**G-Research Seminar** 

#### Imperial College London

# Reliable training of GP approximations

#### Mark van der Wilk

Department of Computing Imperial College London ✓ @markvanderwilk m.vdwilk@imperial.ac.uk

Jun 2, 2021

#### Goals

Main goal:

Uncertainty-aware regression for making **robust** predictions that you can **act** on.

#### Goals

Main goal:

# Uncertainty-aware regression for making **robust** predictions that you can **act** on.

The goal of the work I will discuss today:

High-quality GP approximations for **big data** that you can train **without human intervention**.

#### Goals

Main goal:

Uncertainty-aware regression for making **robust** predictions that you can **act** on.

The goal of the work I will discuss today:

# High-quality GP approximations for **big data** that you can train **without human intervention**.

... And when I mean "high-quality", I mean that we really want to squeeze every last bit of predictive performance out of the data.

#### Overview

#### The What & Why of Gaussian Processes

Variational Inference for Big Data

Theory and Practice of Training Variational GPs

Robust Conjugate Gradient Methods

Wrap-up

#### Neural networks are basis function models



$$f(\mathbf{x}) = \sum_{b=1}^{B} w_b \phi_b(\mathbf{x}) = \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x})$$
$$\phi_b(\mathbf{x}) = \sigma \left(\sum_{d=1}^{D} \tilde{w}_d x_d\right) = \sigma(\tilde{\mathbf{w}}^{\mathsf{T}} \mathbf{x})$$

Bayesian Neural Networks are a prior over functions

Placing priors on **w** gives us a distribution over functions:



$$f(\mathbf{x}) = \sum_{b=1}^{B} w_b \phi_b(\mathbf{x}) = \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x})$$
$$\mathbf{w} \sim \mathcal{N} \left(\mathbf{w}; \mathbf{0}, \sigma_w^2 I\right)$$

#### Bayesian advantages



Using the prior, we can obtain the **posterior**:

$$p(\mathbf{w}|\mathbf{y},\theta) = \frac{\prod_{n} p(y_{n}|\mathbf{w},\theta)p(\mathbf{w}|\theta)}{p(\mathbf{y}|\theta)}$$

Advantages:

- Posterior quantifies our uncertainty
- Marginal likelihood  $p(\mathbf{y}|\theta)$  gives us the hyperparameters



- How to set basis function locations?
- How many basis functions?
- Should we be so certain far from the data?



- How to set basis function locations?
- How many basis functions?
- Should we be so certain far from the data?



- How to set basis function locations?
- How many basis functions?
- Should we be so certain far from the data?



- How to set basis function locations?
- How many basis functions?
- Should we be so certain far from the data?



- How to set basis function locations?
- How many basis functions?
- Should we be so certain far from the data?

Solution: Use large number of basis functions



- How to set basis function locations?
- How many basis functions?
- Should we be so certain far from the data?

Solution: Use an infinite number of basis functions?

#### A different representation?

- ► Having many weights is expensive (cost is  $O(W^3)$ )
- Our likelihood only really depends on the function values
   *f*(**x**) = φ(**x**)<sup>T</sup>**w**
- Weights are only a means to specify the prior over functions

#### A different representation?

- Having many weights is expensive (cost is  $\mathcal{O}(W^3)$ )
- Our likelihood only really depends on the function values  $f(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} \mathbf{w}$
- Weights are only a means to specify the prior over functions

Can we specify the distribution over functions directly?

#### Distributions over functions

We can not write densities for entire functions " $p(f_{(\cdot)})$ " However, a **density over the values** of  $f_{(\cdot)}$  at **arbitrary input locations** is just as good.

$$p(f_{(\mathbf{x}_1)}, f_{(\mathbf{x}_2)}, f_{(\mathbf{x}_3)})$$



#### Distributions over functions

This allows us to always find the posterior of observations and any points to predict at.



#### Distributions over functions

This allows us to always find the posterior of observations and any points to predict at.



