



# Fitting difficult GP models

James Hensman

PROWLER.io

October 2017



# Overview

Motivation

Variational Bayes

Variational Bayes with processes

Bonus content: coping with image inputs

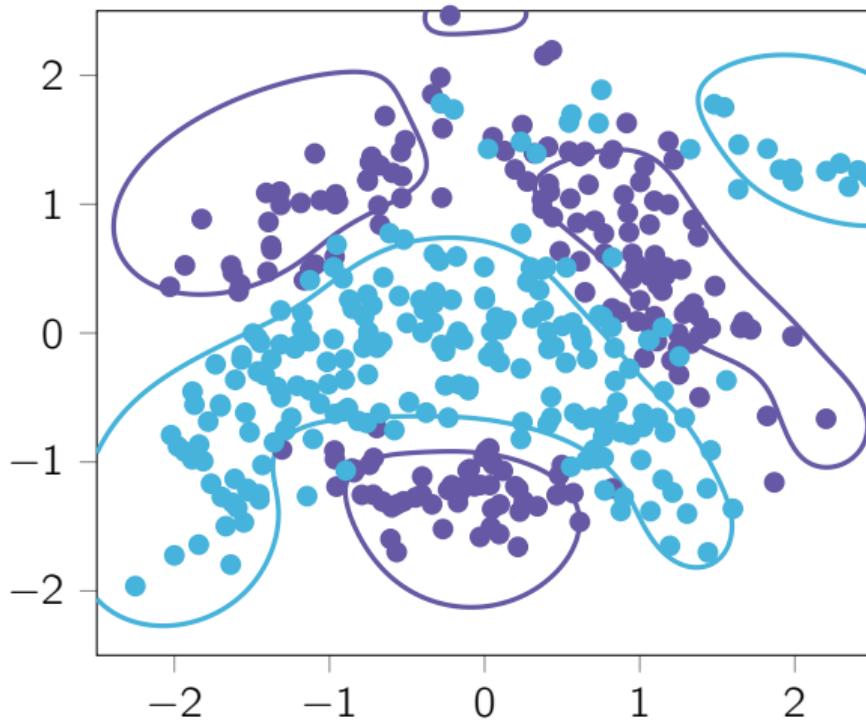




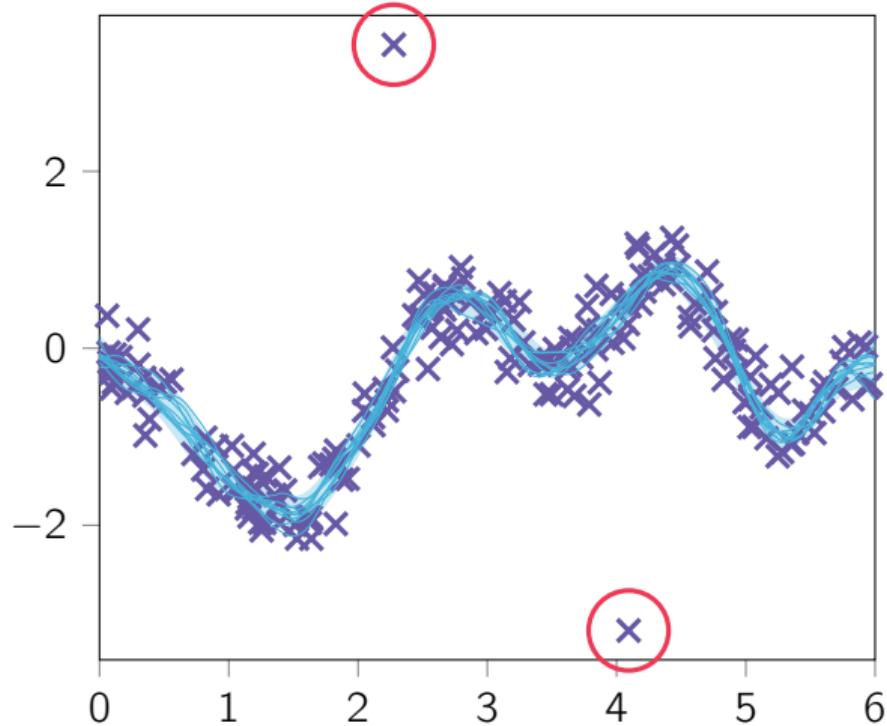
