

Google Brain

**Imperial College  
London**

# Data Augmentation and Infinitely Wide Neural Networks

Mark van der Wilk

Department of Computing  
Imperial College London

 @markvanderwilk  
m.vdwilk@imperial.ac.uk

Sep 21, 2021

# About our research group

- ▶ 2020–: Lecturer (Assistant Prof) at Imperial College London.
- ▶ Currently growing a research group.
- ▶ Research focus:
  - ▶ Gaussian process inference, backed by theory to make it reliable.
  - ▶ Automatic learning of inductive bias in neural networks.
  - ▶ Central question: When should neurons be connected?

## PhD Candidates



Artem Artemev



Jose Pablo Folch



Ruby Sedgwick



Seth Nabarro



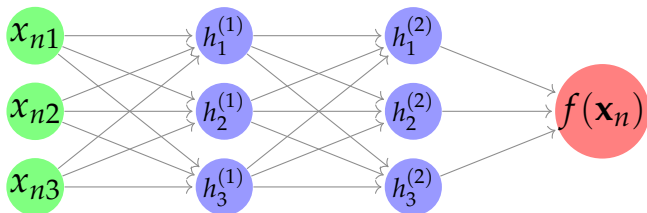
Tycho van der Ouderaa

## How can Sparse GP approximations help with studying infinitely wide NNs?

### Outline:

1. How do Sparse GPs work?
2. How accurate are Sparse GPs?
3. Sparse GPs, Data Augmentation and Invariance

# Recap: Gaussian Processes & Infinite Width NNs



- ▶ Place Gaussian prior distribution on weights.
- ▶ As number of hidden units  $\rightarrow \infty$ , we have (conditions apply)<sup>1</sup>:
  - ▶ Function values  $\{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots\}$  become jointly Gaussian.
  - ▶  $\text{Cov}[f(\mathbf{x}), f(\mathbf{x}')] = k(\mathbf{x}, \mathbf{x}')$ .
  - ▶ Kernel function depends on NN architecture, but can be computed for many!<sup>2</sup>

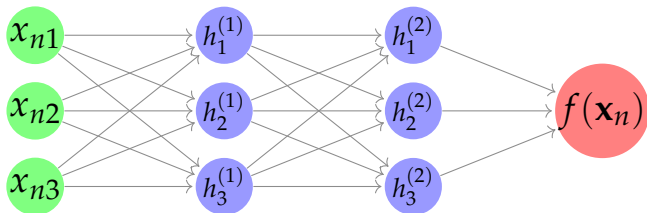
---

<sup>1</sup>Neal (1996); Matthews et al. (2018); Lee et al. (2018)

<sup>2</sup>Garriga-Alonso et al. (2019); Novak et al. (2019); Yang (2019)



## Recap: Gaussian Processes & Infinite Width NNs



- For sets of points  $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{Z} \in \mathbb{R}^{M \times D}$ , we denote the covariance of their function values as

$$\text{Cov}[f(\mathbf{X}), f(\mathbf{Z})] = \mathbf{K}_{\mathbf{XZ}} \in \mathbb{R}^{N \times M}, \quad [\mathbf{K}_{\mathbf{XZ}}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j). \quad (1)$$

- Prior on function values for any set of input points is

$$\begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix} = f(\mathbf{X}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}_{\mathbf{XX}}). \quad (2)$$

## Recap: Gaussian Processes

Given observations through some likelihood  $p(y_n|f(\mathbf{x}_n))$ , find:

1. the distribution of function values at new points  $f(\hat{\mathbf{x}})$ ,
2. the best hyperparameters  $\boldsymbol{\theta}$  of the kernel  $k_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}')$ .

Use Bayes' rule ( $\mathbf{X}'$  includes training and testing points):

$$p(f(\mathbf{X}')|\mathbf{y}) = \frac{\prod_{n=1}^N p(y_n|f(\mathbf{x}_n))p(f(\mathbf{X}')|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})} \quad (3)$$

$$= \underbrace{\frac{p(\mathbf{y}|f(\mathbf{X}))p(f(\mathbf{X}')|\boldsymbol{\theta})}{p(\mathbf{y}|\boldsymbol{\theta})}}_{p(f(\mathbf{X}')|\mathbf{y},\boldsymbol{\theta})} \cdot \underbrace{\frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})}}_{p(\boldsymbol{\theta}|\mathbf{y})} \quad (4)$$

For Gaussian likelihood  $p(y_n|f(\mathbf{x}_n)) = \mathcal{N}(y_n; f(\mathbf{x}_n), \sigma^2)$ :

1.  $p(f(\hat{\mathbf{x}})|\mathbf{y}, \boldsymbol{\theta}) = \mathcal{N}(f(\hat{\mathbf{x}}); \mathbf{K}_{\hat{\mathbf{x}}\mathbf{X}}(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma^2\mathbf{I})^{-1}\mathbf{y}, \dots)$
2.  $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \log p(\mathbf{y}|\boldsymbol{\theta}) = \log \mathcal{N}(\mathbf{y}; 0, \mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma^2\mathbf{I})$