#### A different representation of the neural net prior

Concentrating on the function values at input locations  $X \in \mathbb{R}^{N \times D}$ :

$$p(f(X)) = \int \delta(f(X) - \Phi(X)\mathbf{w})p(\mathbf{w})d\mathbf{w}$$
$$= \mathcal{N}\left(f(X); \mathbf{0}, \sigma^2 \Phi(X)\Phi(X)^{\mathsf{T}}\right)$$
$$\left[\Phi(X)\Phi(X)^{\mathsf{T}}\right]_{nn'} = \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}')$$

#### A different representation of the neural net prior

Concentrating on the function values at input locations  $X \in \mathbb{R}^{N \times D}$ :

$$p(f(X)) = \int \delta(f(X) - \Phi(X)\mathbf{w})p(\mathbf{w})d\mathbf{w}$$
$$= \mathcal{N}\left(f(X); 0, \sigma^2 \Phi(X)\Phi(X)^{\mathsf{T}}\right)$$
$$\left[\Phi(X)\Phi(X)^{\mathsf{T}}\right]_{nn'} = \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}')$$

• We only need to define the **kernel function**.

#### A different representation of the neural net prior

Concentrating on the function values at input locations  $X \in \mathbb{R}^{N \times D}$ :

$$p(f(X)) = \int \delta(f(X) - \Phi(X)\mathbf{w})p(\mathbf{w})d\mathbf{w}$$
$$= \mathcal{N}\left(f(X); 0, \sigma^2 \Phi(X)\Phi(X)^{\mathsf{T}}\right)$$
$$\left[\Phi(X)\Phi(X)^{\mathsf{T}}\right]_{nn'} = \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}')$$

- We only need to define the **kernel function**.
- There are kernels which correspond to **infinite** basis functions.

### Three problems solved



- ► How to set basis function locations? → everywhere!
- ► How may basis functions? → infinite!
- Should we be so certain far from the data? → No, and we don't have to be!

#### Scalability

Gaussian processes are the "**gold-standard**" uncertainty-aware regression method. Training requires:

1. Selecting kernel hyperparameters

$$\boldsymbol{\theta}^{*} = \operatorname*{argmax}_{\boldsymbol{\theta}} \log p(\mathbf{y}|\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta}} \log \mathcal{N}\left(\mathbf{y}; \mathbf{0}, \mathbf{K}_{\mathrm{ff}} + \sigma^{2} \mathbf{I}_{N}\right), \quad (1)$$
$$[\mathbf{K}_{\mathrm{ff}}]_{nn'} = k(\mathbf{x}_{n}, \mathbf{x}_{n'}). \quad (2)$$

2. Calculate predictive mean at new point  $x_*$ 

$$\mu_* = \mathbf{k_{*f}}(\mathbf{K_{ff}} + \sigma^2)^{-1}\mathbf{y}$$
(3)

3. Calculate predictive variance at new point

$$\sigma_*^2 = k_{**} - \mathbf{k}_{*f} \mathbf{K}_{\mathbf{f}\mathbf{f}}^{-1} \mathbf{k}_{*f}^{\mathsf{T}}$$
(4)

#### Overview

The What & Why of Gaussian Processes

#### Variational Inference for Big Data

Theory and Practice of Training Variational GPs

Robust Conjugate Gradient Methods

Wrap-up

### Two problems

Our general model:

Prior: 
$$f_{(\cdot)} \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$$
  
Likelihood:  $y_n \sim p(y_n | f_{(\mathbf{x}_n)})$ 

We have two problems:

- ► Non-conjugacy When  $p(f_{(X)}|\mathbf{y}) = \frac{\prod_n p(y_n|f_{(\mathbf{x}_n)})p(f_{(X)})}{p(\mathbf{y})}$  is non-Gaussian.
- Scale Prior and posterior require O(N<sup>2</sup>) space and O(N<sup>3</sup>) time.