# Motivation

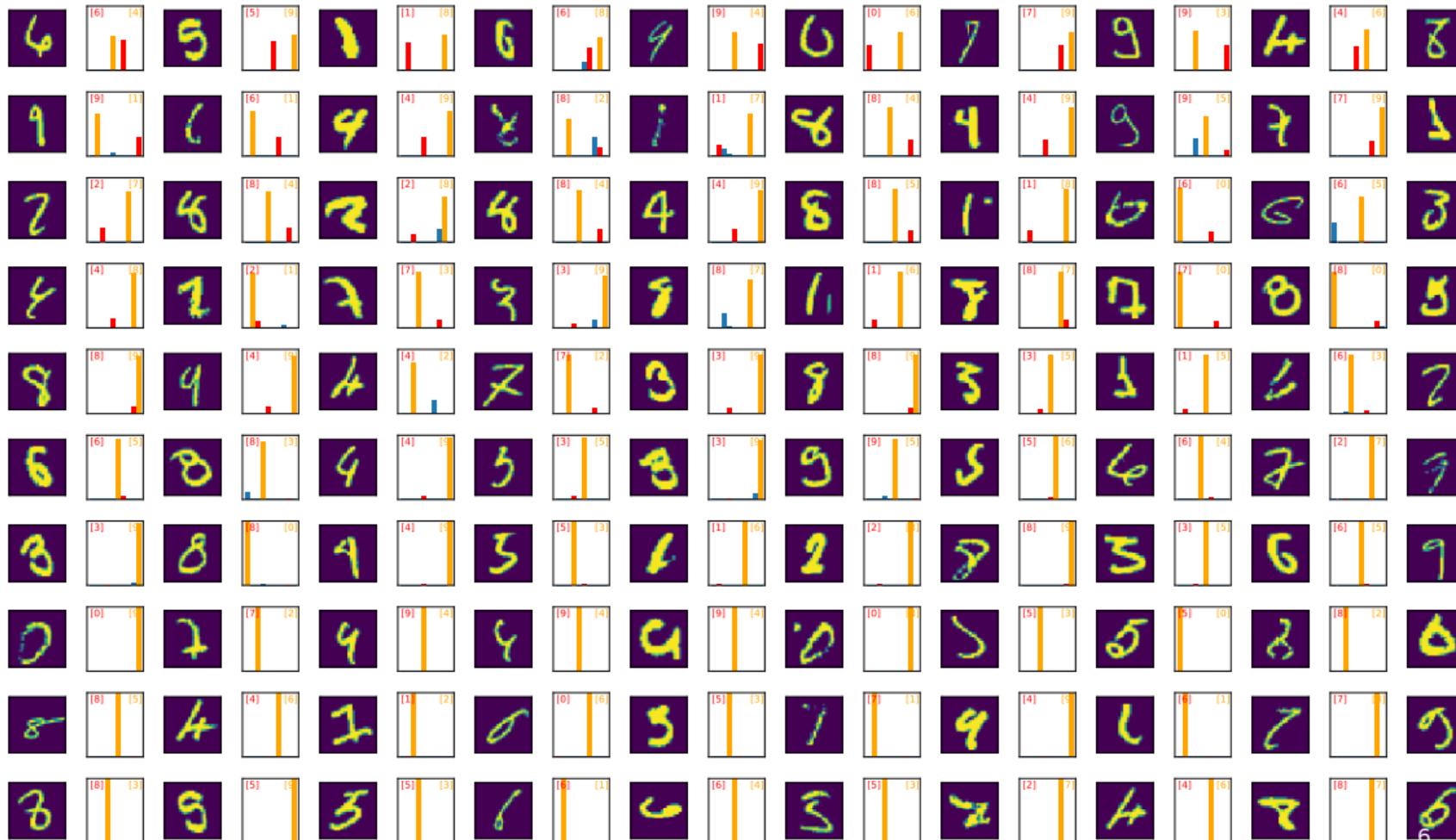


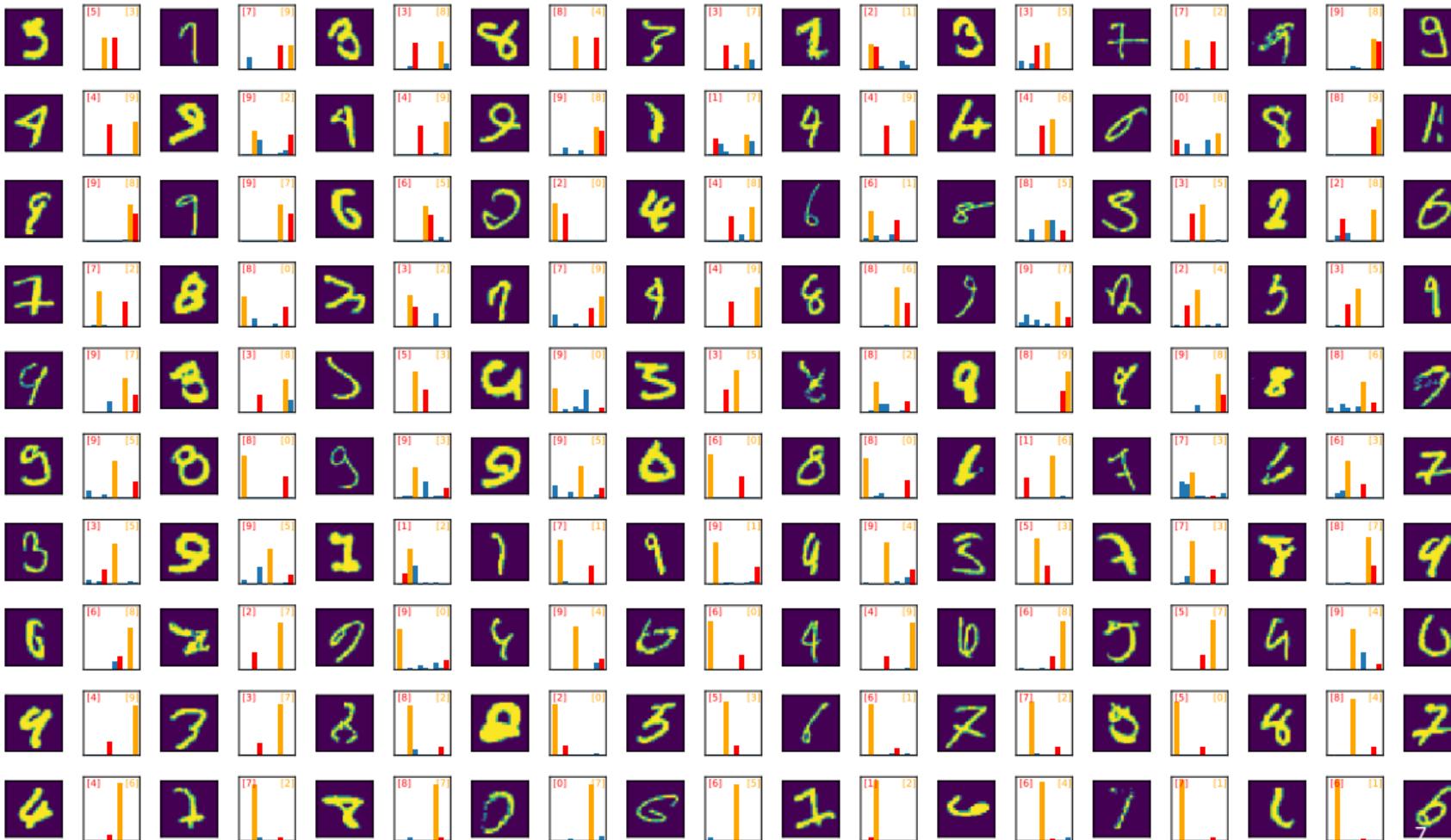
# Non-Gaussian Gaussian