# The Problems

1.  $p(f(\hat{\mathbf{x}})|\mathbf{y}, \boldsymbol{\theta}) = \mathcal{N}(f(\hat{\mathbf{x}}); \mathbf{K}_{\hat{\mathbf{x}}\mathbf{X}}(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma^2\mathbf{I})^{-1}\mathbf{y}, \dots)$
2.  $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \log p(\mathbf{y}|\boldsymbol{\theta}) = \log \mathcal{N}(\mathbf{y}; 0, \mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma^2\mathbf{I})$

Scalability limited by  $N \times N$  kernel matrix:

1. Storing  $N \times N$  matrix requires  $O(N^2)$  memory.
2. Inverting / log determinant takes  $O(N^3)$  time.
3. Time for calculating  $\mathbf{K}_{\mathbf{X}\mathbf{X}}$  asymptotically scales as  $O(N^2)$   
... but with huge constant, so this is the real bottleneck!

The GP side of my research develops solutions which have [guarantees on quality](#) and are [automatic](#).

# Solutions: Speed up Linear Algebra

Conjugate Gradient based solutions (Gibbs and Mackay 1997; Wang et al. 2019; Artemev, Burt, and van der Wilk, 2021)

- ▶ Speeds up inverse/logdet to  $O(N^2I)$  (how many iterations?)
- ▶ Still requires full  $\mathbf{K}_{XX}$ :  $\frac{1}{2}N^2 + N$  🙄

Nyström based solutions (Smola and Schölkopf, 2000; Williams and Seeger, 2001)

- ▶ Speeds up inverse/logdet to  $O(NM^2)$  (how big is  $M$ ?)
- ▶ Only requires  $(N + 1)M + \frac{1}{2}M^2$  kernel evals 👍

# Nyström Approximation

We want to compute 3 quantities:

1.  $c - \frac{1}{2} \log |\mathbf{K}_{XX} + \sigma^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$  (marginal likelihood)
2.  $\mathbf{K}_{\hat{\mathbf{x}}X} (\mathbf{K}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$  (pred mean)
3.  $\mathbf{K}_{\hat{\mathbf{x}}\hat{\mathbf{x}}} - \mathbf{K}_{\hat{\mathbf{x}}X} (\mathbf{K}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{X\hat{\mathbf{x}}}$  (pred variance)

Straightforward Nyström suggests:

- ▶ Select a set  $\mathbf{Z}$  with  $|\mathbf{Z}| = M \ll N$  training points
- ▶ Construct the approximation  $\mathbf{K}_{XX} \approx \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX}$
- ▶ Use Woodbury for cheap inverse approximation:

$$\begin{aligned} (\mathbf{K}_{XX} + \sigma^2 \mathbf{I})^{-1} &\approx (\mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} + \sigma^2 \mathbf{I})^{-1} \\ &= \sigma^{-2} \mathbf{I} - \sigma^{-4} \mathbf{K}_{XZ} (\mathbf{K}_{ZZ} + \sigma^{-2} \mathbf{K}_{ZX} \mathbf{K}_{XZ})^{-1} \mathbf{K}_{ZX} \end{aligned}$$

# Nyström and Inducing Variables

- ▶ Predicted variances can be negative 🙄🙄🙄
- ▶ How good is the approximation?
- ▶ When is  $M$  large enough?
- ▶ How to select the points in  $Z$ ?

In a single framework, variational **inducing variable** approximations elegantly gives:

- ▶ valid posterior approximations,
- ▶ a quantification of the quality of the approximation,
- ▶ a way to determine when  $M$  is sufficiently large,
- ▶ methods for selecting points in  $Z$ ,

as well as an approximation of  $K_{XX}$  based on Nyström.

# Variational Inference for Gaussian Processes

Problem: Computational scaling of posterior and marglik.

Three steps of Variational Inference:

1. Introduce a tractable family of variational distributions:  
GP posteriors for arbitrary Gaussian likelihoods  $\tilde{q}(\tilde{\mathbf{y}}|f(\mathbf{Z}))$

$$q(f(\hat{\mathbf{x}}), f(\mathbf{X}), f(\mathbf{Z})) = \frac{\tilde{q}(\tilde{\mathbf{y}}|f(\mathbf{Z}))p(f(\mathbf{Z}), f(\hat{\mathbf{x}}), f(\mathbf{X}))}{\tilde{q}(\tilde{\mathbf{y}})} \quad (5)$$

$$= p(f(\hat{\mathbf{x}}), f(\mathbf{X})|f(\mathbf{Z}))q(\mathbf{Z}) \quad (6)$$

2. Construct  $\mathcal{L}$  such that<sup>3</sup>  $\mathcal{L} = \log p(\mathbf{y}) - \text{KL}[q(f)||p(f|\mathbf{y})]$

$$\mathcal{L} = \sum_{n=1}^N \mathbb{E}_{q(f(\mathbf{x}_n))} [\log p(y_n|f(\mathbf{x}_n))] - \text{KL}[q(f(\mathbf{Z}))||p(f(\mathbf{Z}))] \quad (7)$$