# Single solution

Variational inference:

- 1. Define a tractable set of posteriors  $q_{\eta}(f)$
- 2. Minimise the KL divergence to the true posterior  $\operatorname{KL}[q_{\eta}(f)||p(f|\mathbf{y})]$

# Single solution

Variational inference:

- 1. Define a tractable set of posteriors  $q_{\eta}(f)$
- 2. Minimise the KL divergence to the true posterior  $KL[q_{\eta}(f)||p(f|\mathbf{y})]$

Procedure:

- 1. Formulate the lower bound  $\mathcal{L} \leq \log p(\mathbf{y}|\theta)$ , such that  $\mathcal{L} + \mathrm{KL}[q_{\eta}||p] = \log p(\mathbf{y}|\theta)$
- 2. Maximise w.r.t. variational parameters  $\eta$
- 3. Also maximise w.r.t. hyperparameters  $\theta$ , for approximate maximum likelihood

# Set of approximate posteriors

#### A small GP is a tractable GP:



- A single noiseless point constrains the prior similarly to many noisy ones
- Without the redundancy, we can get very close with far fewer points

#### Set of approximate posteriors



$$q_{\eta}(f) = \frac{\prod_{m} \mathcal{N}\left(\tilde{y}_{m}; f_{(\mathbf{z}_{m})}, \tilde{\sigma}^{2}\right) p(f)}{p(\tilde{\mathbf{y}})}$$
$$= \frac{\prod_{m} \mathcal{N}\left(\tilde{y}_{m}; f_{(\mathbf{z}_{m})}, \tilde{\sigma}^{2}\right) p\left(f_{(Z)}\right)}{p(\tilde{\mathbf{y}})} p\left(f|f_{(Z)}\right)$$
$$= p\left(f|f_{(Z)}\right) p\left(f_{(Z)}|\tilde{\mathbf{y}}\right)$$

 $\eta = \{Z, \tilde{y}_m, \tilde{\sigma}\}$ 

How Accurate Gaussian Processes can help Deep Learning

### Set of approximate posteriors



Can further relax the assumption and choose an arbitrary fully correlated Gaussian likelihood:

$$q_{\eta}(f) = \frac{\prod_{m} \mathcal{N}\left(\tilde{y}_{m}; f_{(\mathbf{z}_{m})}, \Sigma\right) p(f)}{p(\tilde{\mathbf{y}})}$$
$$= \mathcal{N}\left(f_{(Z)}; \mu, \Sigma\right) p(f|f_{(Z)}) = p\left(f|f_{(Z)}\right) q(f_{(Z)})$$

 $\eta = \{Z, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ 

$$\operatorname{KL}[q(f)||p(f|\mathbf{y})] = \operatorname{KL}\left[p\left(f|f_{(X)}, f_{(Z)}\right)q\left(f_{(X)}, f_{(Z)}\right) || \\ p\left(f|f_{(X)}, f_{(Z)}\right)p\left(f_{(X)}, f_{(Z)}|\mathbf{y}\right)\right]$$

$$\operatorname{KL}[q(f)||p(f|\mathbf{y})] = \operatorname{KL}\left[p\left(f|f_{(X)}, f_{(Z)}\right)q\left(f_{(X)}, f_{(Z)}\right) || \\ p\left(f|f_{(X)}, f_{(Z)}\right)p\left(f_{(X)}, f_{(Z)}|\mathbf{y}\right)\right]$$

$$\mathrm{KL} = \log p(\mathbf{y}) - \mathbb{E}_{q(f_{(X)}, f_{(Z)})} \left[ \log \frac{\prod_{n} p(y_{n}|f_{(\mathbf{x}_{n})}) p(f_{(X)}|f_{(Z)}) p(f_{(Z)})}{p(f_{(X)}|f_{(Z)}) q(f_{(Z)})} \right]$$

$$\mathrm{KL} = \log p(\mathbf{y}) - \mathbb{E}_{q(f_{(X)}, f_{(Z)})} \left[ \log \frac{\prod_{n} p(y_{n}|f_{(\mathbf{x}_{n})}) p(f_{(X)}|f_{(Z)}) p(f_{(Z)})}{p(f_{(X)}|f_{(Z)}) q(f_{(Z)})} \right]$$