processes



# Non-Gaussian Gaussian processes

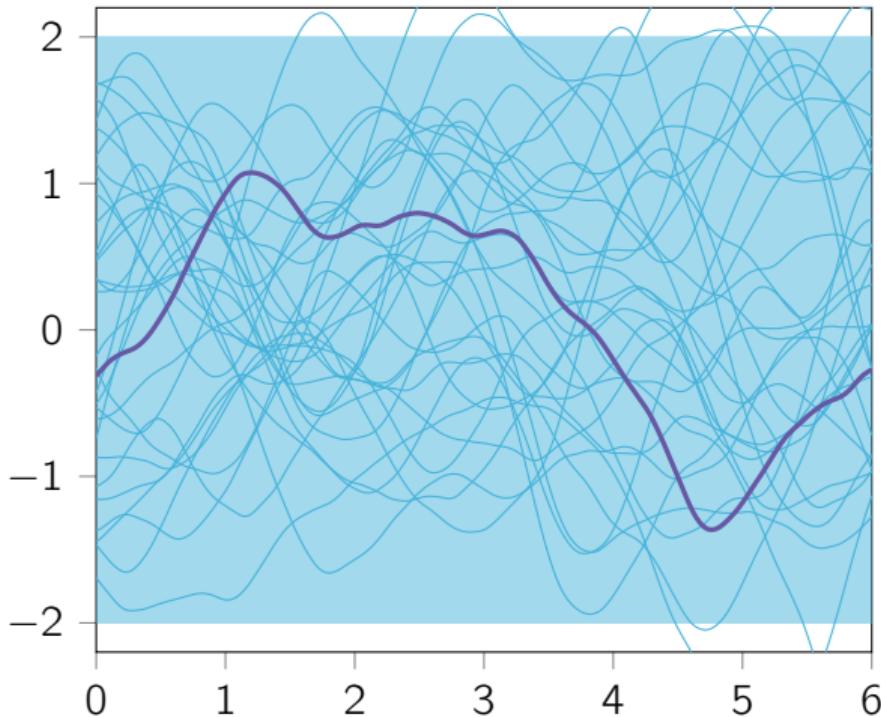


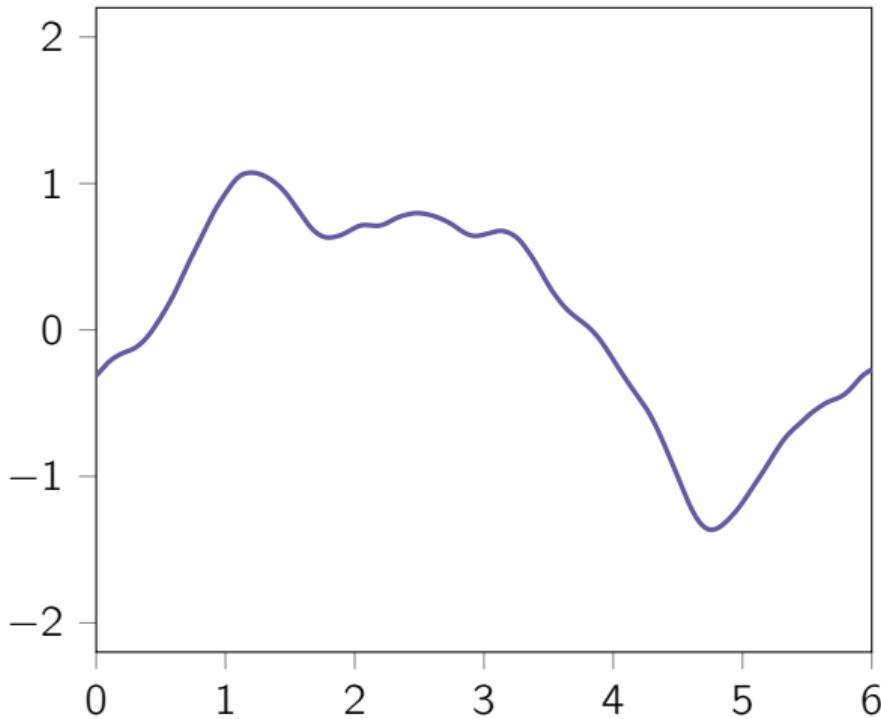


