 1 in Knoblauch et al. [2019] holds for infinite dimensional parameter spaces
- We can target

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big), \tag{8}$$

for inference where $\mathbb{D}$ is an arbitrary divergence.
- How to define priors and variational measures $\mathbb{P}^F$ and $\mathbb{Q}^F$ in infinite dimensions?

UNIVERSITY OF
OXFORD

# Gaussian Measures on Hilbert spaces

# Gaussian Measures on Hilbert spaces

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and $(H, \langle \cdot, \cdot \rangle)$ be a Hilbert space.

# Gaussian Measures on Hilbert spaces

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and $(H, \langle \cdot, \cdot \rangle)$ be a Hilbert space.

## Definition (Gaussian Random Element)

A random mapping $F : \Omega \to H$ is called Gaussian random element (GRE) if and only if

$$\langle F, h \rangle : \Omega \to \mathbb{R} \tag{9}$$

is a scalar Gaussian variable for every $h \in H$.

# Gaussian Measures on Hilbert spaces

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and $(H, \langle \cdot, \cdot \rangle)$ be a Hilbert space.

## Definition (Gaussian Random Element)

A random mapping $F : \Omega \to H$ is called Gaussian random element (GRE) if and only if

$$\langle F, h \rangle : \Omega \to \mathbb{R} \tag{9}$$

is a scalar Gaussian variable for every $h \in H$.

The mean element of $F$ is defined as

$$m := \mathbb{E}[F] := \int F(\omega) \, d\mathbb{P}(\omega) \in H$$

# Gaussian Measures on Hilbert spaces

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and $(H, \langle \cdot, \cdot \rangle)$ be a Hilbert space.

## Definition (Gaussian Random Element)

A random mapping $F : \Omega \to H$ is called Gaussian random element (GRE) if and only if

$$\langle F, h \rangle : \Omega \to \mathbb{R} \tag{9}$$

is a scalar Gaussian variable for every $h \in H$.

The mean element of $F$ is defined as

$$m := \mathbb{E}[F] := \int F(\omega) \, d\mathbb{P}(\omega) \in H \tag{10}$$

and the (linear) covariance operator $C : H \to H$ of $F$ is defined as

$$C(h) := \int \langle F(\omega), h \rangle F(\omega) \, d\mathbb{P}(\omega) - \langle m, h \rangle m, \; h \in H.$$

# Gaussian Measures on Hilbert spaces

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and $(H, \langle \cdot, \cdot \rangle)$ be a Hilbert space.

## Definition (Gaussian Random Element)

A random mapping $F : \Omega \to H$ is called Gaussian random element (GRE) if and only if

$$\langle F, h \rangle : \Omega \to \mathbb{R} \tag{9}$$

is a scalar Gaussian variable for every $h \in H$.

The mean element of $F$ is defined as

$$m := \mathbb{E}[F] := \int F(\omega) \, d\mathbb{P}(\omega) \in H \tag{10}$$

and the (linear) covariance operator $C : H \to H$ of $F$ is defined as

$$C(h) := \int \langle F(\omega), h \rangle F(\omega) \, d\mathbb{P}(\omega) - \langle m, h \rangle m, \ h \in H. \tag{11}$$

UNIVERSITY OF
OXFORD

# Gaussian Measures on Hilbert spaces

# Gaussian Measures on Hilbert spaces

By properties of the Bochner integral:

$$\langle F, h \rangle \sim \mathcal{N}\big(\langle m, h \rangle, \langle Ch, h \rangle\big), \tag{12}$$

for any $h \in H$.

# Gaussian Measures on Hilbert spaces

By properties of the Bochner integral:

$$\langle F, h \rangle \sim \mathcal{N}\big(\langle m, h \rangle, \langle Ch, h \rangle\big), \tag{12}$$

for any $h \in H$. Write $F \sim \mathcal{N}(m, C)$ for a GRE with mean element $m \in H$ and covariance operator $C$.