$$\mathrm{KL} = \log p(\mathbf{y}) - \mathbb{E}_{q(f(X), f(Z))} \left[ \log \frac{\prod_{n} p(y_{n}|f_{(\mathbf{x}_{n})}) p(f_{(\mathcal{X})}|f_{(Z)}) p(f_{(Z)})}{\underline{p}(f_{(\mathcal{X})}|f_{(Z)}) q(f_{(Z)})} \right]$$

$$\therefore \mathcal{L} = \sum_{n} \mathbb{E}_{p(f_{(\mathbf{x}_{n})}|f_{(Z)})q(f_{(Z)})} \left[\log p(y_{n}|f_{(\mathbf{x}_{n})})\right] - \mathrm{KL}\left[q(f_{(Z)})||p(f_{(Z)})\right]$$
#### Variational lower bound

$$\mathrm{KL} = \log p(\mathbf{y}) - \mathbb{E}_{q(f_{(X)}, f_{(Z)})} \left[ \log \frac{\prod_{n} p(y_{n}|f_{(\mathbf{x}_{n})}) p(f_{(X)}|f_{(Z)}) p(f_{(Z)})}{\underline{p}(f_{(X)}|f_{(Z)}) q(f_{(Z)})} \right]$$

$$\therefore \mathcal{L} = \sum_{n} \mathbb{E}_{p(f(\mathbf{x}_n)|f(z))q(f(z))} \Big[ \log p(y_n|f(\mathbf{x}_n)) \Big] - \mathrm{KL}\Big[q(f(z))||p(f(z))\Big]$$

- Maximising  $\mathcal{L}$  minimises the KL (gets us closer to the true posterior).
- Analytically tractable (with 1D quadrature)

#### Variational lower bound

$$\mathrm{KL} = \log p(\mathbf{y}) - \mathbb{E}_{q(f_{(X)}, f_{(Z)})} \left[ \log \frac{\prod_{n} p(y_{n}|f_{(\mathbf{x}_{n})}) p(f_{(X)}|f_{(Z)}) p(f_{(Z)})}{\underline{p}(f_{(X)}|f_{(Z)}) q(f_{(Z)})} \right]$$

$$\therefore \mathcal{L} = \sum_{n} \mathbb{E}_{p(f(\mathbf{x}_n)|f(Z))} q(f(Z)) \left[ \log p(y_n|f(\mathbf{x}_n)) \right] - \mathrm{KL} \left[ q(f(Z)) || p(f(Z)) \right]$$

- Maximising *L* minimises the KL (gets us closer to the true posterior).
- Analytically tractable (with 1D quadrature)
- Low-rank approximate GP (only needs  $M \times M$  inversion)

#### Variational lower bound

$$\mathrm{KL} = \log p(\mathbf{y}) - \mathbb{E}_{q(f_{(X)}, f_{(Z)})} \left[ \log \frac{\prod_{n} p(y_{n}|f_{(\mathbf{x}_{n})}) p(f_{(X)}|f_{(Z)}) p(f_{(Z)})}{\underline{p}(f_{(X)}|f_{(Z)}) q(f_{(Z)})} \right]$$

$$\therefore \mathcal{L} = \sum_{n} \mathbb{E}_{p(f(\mathbf{x}_n)|f(Z))} q(f(Z)) \left[ \log p(y_n|f(\mathbf{x}_n)) \right] - \mathrm{KL} \left[ q(f(Z)) || p(f(Z)) \right]$$

- Maximising *L* minimises the KL (gets us closer to the true posterior).
- Analytically tractable (with 1D quadrature)
- Low-rank approximate GP (only needs  $M \times M$  inversion)

# What should we do in practice?



Practical training requires finding  $\theta$ , *Z*,  $\mu$ ,  $\Sigma$  by maximising  $\mathcal{L}(\theta, Z, \mu, \Sigma)$ .

- 1. How to initialise the parameters? How many inducing points?
- 2. How should we optimise  $\mathcal{L}$ ?
- 3. How small can we make the KL gap to the true posterior?