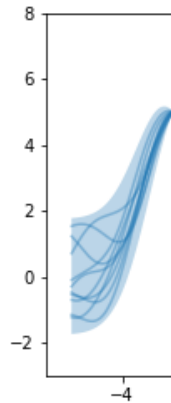
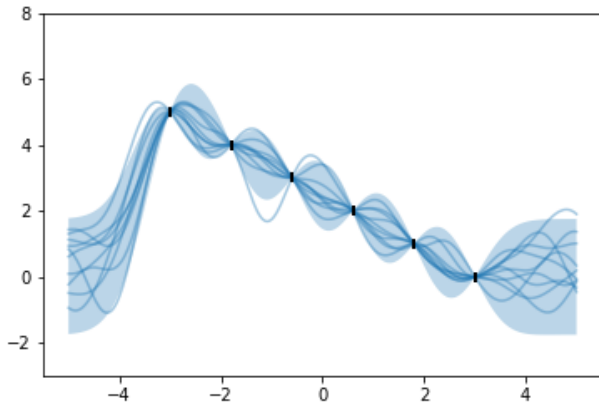
3. Minimise KL divergence by maximising  $\mathcal{L}$ !

---

<sup>3</sup>Hensman et al. (2013); Matthews (2016)

# Understanding Inducing Points

We can control  $Z$ , and  $\mu, \Sigma$  of  $q(f(Z))$ .





# Sparse GP Regression

Gaussian noise regression works particularly well, since the optimal  $\mu, \Sigma$  can be found in closed form<sup>4</sup>, giving

$$\mathcal{L} = \log \mathcal{N}(\mathbf{y}; 0, \mathbf{K}_{\mathbf{XZ}} \mathbf{K}_{\mathbf{ZZ}}^{-1} \mathbf{K}_{\mathbf{ZX}} + \sigma^2 \mathbf{I}) - \frac{1}{2\sigma^2} \text{Tr}(\mathbf{K}_{\mathbf{XX}} - \mathbf{K}_{\mathbf{XZ}} \mathbf{K}_{\mathbf{ZZ}}^{-1} \mathbf{K}_{\mathbf{ZX}}) \quad (8)$$

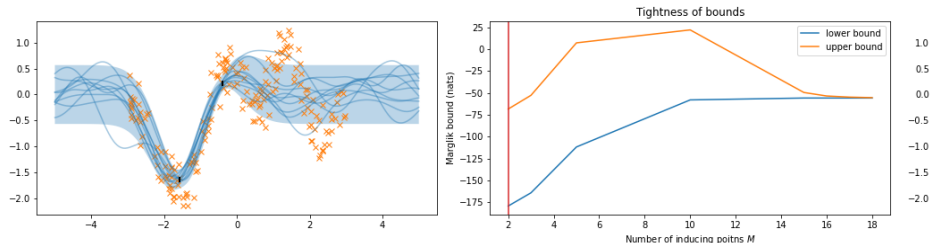
The ELBO helps us select **every** free parameter of the method!

- ▶ Q: How to select hyperparameters?  
A: Maximise  $\mathcal{L}$ . No overfitting, since it's a lower bound.
- ▶ Q: How to select inducing inputs  $\mathbf{Z}$ ?  
A: Maximise  $\mathcal{L}$  only reduces KL to true posterior.
- ▶ Q: When do we have enough inducing points?  
A: Once  $\mathcal{L}$  stops increasing (we also have upper bound).

---

<sup>4</sup>Titsias (2009)

# Demo



We jointly optimise  $\mathcal{L}$  w.r.t. its two free parameters:  $\mathbf{Z}, \boldsymbol{\theta}$ .

- ▶ Approximation and fit are poor when  $M$  is too small.
- ▶ ELBO converges with  $M \ll N$ .
- ▶ Upper bound<sup>5</sup> converges later to confirm good quality.

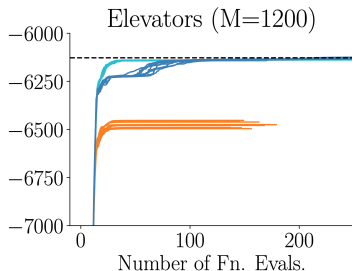
<sup>5</sup>See Titsias (2014), or Burt et al. (2020) for a discussion.

# Theory Gives Solutions

In “Convergence of Sparse Variational Inference in Gaussian Processes Regression” (Burt et al., 2020) we

- ▶ discussed a gradient-free inducing point initialisation scheme
- ▶ proved that it would give arbitrarily accurate results as  $N \rightarrow \infty$
- ▶ proved that the asymptotic complexity was reasonable  $O(N(\log N)^{2D}(\log \log N)^2)$  for SqExp, barely above linear<sup>6</sup>

Practical implications for the Titsias (2009) method:



<sup>6</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.

# Sparse GPs: Summary

Recipe for Sparse Variational GP Regression<sup>7</sup>:

1. Select initial number of inducing points  $M$  to try.
2. Select  $\mathbf{Z}$  with the greedy variance method (Burt et al., 2020).
3. Optional: Optimise  $\mathcal{L}$  w.r.t.  $\boldsymbol{\theta}$ .
4. Stop if upper-lower gap is small, or if improvement in  $\mathcal{L}$  is small. Otherwise repeat from step 2.