# Gaussian Measures on Hilbert spaces

By properties of the Bochner integral:

$$\langle F, h \rangle \sim \mathcal{N}\big(\langle m, h \rangle, \langle Ch, h \rangle\big), \tag{12}$$

for any $h \in H$. Write $F \sim \mathcal{N}(m, C)$ for a GRE with mean element $m \in H$ and covariance operator $C$.

- $C : H \to H$ of a GRE is a positive self-adjoint trace-class operator.

# Gaussian Measures on Hilbert spaces

By properties of the Bochner integral:

$$\langle F, h \rangle \sim \mathcal{N}\big(\langle m, h \rangle, \langle Ch, h \rangle\big), \tag{12}$$

for any $h \in H$. Write $F \sim \mathcal{N}(m, C)$ for a GRE with mean element $m \in H$ and covariance operator $C$.

- $C : H \to H$ of a GRE is a positive self-adjoint trace-class operator.
- For arbitrary $m \in H$ and arbitrary $C$ positive, self-adjoint and trace-class there exists a GRE such that $F \sim \mathcal{N}(m, C)$.

UNIVERSITY OF
OXFORD

# Gaussian Measures on Hilbert spaces

By properties of the Bochner integral:

$$\langle F, h \rangle \sim \mathcal{N}\big(\langle m, h \rangle, \langle Ch, h \rangle\big), \tag{12}$$

for any $h \in H$. Write $F \sim \mathcal{N}(m, C)$ for a GRE with mean element $m \in H$ and covariance operator $C$.

- $C : H \to H$ of a GRE is a positive self-adjoint trace-class operator.
- For arbitrary $m \in H$ and arbitrary $C$ positive, self-adjoint and trace-class there exists a GRE such that $F \sim \mathcal{N}(m, C)$.

## Definition (Gaussian Measure)

Let $F \sim \mathcal{N}(m, C)$ be a GRE. Then $P$ defined as

$$P(A) := \mathbb{P}^F(A) := \mathbb{P}(F \in A) \tag{13}$$

for any (measurable) $A \subset H$ is called a Gaussian measure.

# Contents

UNIVERSITY OF
OXFORD

# Model description

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big)$$

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log \mathrm{p}(\mathrm{y}|\mathrm{F})\big] + \mathbb{D}\big(\mathbb{Q}^{\mathrm{F}}, \mathbb{P}^{\mathrm{F}}\big) \qquad (14)$$

Gaussian Wasserstein Inference:

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big) \tag{14}$$

Gaussian Wasserstein Inference:

- $E = L^2(\mathcal{X}, \rho, \mathbb{R}) := \big\{ f : \mathcal{X} \to \mathbb{R} \mid \int |f(x)|^2 \, d\rho(x) < \infty \big\}$ with $\rho$ input distribution on $\mathcal{X}$

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big) \tag{14}$$

Gaussian Wasserstein Inference:

- $E = L^2(\mathcal{X}, \rho, \mathbb{R}) := \big\{ f : \mathcal{X} \to \mathbb{R} \mid \int |f(x)|^2 \, d\rho(x) < \infty \big\}$ with $\rho$ input distribution on $\mathcal{X}$
- $P := \mathbb{P}^F \sim \mathcal{N}\big(m_P, C_P\big)$

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big) \tag{14}$$

Gaussian Wasserstein Inference:

- $E = L^2(\mathcal{X}, \rho, \mathbb{R}) := \big\{f : \mathcal{X} \to \mathbb{R} \mid \int |f(x)|^2 \, d\rho(x) < \infty\big\}$ with $\rho$ input distribution on $\mathcal{X}$
- $P := \mathbb{P}^F \sim \mathcal{N}\big(m_P, C_P\big)$
- $Q := \mathbb{Q}^F \sim \mathcal{N}\big(m_Q, C_Q\big)$

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big) \tag{14}$$

Gaussian Wasserstein Inference:

- $E = L^2(\mathcal{X}, \rho, \mathbb{R}) := \big\{ f : \mathcal{X} \to \mathbb{R} \mid \int |f(x)|^2 \, d\rho(x) < \infty \big\}$ with $\rho$ input distribution on $\mathcal{X}$
- $P := \mathbb{P}^F \sim \mathcal{N}\big(m_P, C_P\big)$
- $Q := \mathbb{Q}^F \sim \mathcal{N}\big(m_Q, C_Q\big)$
- $\mathbb{D}(\cdot, \cdot) = W_2(\cdot, \cdot)$ with $W_2$ given as Wasserstein-distance

# Model description

Recall the generalised loss:

$$\mathcal{L} := -\mathbb{E}_{\mathbb{Q}}\big[\log p(y|F)\big] + \mathbb{D}\big(\mathbb{Q}^F, \mathbb{P}^F\big) \tag{14}$$

Gaussian Wasserstein Inference:

- $E = L^2(\mathcal{X}, \rho, \mathbb{R}) := \big\{ f : \mathcal{X} \to \mathbb{R} \mid \int |f(x)|^2 \, d\rho(x) < \infty \big\}$ with $\rho$ input distribution on $\mathcal{X}$
- $P := \mathbb{P}^F \sim \mathcal{N}\big(m_P, C_P\big)$
- $Q := \mathbb{Q}^F \sim \mathcal{N}\big(m_Q, C_Q\big)$
- $\mathbb{D}(\cdot, \cdot) = W_2(\cdot, \cdot)$ with $W_2$ given as Wasserstein-distance

with:

$$C_P g := \int k(\cdot, x') g(x') \, d\rho(x'), \qquad C_Q g := \int r(\cdot, x') g(x') \, d\rho(x') \tag{15}$$

for all $g \in L^2(\mathcal{X}, \rho, \mathbb{R})$ where k and r are trace-class kernels.

# Regression

# Regression

For regression:

$$p(y|F) := \prod_{n=1}^{N} p(y_n|F) := \prod_{n=1}^{N} \mathcal{N}(y_n \mid F(x_n), \sigma^2), \tag{16}$$

where $\sigma^2 > 0$.

UNIVERSITY OF
OXFORD

# Regression

For regression:

$$p(y|F) := \prod_{n=1}^{N} p(y_n|F) := \prod_{n=1}^{N} \mathcal{N}(y_n \mid F(x_n), \sigma^2), \qquad (16)$$

where $\sigma^2 > 0$.

The Wasserstein distance is tractable [Gelbrich, 1990]:

$$W_2^2(P, Q) = \|m_P - m_Q\|_2^2 + \mathrm{tr}(C_P) + \mathrm{tr}(C_Q) - 2 \cdot \mathrm{tr}\left[\left(C_P^{1/2} C_Q C_P^{1/2}\right)^{1/2}\right], \quad (17)$$

where $\mathrm{tr}(\cdot)$ denotes the trace of an operator and $C_P^{1/2}$ is the square root of the positive, self-adjoint operator $C_P$.

UNIVERSITY OF
OXFORD

# Approximation of Wasserstein distance

# Approximation of Wasserstein distance

Let $\widehat{\rho} := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}$

# Approximation of Wasserstein distance

Let $\widehat{\rho} := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}$ and notice that

$$\|m_P - m_Q\|_2^2 = \int \left( m_P(x) - m_Q(x) \right)^2 d\rho(x) \tag{18}$$

$$\approx \frac{1}{N} \sum_{n=1}^{N} \left( m_P(x_n) - m_Q(x_n) \right)^2 \tag{19}$$

# Approximation of Wasserstein distance

Let $\widehat{\rho} := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}$ and notice that

$$\|m_P - m_Q\|_2^2 = \int \big(m_P(x) - m_Q(x)\big)^2 \, d\rho(x) \tag{18}$$

$$\approx \frac{1}{N} \sum_{n=1}^{N} \big(m_P(x_n) - m_Q(x_n)\big)^2 \tag{19}$$

Further:

UNIVERSITY OF
OXFORD

# Approximation of Wasserstein distance

Let $\widehat{\rho} := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}$ and notice that

$$\|m_P - m_Q\|_2^2 = \int \left(m_P(x) - m_Q(x)\right)^2 \mathrm{d}\rho(x) \qquad (18)$$

$$\approx \frac{1}{N} \sum_{n=1}^{N} \left(m_P(x_n) - m_Q(x_n)\right)^2 \qquad (19)$$

Further:

$$\mathrm{tr}(C_P) = \int k(x,x) \, \mathrm{d}\rho(x) \approx \frac{1}{N} \sum_{n=1}^{N} k(x_n, x_n) \qquad (20)$$

$$\mathrm{tr}(C_Q) = \int r(x,x) \, \mathrm{d}\rho(x) \approx \frac{1}{N} \sum_{n=1}^{N} r(x_n, x_n) \qquad (21)$$

UNIVERSITY OF
OXFORD

# Approximation of Wasserstein distance

# Approximation of Wasserstein distance

The last term can be approximated as

$$\mathrm{tr}\left[\left(C_P^{1/2} C_Q C_P^{1/2}\right)^{1/2}\right] \approx \frac{1}{\sqrt{NN_S}} \sum_{s=1}^{N_S} \sqrt{\lambda_s\big(r(X_S, X)k(X, X_S)\big)}, \qquad (22)$$

UNIVERSITY OF
OXFORD

# Approximation of Wasserstein distance

The last term can be approximated as

$$\mathrm{tr}\left[\left(C_P^{1/2} C_Q C_P^{1/2}\right)^{1/2}\right] \approx \frac{1}{\sqrt{N N_S}} \sum_{s=1}^{N_S} \sqrt{\lambda_s\big(r(X_S, X) k(X, X_S)\big)}, \qquad (22)$$

where $X_S := (x_{S,1}, \ldots, x_{S,N_S})$, $N_S \in \mathbb{N}$ with:

$$X_{S,1}, \ldots, X_{S,N_S} \overset{\text{ind.}}{\sim} \hat{\rho} \qquad (23)$$

$$r(X_S, X) := \big(r(x_{S,s}, x_n)\big)_{s,n} \qquad (24)$$

$$k(X, X_S) := \big(k(x_n, x_{S,s})\big)_{n,s} \qquad (25)$$

and $\lambda_s\big(r(X_S, X) k(X, X_S)\big)$ denotes the s-th eigenvalue of the matrix $r(X_S, X) k(X, X_S) \in \mathbb{R}^{N_S \times N_S}$.

# Final Loss: Regression

# Final Loss: Regression

The final loss:

$$\mathcal{L} = \mathrm{L} + \widehat{\mathrm{W}}^2 \tag{26}$$

# Final Loss: Regression

The final loss:

$$\mathcal{L} = L + \widehat{W}^2 \tag{26}$$

with:

$$L := \frac{N}{2} \log(2\pi\sigma^2) + \sum_{n=1}^{N} \frac{\left(y_n - m_Q(x_n)\right)^2 + r(x_n, x_n)}{2\sigma^2} \tag{27}$$

$$\hat{W}^2 := \frac{1}{N} \sum_{n=1}^{N} \left(m_P(x_n) - m_Q(x_n)\right)^2 + \frac{1}{N} \sum_{n=1}^{N} k(x_n, x_n) \tag{28}$$

$$+ \frac{1}{N} \sum_{n=1}^{N} r(x_n, x_n) - \frac{2}{\sqrt{NN_S}} \sum_{s=1}^{N_S} \sqrt{\lambda_s\left(r(X_S, X)k(X, X_S)\right)}, \tag{29}$$

# Final Loss: Regression

The final loss:

$$\mathcal{L} = L + \widehat{W}^2 \tag{26}$$

with:

$$L := \frac{N}{2} \log(2\pi\sigma^2) + \sum_{n=1}^{N} \frac{(y_n - m_Q(x_n))^2 + r(x_n, x_n)}{2\sigma^2} \tag{27}$$

$$\hat{W}^2 := \frac{1}{N} \sum_{n=1}^{N} (m_P(x_n) - m_Q(x_n))^2 + \frac{1}{N} \sum_{n=1}^{N} k(x_n, x_n) \tag{28}$$

$$+ \frac{1}{N} \sum_{n=1}^{N} r(x_n, x_n) - \frac{2}{\sqrt{NN_S}} \sum_{s=1}^{N_S} \sqrt{\lambda_s(r(X_S, X)k(X, X_S))}, \tag{29}$$