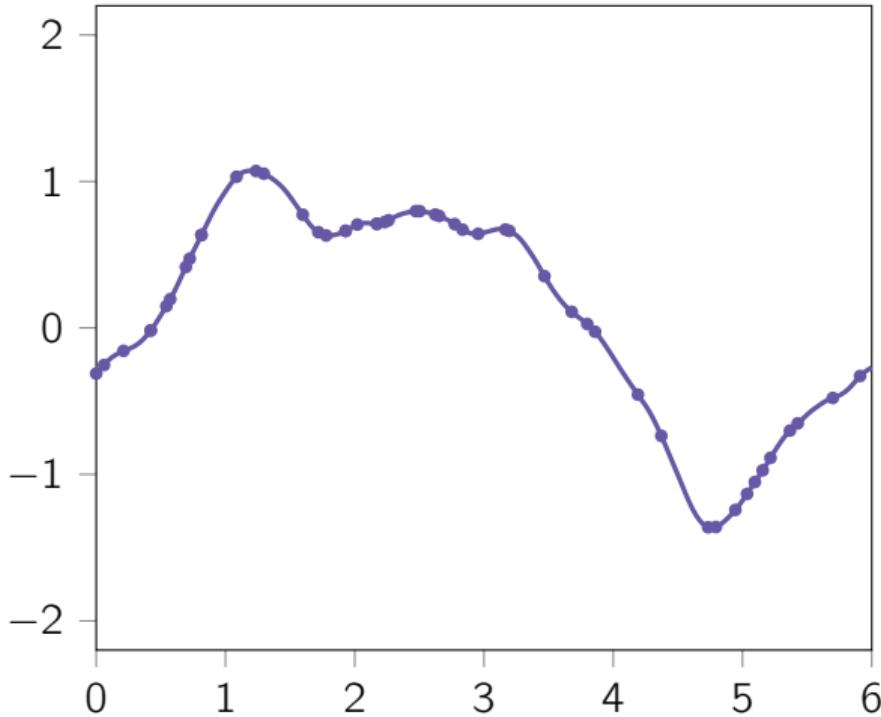


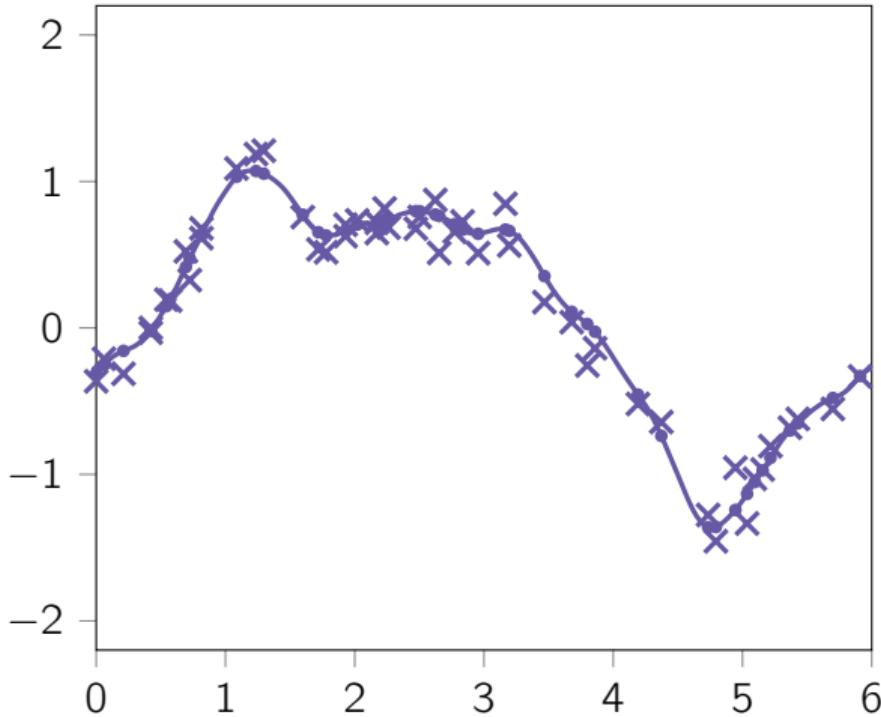
How are data generated in Gaussian process regression?