#### Overview

The What & Why of Gaussian Processes

Variational Inference for Big Data

#### Theory and Practice of Training Variational GPs

Robust Conjugate Gradient Methods

Wrap-up

#### Two methods

Approach 1 (Titsias, 2009):

- Notice that the optimisation for μ, Σ is quadratic, so solve for this in closed form.
- This defines a new bound  $\mathcal{L}'(\boldsymbol{\theta}, Z)$ , with cost  $O(NM^2)$ .
- Gradient-based optimisation w.r.t.  $\theta$  and perhaps *Z*, using BFGS.

Approach 2 (Hensman et al., 2013):

- Compute  $\mathcal{L}$  and its gradients on a minibatch, with cost  $O(BM^2 + M^3)$ .
- Optimise  $\theta$ , *Z*,  $\mu$ ,  $\Sigma$  with stochastic optimiser (e.g. Adam).

# Training guidelines

Before our work, the guidelines were:

- Initialise *Z* by K-means or a random subset of the training inputs.
- List of "tricks" to try (human in the loop).
- **No guarantees** on how many inducing points are needed.

In practice:

- Optimisation would never converge. Waiting 10× longer would always give marginally better results.
- We never got *very* close to the true posterior with a sparse model.
- Approach 2 worked better, because you could do many iterations more quickly.

What was the problem?

- Problem existed in both approaches, so the problem wasn't  $\mu$ ,  $\Sigma$ .
- $\implies$  Problem must be selection of inducing inputs Z!

Draw on board: Impact of selecting bad inducing points (far from data, nearby each other).

**Rates of Convergence for Sparse Variational Gaussian Process Regression** 

David R. Burt<sup>1</sup> Carl Edward Rasmussen<sup>12</sup> Mark van der Wilk<sup>2</sup>

- We set out to analyse how many inducing inputs were needed to get  $KL \rightarrow 0$  as  $N \rightarrow \infty$ .
- Could not prove bounds for uniform subsampling.
- *Could* prove bounds for DPP selection!

Result (assuming inputs in a bounded region):

- Select inducing inputs Z with (approximate) M-DPP.
- We give a rate for *M* to increase with *N*, so that  $KL \rightarrow 0$ .
- Reasonable computational cost, e.g. O(N(log N)<sup>2D</sup>(log log N)<sup>2</sup>)) for SqExp, barely above linear<sup>1</sup>

<sup>1</sup>Recall that  $O(N(\log N)^D) = O(N^{1+\epsilon})$  for any  $D \in \mathbb{N}$  and  $\epsilon > 0$ , and that D is fixed in our problem.

# Theory to the Rescue



- *M*-DPP depends on kernel of the GP.
- *M*-DPP initialisation leads to more equally spread out inducing points than uniform.
- Theory says you don't need to optimise inducing inputs.

# This suggests a different training procedure!

# New Training Procedure

In the journal version (Burt et al., 2020) we recommend:

- 1. Use Approach 1 (Titsias, 2009), so you don't have to worry about optimising  $\mu$ ,  $\Sigma$ .
- 2. Initialise inducing inputs with *M*-DPP approximation ("greedy variance" method).
- 3. Use BFGS to optimise *only* hyperparameters  $\theta$ .
- 4. Once BFGS converges<sup>2</sup>, repeat from step 2.



# What about minibatching?



Courtesy of Maëlhann Rozé, MEng student

If you correctly initialise your inducing points:

- Stochastic optimisation is supposed to be faster than full batch... but isn't!
- Adam becomes a terrible stochastic optimiser for GPs!
- Past recommendations for SVGP were likely due to all methods performing sub-optimally due to bad inducing input placement.

- We have a *proof* that inducing points are *arbitrarily exact* at *reasonable cost* as  $N \rightarrow \infty$ .
- This leads to a new training procedure which reliably obtains better results.
- Particularly important if you care about getting every last bit of performance out of your method.

Is GP inference solved?

- We have a *proof* that inducing points are *arbitrarily exact* at *reasonable cost* as  $N \rightarrow \infty$ .
- This leads to a new training procedure which reliably obtains better results.
- Particularly important if you care about getting every last bit of performance out of your method.