UNIVERSITY OF
OXFORD

# Final Loss: Properties

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:

UNIVERSITY OF OXFORD

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - N evaluations of $m_Q$ and $m_P$

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - N evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N$ evaluations of $r$ and $k$

UNIVERSITY OF
OXFORD

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - $N$ evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N$ evaluations of $r$ and $k$
  - $\mathcal{O}(N + N_S^2 N + N_S^3)$ operations for the eigenvalue problem

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - $N$ evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N$ evaluations of $r$ and $k$
  - $\mathcal{O}(N + N_S^2 N + N_S^3)$ operations for the eigenvalue problem
- One evaluation of $\mathcal{L}$ in batch-mode requires:

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
    - $N$ evaluations of $m_Q$ and $m_P$
    - $N_S \cdot N$ evaluations of $r$ and $k$
    - $\mathcal{O}(N + N_S^2 N + N_S^3)$ operations for the eigenvalue problem
- One evaluation of $\mathcal{L}$ in batch-mode requires:
    - $N_B$ evaluations of $m_Q$ and $m_P$

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - N evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N$ evaluations of $r$ and $k$
  - $\mathcal{O}(N + N_S^2 N + N_S^3)$ operations for the eigenvalue problem
- One evaluation of $\mathcal{L}$ in batch-mode requires:
  - $N_B$ evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N_B$ evaluations of $r$ and $k$

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - N evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N$ evaluations of $r$ and $k$
  - $\mathcal{O}(N + N_S^2 N + N_S^3)$ operations for the eigenvalue problem
- One evaluation of $\mathcal{L}$ in batch-mode requires:
  - $N_B$ evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N_B$ evaluations of $r$ and $k$
  - $\mathcal{O}(N_B + N_S^2 N_B + N_S^3)$ operations for the eigenvalue problem

# Final Loss: Properties

- $\mathcal{L}$ is tractable for any $m_P, m_Q, k$ and $r$
- One evaluation of $\mathcal{L}$ requires:
  - $N$ evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N$ evaluations of $r$ and $k$
  - $\mathcal{O}(N + N_S^2 N + N_S^3)$ operations for the eigenvalue problem
- One evaluation of $\mathcal{L}$ in batch-mode requires:
  - $N_B$ evaluations of $m_Q$ and $m_P$
  - $N_S \cdot N_B$ evaluations of $r$ and $k$
  - $\mathcal{O}(N_B + N_S^2 N_B + N_S^3)$ operations for the eigenvalue problem
  - $\longrightarrow$ very scalable for typical $N_S, N_B << N$, e.g. $N_S = N_B = 100$

UNIVERSITY OF
OXFORD

# Recovering Other Methods

# Recovering Other Methods

- Stochastic Variational Gaussian processes (SVGP) [Titsias, 2009]:

$$m_Q(x) := m_P(x) + \sum_{m=1}^{M} \beta_m k_m(x) \tag{30}$$

$$r(x, x') := k(x, x') - k_Z(x)^T k(Z, Z)^{-1} k_Z(x) + k_Z(x)^T \Sigma k_Z(x), \tag{31}$$

where $\beta = (\beta_1, \ldots, \beta_M) \in \mathbb{R}^M$ and $\Sigma \in \mathbb{R}^{M \times M}$ are variational parameters. Further $Z = (Z_1, \ldots, Z_M)$ with $\{Z_m\}_{m=1}^{M} \overset{iid}{\sim} \widehat{\rho}$.