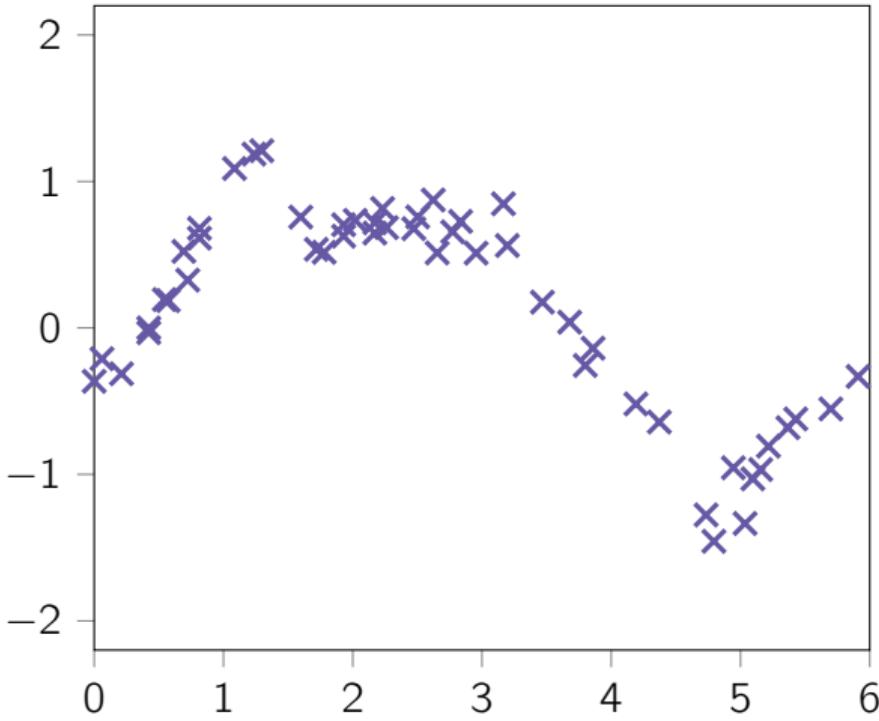


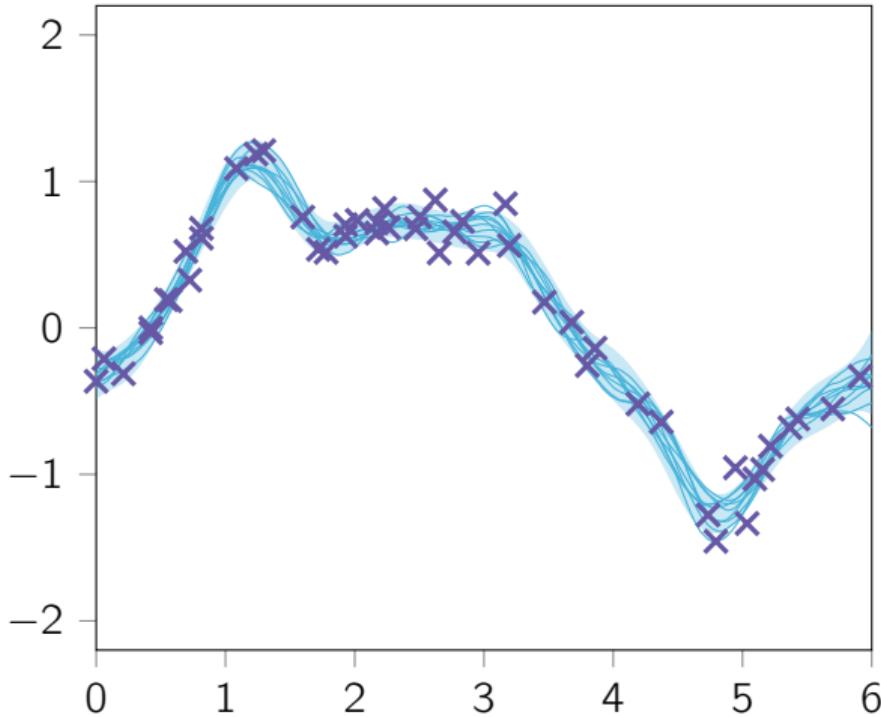






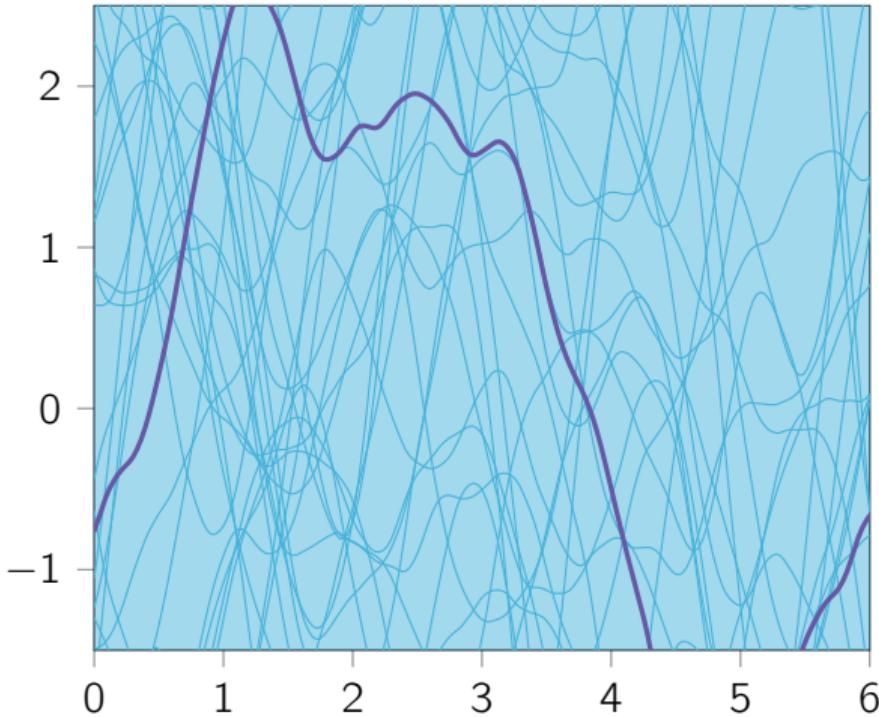