---

<sup>7</sup>as recommended in Burt, Rasmussen, and van der Wilk (2020)

# Sparse GPs: Conclusion

Sparse Gaussian Process approximations provide a **unified** way to approximate GPs:

- ▶ Correct and consistent posterior approximations.
- ▶ Single objective function can be used for setting all parameters.
- ▶ Measurable quality of approximation.
- ▶ For certain kernels, guarantees of good and cheap approximation as  $N \rightarrow \infty$  (+conditions).
- ▶ Burt et al. (2020) hints at a link between generalisation and approximation sparsity/quality. This would be very interesting to investigate in the context of infinite NN kernels.

# Modelling Assumptions

Goal: Learn some mapping  $f : \mathcal{X} \rightarrow \mathcal{Y}$ .

Assumptions about  $f$  influence generalisation performance:

- ▶ Fully connected vs convolutional?
- ▶ How smooth is the function?
- ▶ Data augmentation? I.e. what transformations leave the label unchanged?

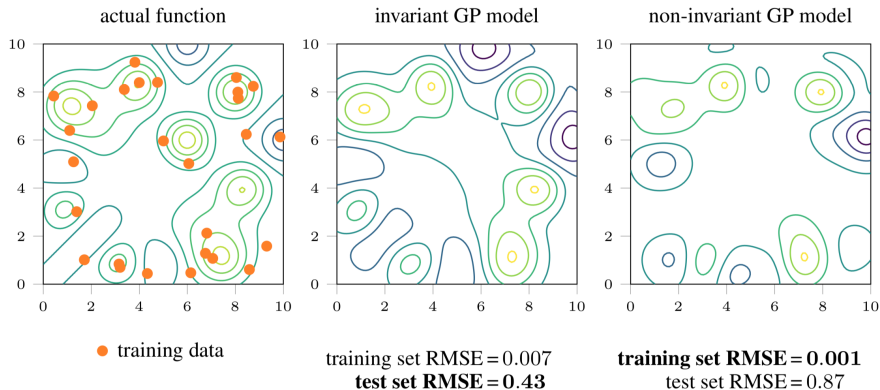
Central question:

How can changes on the input affect the output?

The fewer unnecessary degrees of variation we have, the better we will generalise.

▶▶ Goal: Find the right degrees of freedom as well as learning  $f$

# Example: Symmetry



- Learn symmetric function
- Pick either symmetrically constrained model or flexible model
- Symmetric model generalises better

# Data Augmentation and Invariances

Data augmentations express the **knowledge** about  $f(\cdot)$  that the output doesn't change in response to changes in the input. This is **invariance**.

We can consider strict invariances:

$$f(\mathbf{x}) = f(t(\mathbf{x})) \quad \forall \mathbf{x} \in \mathcal{X} \quad \forall t \in \mathcal{T} \quad (9)$$

or softer invariance:

$$P\left(\left(f(\mathbf{x}) - f(t(\mathbf{x}))\right)^2 > L\right) < \epsilon \quad \forall \mathbf{x} \in \mathcal{X} \quad t \sim p(t) \quad (10)$$



# Data Augmentation and Invariance

## Questions:

- ▶ How should we incorporate knowledge of invariances/DA into our models? (Particularly in the Bayesian context!)
- ▶ How can we select the right invariance/DA if we do not know it a priori?

In “Learning Invariances using the Marginal Likelihood” (v.d.Wilk et al., 2018) we

- ▶ Provide a clear formulation of how this can be done in a Bayesian context.<sup>8</sup>
- ▶ Provide a practical procedure for learning invariances using gradient descent in GPs.

---

<sup>8</sup><https://statmodeling.stat.columbia.edu/2019/12/02/a-bayesian-view-of-data-augmentation/>

# Model Selection according to Bayes

Model selection from a Bayesian point of view:

$$\begin{aligned} p(f, \boldsymbol{\theta} | \mathbf{y}) &= \frac{p(\mathbf{y} | f) p(f | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y})} \\ &= \underbrace{\frac{p(\mathbf{y} | f) p(f | \boldsymbol{\theta})}{p(\mathbf{y} | \boldsymbol{\theta})}}_{p(f | \mathbf{y}, \boldsymbol{\theta})} \underbrace{\frac{p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y})}}_{p(\boldsymbol{\theta} | \mathbf{y})} \end{aligned}$$

Key quantity for model selection is the [marginal likelihood](#)

$$p(\mathbf{y} | \boldsymbol{\theta}) = \int p(\mathbf{y} | f) p(f | \boldsymbol{\theta}) d\boldsymbol{\theta}$$

By handling our uncertainty on  $f(\cdot)$  in a Bayesian way, we also get the marginal likelihood for model selection.