#### Is GP inference solved? **No!**

 Given a dataset of fixed size, required number of inducing points may still be too large.

- We have a *proof* that inducing points are *arbitrarily exact* at *reasonable cost* as  $N \rightarrow \infty$ .
- This leads to a new training procedure which reliably obtains better results.
- Particularly important if you care about getting every last bit of performance out of your method.

#### Is GP inference solved? No!

- Given a dataset of fixed size, required number of inducing points may still be too large.
- Inducing points may not give the best speed-accuracy trade-off.

- We have a *proof* that inducing points are *arbitrarily exact* at *reasonable cost* as  $N \rightarrow \infty$ .
- This leads to a new training procedure which reliably obtains better results.
- Particularly important if you care about getting every last bit of performance out of your method.

#### Is GP inference solved? No!

- Given a dataset of fixed size, required number of inducing points may still be too large.
- Inducing points may not give the best speed-accuracy trade-off.
- Bounded/iid inputs are only one special case.

- We have a *proof* that inducing points are *arbitrarily exact* at *reasonable cost* as  $N \rightarrow \infty$ .
- This leads to a new training procedure which reliably obtains better results.
- Particularly important if you care about getting every last bit of performance out of your method.

#### Is GP inference solved? **No!**

- Given a dataset of fixed size, required number of inducing points may still be too large.
- Inducing points may not give the best speed-accuracy trade-off.
- Bounded/iid inputs are only one special case.
- We need to make running this **super easy** for people.

#### Overview

The What & Why of Gaussian Processes

Variational Inference for Big Data

Theory and Practice of Training Variational GPs

#### Robust Conjugate Gradient Methods

Wrap-up

Alternative approach to inducing points. Tries to find  $\mathbf{K}^{-1}\mathbf{v}$  by solving

$$\underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2} \mathbf{x}^{\mathsf{T}} \mathbf{K} \mathbf{x} - \mathbf{x}^{\mathsf{T}} \mathbf{b}$$
(5)

- **Conjugate Gradients** gives iterative solution that is exact in the limit.
- May give better speed-accuracy trade-off than inducing points.
- Has given genuinely impressive results<sup>3</sup>:

<sup>&</sup>lt;sup>3</sup>The results and the scale are genuinely impressive. However I disagree that they can be called exact.

Alternative approach to inducing points. Tries to find  $\mathbf{K}^{-1}\mathbf{v}$  by solving

$$\underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2} \mathbf{x}^{\mathsf{T}} \mathbf{K} \mathbf{x} - \mathbf{x}^{\mathsf{T}} \mathbf{b}$$
(5)

- Conjugate Gradients gives iterative solution that is exact in the limit.
- May give better speed-accuracy trade-off than inducing points.
- Has given genuinely impressive results<sup>3</sup>:

#### **Exact Gaussian Processes on a Million Data Points**

Ke Alexander Wang<sup>1\*</sup> Geoff Pleiss<sup>1\*</sup> Jacob R. Gardner<sup>2</sup> Stephen Tyree<sup>3</sup> Kilian Q. Weinberger<sup>1</sup> Andrew Gordon Wilson<sup>1,4</sup> <sup>1</sup>Cornell University, <sup>2</sup>Uber AI Labs, <sup>3</sup>NVIDIA, <sup>4</sup>New York University

<sup>&</sup>lt;sup>3</sup>The results and the scale are genuinely impressive. However I disagree that they can be called exact.

One should expect the training procedure of an "exact" method, to also be "exact".

One should expect the training procedure of an "exact" method, to also be "exact". Open questions about CG approach of Wang et al. (2019):

One should expect the training procedure of an "exact" method, to also be "exact".

Open questions about CG approach of Wang et al. (2019):

• How is (the convergence of) hyperparameters affected by error in the estimation of the gradients?

One should expect the training procedure of an "exact" method, to also be "exact".