# Recovering Other Methods

- Stochastic Variational Gaussian processes (SVGP) [Titsias, 2009]:

$$m_Q(x) := m_P(x) + \sum_{m=1}^{M} \beta_m k_m(x) \tag{30}$$

$$r(x, x') := k(x, x') - k_Z(x)^T k(Z, Z)^{-1} k_Z(x) + k_Z(x)^T \Sigma k_Z(x), \tag{31}$$

where $\beta = (\beta_1, \ldots, \beta_M) \in \mathbb{R}^M$ and $\Sigma \in \mathbb{R}^{M \times M}$ are variational parameters. Further $Z = (Z_1, \ldots, Z_M)$ with $\{Z_m\}_{m=1}^{M} \overset{iid}{\sim} \widehat{\rho}$.

- Decoupled SVGPs [Cheng and Boots, 2017]: Same kernel r as in SVGP but mean

$$m_Q(x) := m_P(x) + \sum_{n=1}^{\tilde{N}} \beta_n k_n(x), \tag{32}$$

where $\widetilde{N} > M$.

# GWI-net

# GWI-net

Use neural net for posterior mean

# GWI-net

Use neural net for posterior mean

- Let $L \in \mathbb{N}$ be the number of hidden layers.
- Let $D_\ell$, $\ell = 0, \ldots, L + 1$ be the width of layer $\ell$ with $D_0 := D$.
- Define $g^1(x) := W^1 x + b^1$ and further

$$h^\ell(x) := \phi\big(g^\ell(x)\big), \tag{33}$$
$$g^{\ell+1}(x) := W^{\ell+1} h^\ell(x) + b^{\ell+1} \tag{34}$$

for $x \in \mathcal{X}$ where $\phi$ is an activation function.

# GWI-net

Use neural net for posterior mean

- Let $L \in \mathbb{N}$ be the number of hidden layers.
- Let $D_\ell$, $\ell = 0, \ldots, L + 1$ be the width of layer $\ell$ with $D_0 := D$.
- Define $g^1(x) := W^1 x + b^1$ and further

$$h^\ell(x) := \phi\big(g^\ell(x)\big), \tag{33}$$

$$g^{\ell+1}(x) := W^{\ell+1} h^\ell(x) + b^{\ell+1} \tag{34}$$

  for $x \in \mathcal{X}$ where $\phi$ is an activation function.

- Define

$$m_Q(x) := g^{L+1}(x) \tag{35}$$

  for $x \in \mathcal{X}$.

# GWI-net

Use neural net for posterior mean

- Let $L \in \mathbb{N}$ be the number of hidden layers.
- Let $D_\ell$, $\ell = 0, \ldots, L+1$ be the width of layer $\ell$ with $D_0 := D$.
- Define $g^1(x) := W^1 x + b^1$ and further

$$h^\ell(x) := \phi\big(g^\ell(x)\big), \tag{33}$$
$$g^{\ell+1}(x) := W^{\ell+1} h^\ell(x) + b^{\ell+1} \tag{34}$$

  for $x \in \mathcal{X}$ where $\phi$ is an activation function.
- Define
$$m_Q(x) := g^{L+1}(x) \tag{35}$$

  for $x \in \mathcal{X}$.

and the SVGP kernel r in (31) for the posterior covariance.