# Model Selection: Procedure

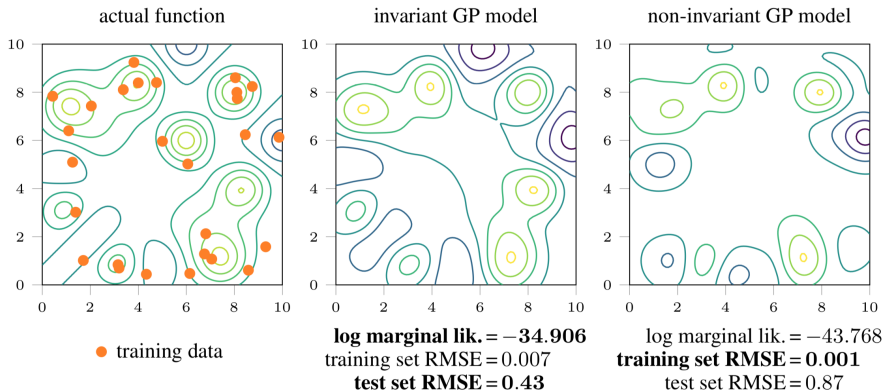
Our desired simplified procedure:

- ▶ Place prior on  $f$  with invariances described by  $\theta$ .
- ▶ Find posterior over functions  $p(f | \mathbf{y}, \theta)$ .
- ▶ Perform Maximum Likelihood on  $p(\mathbf{y} | \theta)$ .

This is **more safe** from over-fitting that performing Maximum Likelihood on  $f, \theta$  together.

(If you're sceptical, ask me for an example at the end.)

# Marginal Likelihood



$$\begin{aligned}\log p(\mathbf{y} | \boldsymbol{\theta}) &= \log p(y_1 | \boldsymbol{\theta}) + \log p(y_2 | y_1, \boldsymbol{\theta}) + \log p(y_3 | \{y_i\}_{i=1}^2, \boldsymbol{\theta}) \dots \\ &= \sum_{n=1}^N \log p(y_n | \{y_i\}_{i=1}^{n-1}, \boldsymbol{\theta})\end{aligned}$$

# Practicalities

Procedure so far is completely general and abstract. We need to:

- ▶ Choose our model class through the prior  $p(f | \theta)$
- ▶ Parameterise invariances in the prior through  $\theta$
- ▶ Show how to calculate  $p(f | \mathbf{y}, \theta)$  and  $p(\mathbf{y} | \theta)$

We choose

- ▶ Gaussian process priors  $p(f | \theta)$
- ▶ A construction of invariant GPs following Kondor (2008) and Ginsbourger et al. (2012)
- ▶ Variational approximation for posterior and marginal likelihood (Titsias 2009; Hensman et al. 2013)

# Invariant Gaussian Processes

Easy to place Gaussian process priors on non-invariant functions:

$$g(\cdot) \sim \mathcal{GP}(0, k_g(\cdot, \cdot')), \quad g : \mathbb{R}^D \rightarrow \mathbb{R}, \quad k_g : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}. \quad (11)$$

Can construct an invariant  $f(\cdot)$  by summing over the **orbit** of the group of transformations we want to be invariant to.

$$f(\cdot) = \sum_{\mathbf{x}_a \in \mathcal{O}(\mathbf{x})} g(\mathbf{x}_a) \quad (12)$$

Example:  $\mathcal{T}$  group of all rotations

► Orbit of image is set of images rotated by all angles

# Insensitivity

Parameterising orbits is hard, so we relax strict invariance constraint, and sum over arbitrary sets  $\mathcal{A}(\mathbf{x})$

$$f(\cdot) = \sum_{\mathbf{x}_a \in \mathcal{A}(\mathbf{x})} g(\mathbf{x}_a) \quad (13)$$

Parameterising sets is also cumbersome, so we take the infinite limit, to get an expectation

$$f(\cdot) = \int g(\mathbf{x}_a) p(\mathbf{x}_a | \mathbf{x}) d\mathbf{x}_a \quad (14)$$

- ▶ No longer strictly invariant
- ▶ Instead (roughly) a limit on  $P((f(\mathbf{x}_a) - f(\mathbf{x}))^2 > L)$
- ▶ Can interpolate between non-invariant and invariant

# Insensitive Kernels

The summation construction implies a kernel over  $f(\cdot)$ :

$$g(\cdot) \sim \mathcal{GP}(0, k_g(\cdot, \cdot'))$$

$$f(\cdot) \sim \mathcal{GP}(0, k_f(\cdot, \cdot')) \quad (\text{by linearity})$$

$$\begin{aligned} k_f(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_g[f(\mathbf{x})f(\mathbf{x}')] \\ &= \mathbb{E}_g \left[ \left( \int g(\mathbf{x}_a) p(\mathbf{x}_a | \mathbf{x}) d\mathbf{x}_a \right) \left( \int g(\mathbf{x}'_a) p(\mathbf{x}'_a | \mathbf{x}') d\mathbf{x}'_a \right) \right] \\ &= \iint \mathbb{E}_g[g(\mathbf{x}_a)g(\mathbf{x}'_a)] p(\mathbf{x}_a | \mathbf{x}) p(\mathbf{x}'_a | \mathbf{x}') d\mathbf{x}_a d\mathbf{x}'_a \\ &= \iint k_g(\mathbf{x}_a, \mathbf{x}'_a) p(\mathbf{x}_a | \mathbf{x}) p(\mathbf{x}'_a | \mathbf{x}') d\mathbf{x}_a d\mathbf{x}'_a \end{aligned}$$

Parameterise insensitivity by parameterising  $p_{\theta}(\mathbf{x}_a | \mathbf{x})$ !