Open questions about CG approach of Wang et al. (2019):

- How is (the convergence of) hyperparameters affected by error in the estimation of the gradients?
  - Wang et al. (2019) directly approximate the gradients with CG. With error introduced, it is not clear whether following them will lead to convergence.

One should expect the training procedure of an "exact" method, to also be "exact".

Open questions about CG approach of Wang et al. (2019):

- How is (the convergence of) hyperparameters affected by error in the estimation of the gradients?
  - Wang et al. (2019) directly approximate the gradients with CG. With error introduced, it is not clear whether following them will lead to convergence.
- What convergence tolerances should be used to obtain good accuracy-speed tradeoff?

One should expect the training procedure of an "exact" method, to also be "exact".

Open questions about CG approach of Wang et al. (2019):

- How is (the convergence of) hyperparameters affected by error in the estimation of the gradients?
  - Wang et al. (2019) directly approximate the gradients with CG. With error introduced, it is not clear whether following them will lead to convergence.
- What convergence tolerances should be used to obtain good accuracy-speed tradeoff?
  - An "arbitrarily exact" method would require a well-defined rule, a proof, and an asymptotic computational cost analysis, as inducing points currently has (Burt et al., 2019, 2020)

This has practical consequences, with behaviour that you would not expect from an exact method:



This said, CG is a good idea, and inspired by the results of Wang et al. (2019), we (+ Artem Artemev, David Burt) propose some improvements.

 VI measures quality of approximation and hyperparameters in a single objective: The ELBO.

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

Can we do something similar for Conjugate Gradient methods?

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

Can we do something similar for Conjugate Gradient methods? Yes!

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

Can we do something similar for Conjugate Gradient methods? Yes!

We develop the Conjugate Gradient Lower Bound (CGLB).

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

Can we do something similar for Conjugate Gradient methods? Yes!

We develop the Conjugate Gradient Lower Bound (CGLB).

The partial solution v to the CG optimisation problem to find K<sup>-1</sup>y is unified with the hyperparameter objective
 ⇒ easier to guarantee convergence

$$\theta^*, \mathbf{v}^* = \operatorname{argmax}_{\theta, \mathbf{v}} L(\theta, \mathbf{v}), \text{ with } \mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} L(\theta, \mathbf{v}) = \mathbf{K}^{-1} \mathbf{y}, \forall \theta.$$

- VI measures quality of approximation and hyperparameters in a single objective: The ELBO.
- Makes it easy to design convergent algorithms: Just optimise the ELBO!
- Q: How to set parameters? A: Just optimise the ELBO!

Can we do something similar for Conjugate Gradient methods? Yes!

We develop the Conjugate Gradient Lower Bound (CGLB).

- The partial solution v to the CG optimisation problem to find K<sup>-1</sup>y is unified with the hyperparameter objective
   ⇒ easier to guarantee convergence
- Additional upper bound to automatically determine number of CG iterations.

 $\boldsymbol{\theta}^*, \mathbf{v}^* = \operatorname{argmax}_{\boldsymbol{\theta}, \mathbf{v}} L(\boldsymbol{\theta}, \mathbf{v}), \text{ with } \mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} L(\boldsymbol{\theta}, \mathbf{v}) = \mathbf{K}^{-1} \mathbf{y}, \forall \boldsymbol{\theta}.$ 



- Fewer iterations of CG  $\implies$  faster.
- More guarantees for optimisation  $\implies$  better performance.
# Takeaways

- If you don't manually adjust parameters in "Iterative GP" (Wang et al., 2019), you can get sub-optimal performance.
- We (Artemev et al., 2021) introduce a CG-based approximation where all parameters can be automatically set through optimisation.
- This gives better and faster predictions, without requiring human intervention.

#### Overview

The What & Why of Gaussian Processes

Variational Inference for Big Data

Theory and Practice of Training Variational GPs

Robust Conjugate Gradient Methods

#### Wrap-up

# Conclusion

- Inducing point methods have theoretically-backed methods for selecting all parameters.
- With our new training methods, inducing point methods can be highly accurate.
- There are cases though, where sparse approximations simply do not exist.
- Conjugate Gradient methods can help in these situations.
- Be careful with existing methods, even if branded as "exact" (Wang et al., 2019).
- We (Artemev et al., 2021) introduce automatic solutions, with better performance.