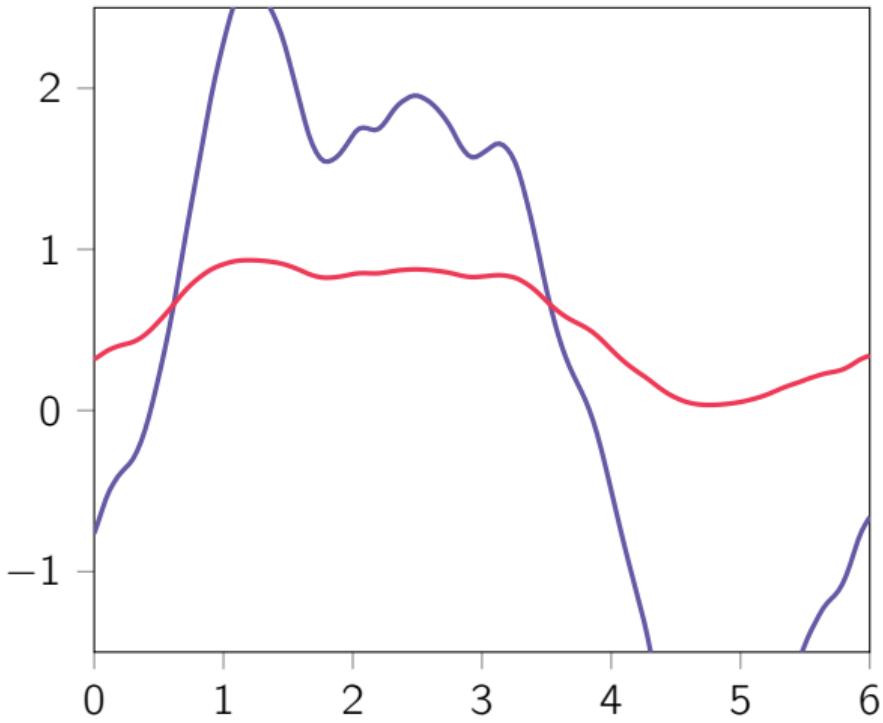


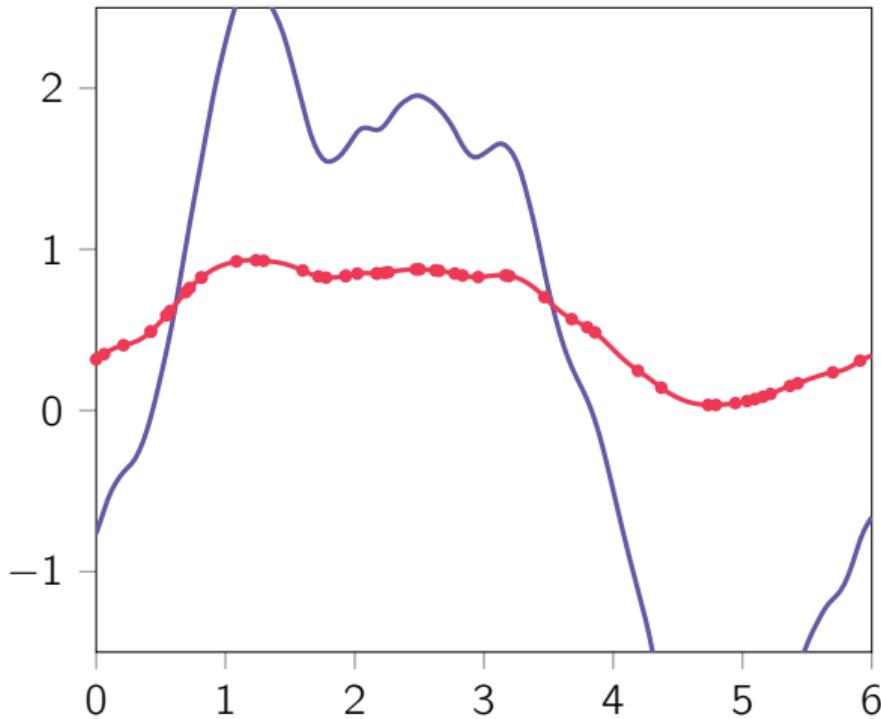


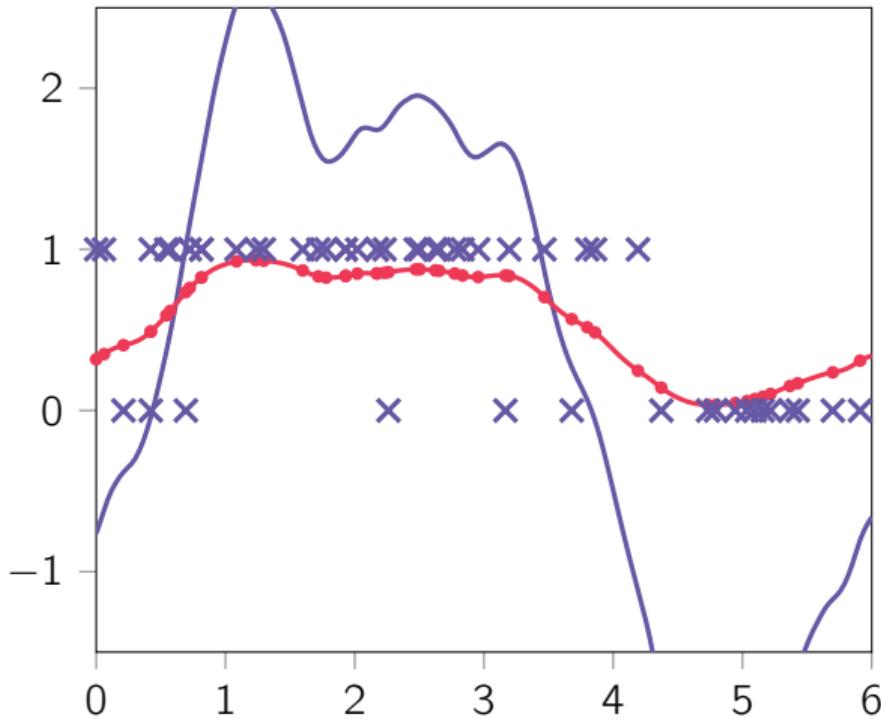