# Interpolating to strict invariance

$$f(\mathbf{x}) \quad g(\mathbf{x}_a), p_{\theta}(\mathbf{x}_a | \mathbf{x})$$

# Overview

We have introduced GP priors with invariance properties controlled by  $p_{\boldsymbol{\theta}}(\mathbf{x}_a | \mathbf{x})$ .

Now we must compute

- ▶ The marginal likelihood  $p(\boldsymbol{\theta} | \mathbf{y})$  and its gradients for selecting the invariance
- ▶ The posterior  $p(f | \mathbf{y}, \boldsymbol{\theta})$  to make predictions

# Computational difficulties

Approximations are necessary in Gaussian process models for the well-known reasons:

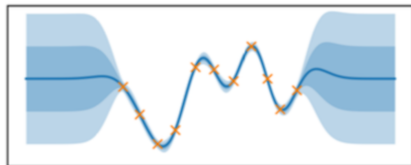
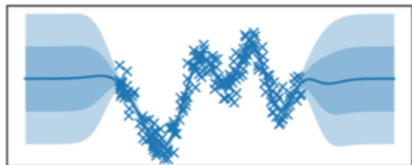
- ▶ Kernel inversions cost  $O(N^3)$
- ▶ Non-conjugate likelihoods (classification)
- ▶ No minibatch training

Here we have an additional problem:

We can't even evaluate the kernel!

# Variational inference

Approximate posterior: Prior conditioned on  $M \ll N$  noisy observations  $f(Z) = \{f(\mathbf{z}_m)\}_{m=1}^M$ :



$$q(f(\mathbf{x})) = \mathcal{N}\left(f(\mathbf{x}); \mathbf{k}_{\mathbf{x}Z} \mathbf{K}_{ZZ}^{-1} \mathbf{m}, k_f(\mathbf{x}, \mathbf{x}) - \mathbf{k}_{\mathbf{x}Z} \mathbf{K}_{ZZ}^{-1} (\mathbf{K}_{ZZ} - \mathbf{S}) \mathbf{K}_{ZZ}^{-1} \mathbf{k}_{Z\mathbf{x}}\right)$$

$$\mathcal{L} = \sum_{n=1}^N \mathbb{E}_{q(f(\mathbf{x}_n))} [\log p(\mathbf{y}_n | f(\mathbf{x}_n))] - \text{KL}[q(f(Z)) || p(f(Z))]$$

Gives: **Approximate posterior** , **lower bound to marginal likelihood**

Solves: **inversion cost** , **non-conjugate likelihoods** , **minibatching**

# Variational Inference

For Gaussian likelihoods

$$\mathbb{E}_{q(f(\mathbf{x}_n))}[\log p(\mathbf{y}_n | f(\mathbf{x}_n))] = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y_n - \mu_n)^2 - \frac{\sigma_n^2}{2\sigma^2}$$

$$q(f(\mathbf{x})) = \mathcal{N}(f(\mathbf{x}); \underbrace{\mathbf{k}_{\mathbf{x}\mathbf{Z}} \mathbf{K}_{\mathbf{ZZ}}^{-1} \mathbf{m}}_{\mu_n}, \underbrace{k_f(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{\mathbf{x}\mathbf{Z}} \mathbf{K}_{\mathbf{ZZ}}^{-1} (\mathbf{K}_{\mathbf{ZZ}} - \mathbf{S}) \mathbf{K}_{\mathbf{ZZ}}^{-1} \mathbf{k}_{\mathbf{Z}\mathbf{x}}}_{\sigma_n^2})$$

With

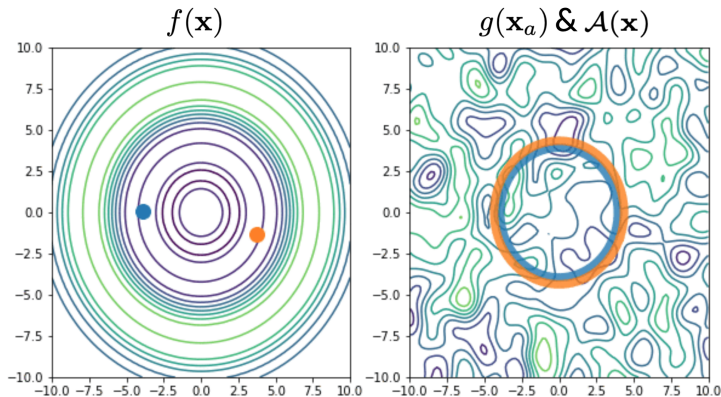
$$[\mathbf{k}_{\mathbf{Z}\mathbf{x}}]_{mn} = \iint k_g(\mathbf{x}_a, \mathbf{x}'_a) p(\mathbf{x}_a | \mathbf{z}_m) p(\mathbf{x}'_a | \mathbf{z}'_m) d\mathbf{x}_a d\mathbf{x}'_a$$

$$[\mathbf{K}_{\mathbf{ZZ}}]_{mm} = \iint k_g(\mathbf{x}_a, \mathbf{x}'_a) p(\mathbf{x}_a | \mathbf{z}_m) p(\mathbf{x}'_a | \mathbf{z}'_m) d\mathbf{x}_a d\mathbf{x}'_a$$

Monte Carlo estimates could help if we didn't have the inverses...