# Contents
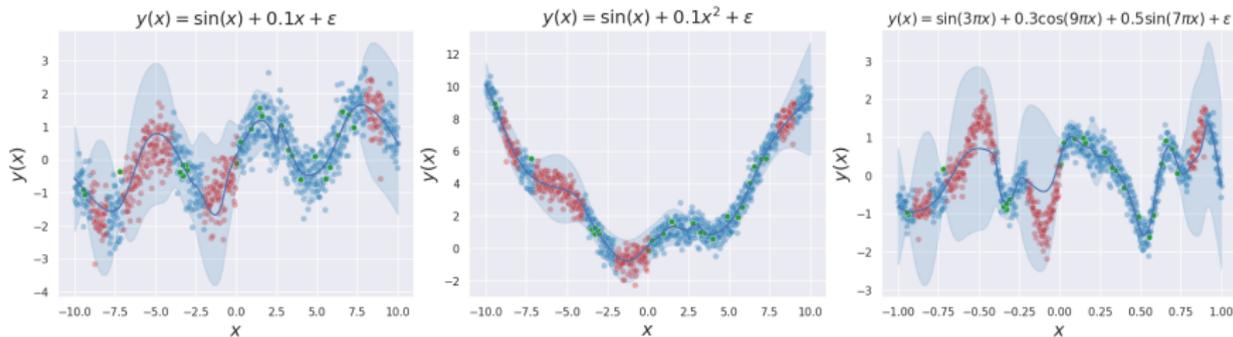
UNIVERSITY OF
OXFORD

# Toy Examples: GWI-net



Figure 1: ■ : Training data     ■ : Unseen data     ■ : Inducing points
We use N = 1000 equidistant points and add white noise with $\epsilon \sim \mathcal{N}(0, 0.5^2)$. The plot shows $m_Q(x) \pm 1.96\sqrt{\mathbb{V}[Y^*(x)|Y]}$ where $\mathbb{V}[Y^*(x)|Y]$ is the posterior predictive variance given as $r(x, x) + \sigma^2$.

# UCI Regression

# UCI Regression

| Dataset | N | D | GWI | | FVI | VIP-BNN | VIP-NP | BBB | VDO | α = 0.5 | FBNN | EXACT GP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | SVGP | DNN-SVGP | | | | | | | | |
| BOSTON | 506 | 13 | 2.8±0.5 | **2.27±0.06** | 2.33±0.04 | 2.45±0.04 | 2.45±0.03 | 2.76±0.04 | 2.63±0.10 | 2.45±0.02 | 2.30±0.10 | 2.46±0.04 |
| CONCRETE | 1030 | 8 | 3.24±0.09 | **2.64±0.06** | 2.88±0.06 | 3.02±0.02 | 3.13±0.02 | 3.28±0.01 | 3.23±0.01 | 3.06±0.03 | 3.09±0.01 | 3.05±0.02 |
| ENERGY | 768 | 8 | 1.81±0.19 | 0.91±0.12 | 0.58±0.05 | **0.56±0.04** | 0.60±0.03 | 2.17±0.02 | 1.13±0.02 | 0.95±0.09 | 0.68±0.02 | 0.54±0.02 |
| KIN8NM | 8192 | 8 | -0.86±0.38 | **-1.2±0.03** | -1.15±0.01 | -1.12±0.01 | -1.05±0.00 | -0.81±0.01 | -0.83±0.01 | -0.92±0.02 | N/A±0.00 | N/A±0.00 |
| POWER | 9568 | 4 | 3.35±0.22 | 2.74±0.02 | **2.69±0.00** | 2.92±0.00 | 2.90±0.00 | 2.83±0.01 | 2.88±0.00 | 2.81±0.00 | N/A±0.00 | N/A±0.00 |
| PROTEIN | 45730 | 9 | **2.84±0.04** | 2.87±0.0 | 2.85±0.00 | 2.87±0.00 | 2.96±0.02 | 3.00±0.00 | 2.99±0.00 | 2.90±0.00 | N/A±0.00 | N/A±0.00 |
| RED WINE | 1588 | 11 | 0.97±0.02 | **0.76±0.08** | 0.97±0.06 | 0.97±0.02 | 1.20±0.04 | 1.01±0.02 | 0.97±0.02 | 1.01±0.02 | 1.04±0.01 | 0.26±0.03 |
| YACHT | 308 | 6 | 2.37±0.55 | 0.29±0.1 | 0.59±0.11 | **-0.02±0.07** | 0.59±0.13 | 1.11±0.04 | 1.22±0.18 | 0.79±0.11 | 1.03±0.03 | 0.10±0.05 |
| NAVAL | 11934 | 16 | **-7.25±0.08** | -6.76±0.1 | -7.21±0.06 | -5.62±0.04 | -4.11±0.00 | -2.80±0.00 | -2.80±0.00 | -2.97±0.14 | -7.13±0.02 | N/A±0.00 |
| Mean Rank | | | 5.5 | **2.06** | 2.22 | 3.33 | 4.94 | 7 | 6.11 | 4.83 | | |

Table 1: The table shows the average test NLL on several UCI regression datasets. We train on random 90% of the data and predict on 10%. This is repeated 10 times and we report mean and standard deviation. The results for our competitors are taken from Ma and Hernández-Lobato [2021].