There is a lot more to do!

Goal: Make GP approximations completely transparent.
 E.g. we don't want to need to manually choose between inducing points or CG. Software is needed for *real* impact.

There is a lot more to do!

- Goal: Make GP approximations completely transparent.
  E.g. we don't want to need to manually choose between inducing points or CG. Software is needed for *real* impact.
- Need better stochastic optimisation for GPs This is necessary for truly big data, or if we want to combine them with deep learning.

There is a lot more to do!

- Goal: Make GP approximations completely transparent.
  E.g. we don't want to need to manually choose between inducing points or CG. Software is needed for *real* impact.
- Need better stochastic optimisation for GPs This is necessary for truly big data, or if we want to combine them with deep learning.
- Deep GPs and DNNs may converge! Recent work: Dutordoir et al. (2021)

There is a lot more to do!

- Goal: Make GP approximations completely transparent.
  E.g. we don't want to need to manually choose between inducing points or CG. Software is needed for *real* impact.
- Need better stochastic optimisation for GPs This is necessary for truly big data, or if we want to combine them with deep learning.
- Deep GPs and DNNs may converge! Recent work: Dutordoir et al. (2021)

I'm currently growing a team at Imperial to work on this, and I'm looking for grants, or industry support for grant applications.

# Summary

Some papers I contributed to on robust automatic GPs

- Understanding Probabilistic Sparse Gaussian Process Approximations Bauer et al. (2016), NeurIPS
- GPflow: A Gaussian process library using TensorFlow Matthews et al. (2017), JMLR
- Convergence of Sparse Variational Inference in GPR Burt et al. (2019, 2020), ICML, JMLR
- Tighter Bounds on the LML of GPR Using CG Artemev et al. (2021), ICML
- Software for automatic & robust training of GPs Work in Progress

## References I

- Artemev, A., Burt, D. R., and van der Wilk, M. (2021). Tighter bounds on the log marginal likelihood of gaussian process regression using conjugate gradients.
- Bauer, M. S., van der Wilk, M., and Rasmussen, C. E. (2016). Understanding probabilistic sparse gaussian process approximations. In <u>Advances in</u> neural information processing systems.
- Burt, D., Rasmussen, C. E., and Van Der Wilk, M. (2019). Rates of convergence for sparse variational Gaussian process regression. In Chaudhuri, K. and Salakhutdinov, R., editors, Proceedings of the 36th International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pages 862–871. PMLR.
- Burt, D. R., Rasmussen, C. E., and van der Wilk, M. (2020). Convergence of sparse variational inference in gaussian processes regression. Journal of Machine Learning Research, 21(131):1–63.

# References II

- Dutordoir, V., Hensman, J., van der Wilk, M., Ek, C. H., Ghahramani, Z., and Durrande, N. (2021). Deep neural networks as point estimates for deep gaussian processes.
- Hensman, J., Fusi, N., and Lawrence, N. D. (2013). Gaussian processes for big data. In <u>Proceedings of the 29th Conference on Uncertainty in</u> Artificial Intelligence (UAI), pages 282–290.
- Matthews, A. G. d. G., van der Wilk, M., Nickson, T., Fujii, K., Boukouvalas, A., León-Villagrá, P., Ghahramani, Z., and Hensman, J. (2017). GPflow: A Gaussian process library using TensorFlow. Journal of Machine Learning Research, 18(40):1–6.
- Titsias, M. K. (2009). Variational learning of inducing variables in sparse Gaussian processes. In Proceedings of the 12th International Conference on Artificial Intelligence and Statistics, pages 567–574.

### References III

Wang, K., Pleiss, G., Gardner, J., Tyree, S., Weinberger, K. Q., and Wilson, A. G. (2019). Exact gaussian processes on a million data points. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, <u>Advances in Neural Information Processing Systems</u>, volume 32. Curran Associates, Inc.