# Interdomain inducing variables

- ▶ The variational distribution is constructed by **conditioning** on “inducing observations”.
- ▶ **Which** random variables we condition on determines the covariances



# Interdomain Variational Inference

For Gaussian likelihoods

$$\mathbb{E}_{q(f(\mathbf{x}_n))}[\log p(\mathbf{y}_n | f(\mathbf{x}_n))] = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y_n - \mu_n)^2 - \frac{\sigma_n^2}{2\sigma^2}$$
$$q(f(\mathbf{x})) = \mathcal{N}(f(\mathbf{x}); \underbrace{\mathbf{k}_{\mathbf{x}\mathbf{Z}} \mathbf{K}_{\mathbf{ZZ}}^{-1} \mathbf{m}}_{\mu_n}, \underbrace{k_f(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{\mathbf{x}\mathbf{Z}} \mathbf{K}_{\mathbf{ZZ}}^{-1} (\mathbf{K}_{\mathbf{ZZ}} - \mathbf{S}) \mathbf{K}_{\mathbf{ZZ}}^{-1} \mathbf{k}_{\mathbf{Z}\mathbf{x}}}_{\sigma_n^2})$$

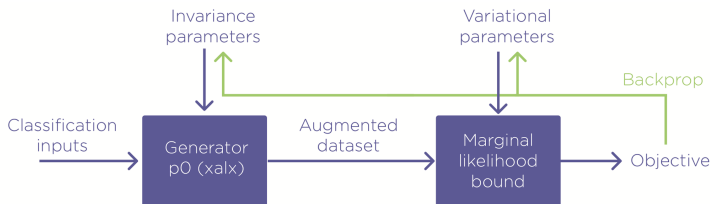
With

$$[\mathbf{k}_{\mathbf{Z}\mathbf{x}}]_{mn} = \int k_g(\mathbf{z}_m, \mathbf{x}_a) p(\mathbf{x}_a | \mathbf{x}) d\mathbf{x}_a$$
$$[\mathbf{K}_{\mathbf{ZZ}}]_{mm'} = k_g(\mathbf{z}_m, \mathbf{z}_{m'})$$

We can now find unbiased estimates of  $\mu_n$  and  $\sigma_n^2$ ! This trick also works with other likelihoods through the Pólya-Gamma trick!

# Procedure

- ▶ Compute ELBO (marginal likelihood lower bound)
- ▶ Back-propagate to **variational** and **invariance** parameters through re-parameterisation
- ▶ Optimise jointly



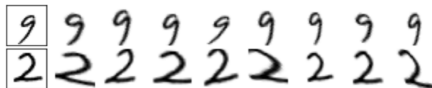
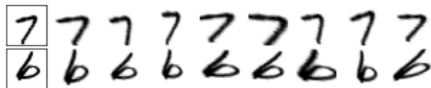
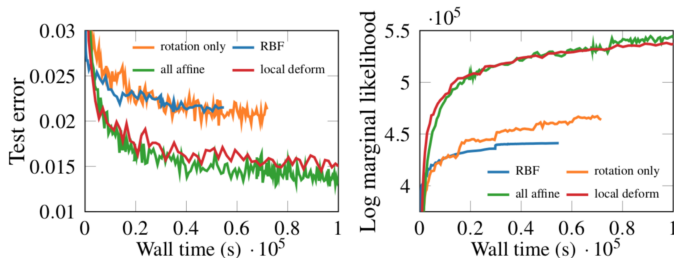
- ▶ No need for even a closed-form evaluation of the kernel  $k_f$
- ▶ Insensitivity distribution  $p(\mathbf{x}_a | \mathbf{x})$  can be **implicit**



# Results

We used various  $p(\mathbf{x}_a | \mathbf{x})$ :

- Affine transformations (parameters: how much rotation / skew / scale to apply)
- Local deformations (parameters: how much deformation, how much to smooth deformations etc)



# Conclusion

- ▶ Can express invariances in kernels (but kernels intractable)
- ▶ Can use the marginal likelihood for learning inductive biases
- ▶ We only need unbiased estimates of kernels to train!
- ▶ This is very much like learning the right data augmentation

Going forward:

- ▶ Embed into deep structures (e.g. deep GPs/NNs, see latest arxiv pre-prints)
- ▶ Could we use infinitely wide NN kernels?

# References

Key references. See paper for more.