# Classification

# Classification

| | FMNIST | | | CIFAR 10 | | |
|---|---|---|---|---|---|---|
| **Model** | Accuracy | NLL | OOD-AUC | Accuracy | NLL | OOD-AUC |
| GWI-net | **93.25 ±0.09** | **0.250 ±0.00** | **0.959 ±0.01** | **83.82 ±0.00** | **0.553 ±0.00** | 0.618 ±0.00 |
| FVI | 91.60±0.14 | 0.254±0.05 | 0.956±0.06 | 77.69 ±0.64 | 0.675±0.03 | 0.883±0.04 |
| MFVI | 91.20±0.10 | 0.343±0.01 | 0.782±0.02 | 76.40±0.52 | 1.372±0.02 | 0.589±0.01 |
| MAP | 91.39±0.11 | 0.258±0.00 | 0.864±0.00 | 77.41±0.06 | 0.690±0.00 | 0.809±0.01 |
| KFAC-LAPLACE | 84.42±0.12 | 0.942±0.01 | 0.945±0.00 | 72.49±0.20 | 1.274±0.01 | 0.548±0.01 |
| RITTER et al. | 91.20±0.07 | 0.265±0.00 | 0.947±0.00 | 77.38±0.06 | 0.661±0.00 | 0.796±0.00 |

Table 2: We report average accuracy, NLL and OOD-AUC on test data for 10 different train/test splits.

UNIVERSITY OF
OXFORD

# References I

Radford M Neal. Bayesian learning for neural networks, volume 118. Springer Science & Business Media, 2012.

Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In International conference on machine learning, pages 1683–1691. PMLR, 2014.

Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In Proceedings of the 28th international conference on machine learning (ICML-11), pages 681–688. Citeseer, 2011.

Andrew Foong, David Burt, Yingzhen Li, and Richard Turner. On the expressiveness of approximate inference in bayesian neural networks. Advances in Neural Information Processing Systems, 33:15897–15908, 2020.

UNIVERSITY OF
OXFORD

# References II

Chao Ma, Yingzhen Li, and José Miguel Hernández-Lobato. Variational
implicit processes. In International Conference on Machine Learning,
pages 4222–4233. PMLR, 2019.

Shengyang Sun, Guodong Zhang, Jiaxin Shi, and Roger Grosse. Functional
variational bayesian neural networks. arXiv preprint arXiv:1903.05779,
2019.

Tim GJ Rudner, Zonghao Chen, and Yarin Gal. Rethinking function-space
variational inference in bayesian neural networks. In Third Symposium
on Advances in Approximate Bayesian Inference, 2020.

Chao Ma and José Miguel Hernández-Lobato. Functional variational
inference based on stochastic process generators. Advances in Neural
Information Processing Systems, 34, 2021.

UNIVERSITY OF
OXFORD

# References III

David R Burt, Sebastian W Ober, Adrià Garriga-Alonso, and Mark van der Wilk. Understanding variational inference in function-space. arXiv preprint arXiv:2011.09421, 2020.

Jeremias Knoblauch, Jack Jewson, and Theodoros Damoulas. Generalized variational inference: Three arguments for deriving new posteriors. arXiv preprint arXiv:1904.02063, 2019.

Matthias Gelbrich. On a formula for the l2 wasserstein metric between measures on euclidean and hilbert spaces. Mathematische Nachrichten, 147(1):185–203, 1990.

Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In Artificial intelligence and statistics, pages 567–574. PMLR, 2009.

Ching-An Cheng and Byron Boots. Variational inference for gaussian process models with linear complexity. Advances in Neural Information Processing Systems, 30, 2017.

UNIVERSITY OF
OXFORD