- ▶ [Learning Invariances using the Marginal Likelihood](#); Mark van der Wilk, Matthias Bauer, ST John, James Hensman; NeurIPS (2018). (Main paper this was about)
- ▶ [Gaussian Processes for Big Data](#); James Hensman, Nicolo Fusi, James D. Hensman; UAI (2013). (Variational bound we use)
- ▶ [Variational Learning of Inducing Variables in Sparse Gaussian Processes](#); Michalis K. Titsias; AISTATS (2009). (First introduction of variational GP approx)
- ▶ [Argumentwise invariant kernels for the approximation of invariant functions](#); David Ginsbourger, Xavier Bay, Olivier Roustant, Laurent Carraro; Annales de la Faculté des Sciences de Toulouse (2012). (Invariant kernels)
- ▶ [Group theoretical methods in machine learning](#); Risi Kondor; PhD thesis (2008). (Invariant kernels)

# Minimising training loss

We're looking for a fit that will **generalise** to new unseen test data.  
Let's minimise the training loss of the posterior mean.

$$\mathcal{L}(\theta, \sigma) = \sum_{n=1}^N \left[ k_{\theta}(\mathbf{x}_n, X) \left( \mathbf{K}_{\theta} + \sigma^2 \mathbf{I} \right)^{-1} \mathbf{y} - y_n \right]^2 \quad (15)$$

$$\{\theta^*, \sigma^*\} = \underset{\theta, \sigma}{\operatorname{argmin}} \mathcal{L}(\theta, \sigma) \quad (16)$$

We can fit anything with a tiny lengthscale and noise variance!

# Marginal likelihood fixes things

Instead, choose hyperparameters by maximising marginal likelihood:

In above  $\mathcal{L}$  is indicated by 'datafit', while 'ELBO' indicates the marginal likelihood.

- ▶ More sensible fit as the marginal likelihood rises
- ▶ Datafit gets worse!

Marginal likelihood trades off  
data fit and model complexity.

# References I

- Artem Artemev, David R. Burt, and Mark van der Wilk. Tighter bounds on the log marginal likelihood of gaussian process regression using conjugate gradients. In [Proceedings of the 38th International Conference on Machine Learning \(ICML\)](#), 2021.
- David R. Burt, Carl Edward Rasmussen, and Mark van der Wilk. Convergence of sparse variational inference in gaussian processes regression. [Journal of Machine Learning Research](#), 21(131):1–63, 2020.
- Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional networks as shallow Gaussian processes. In [7th International Conference on Learning Representations \(ICLR\)](#), 2019.
- Mark Gibbs and David Mackay. Efficient implementation of Gaussian processes. Technical report, Cavendish Laboratory, University of Cambridge, 1997.
- James Hensman, Nicolò Fusi, and Neil D. Lawrence. Gaussian processes for big data. In [Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence \(UAI\)](#), pages 282–290, 2013.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as Gaussian processes. In [6th International Conference on Learning Representation \(ICLR\)](#), 2018.
- Alexander G. de G. Matthews. [Scalable Gaussian Process Inference Using Variational Methods](#). PhD thesis, University of Cambridge, Cambridge, UK, 2016. available at <http://mlg.eng.cam.ac.uk/matthews/thesis.pdf>.
- Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In [6th International Conference on Learning Representations \(ICLR\)](#), 2018.
- Radford M. Neal. [Bayesian learning for neural networks](#), volume 118. Springer, 1996.
- Roman Novak, Lechao Xiao, Yasaman Bahri, Jaehoon Lee, Greg Yang, Jiri Hron, Daniel A. Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are Gaussian processes. In [7th International Conference on Learning Representations \(ICLR\)](#), 2019.

# References II

- Alex J. Smola and Bernhard Schölkopf. Sparse greedy matrix approximation for machine learning. In [17th International Conference on Machine Learning, Stanford, 2000](#), 2000.
- Michalis K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In [Proceedings of the 12th International Conference on Artificial Intelligence and Statistics](#), pages 567–574, 2009.
- Michalis K. Titsias. Variational inference for Gaussian and determinantal point processes (talk slides). In [Workshop on Advances in Variational Inference \(NIPS\)](#), 2014.
- Mark v.d.Wilk, Matthias Bauer, ST John, and James Hensman. Learning invariances using the marginal likelihood. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, [Advances in Neural Information Processing Systems 31](#), pages 9938–9948. Curran Associates, Inc., 2018. URL <http://papers.nips.cc/paper/8199-learning-invariances-using-the-marginal-likelihood.pdf>.
- Ke Wang, Geoff Pleiss, Jacob Gardner, Stephen Tyree, Kilian Q Weinberger, and Andrew Gordon Wilson. Exact gaussian processes on a million data points. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, [Advances in Neural Information Processing Systems](#), volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/01ce84968c6969bdd5d51c5eeaa3946a-Paper.pdf>.
- Christopher Williams and Matthias Seeger. Using the Nyström method to speed up kernel machines. In [Advances in Neural Information Processing Systems 13](#), pages 682–688, 2001.
- Greg Yang. Wide feedforward or recurrent neural networks of any architecture are Gaussian processes. In [Advances in Neural Information Processing Systems 32 \(NeurIPS\)](#). 2019.