How are data generated in Gausian process **classification**?

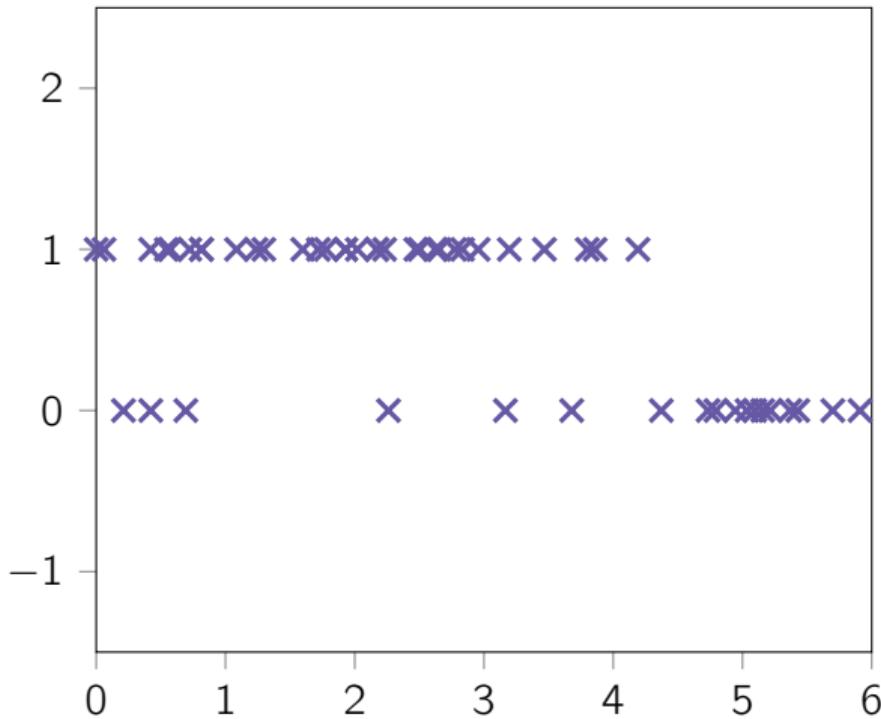


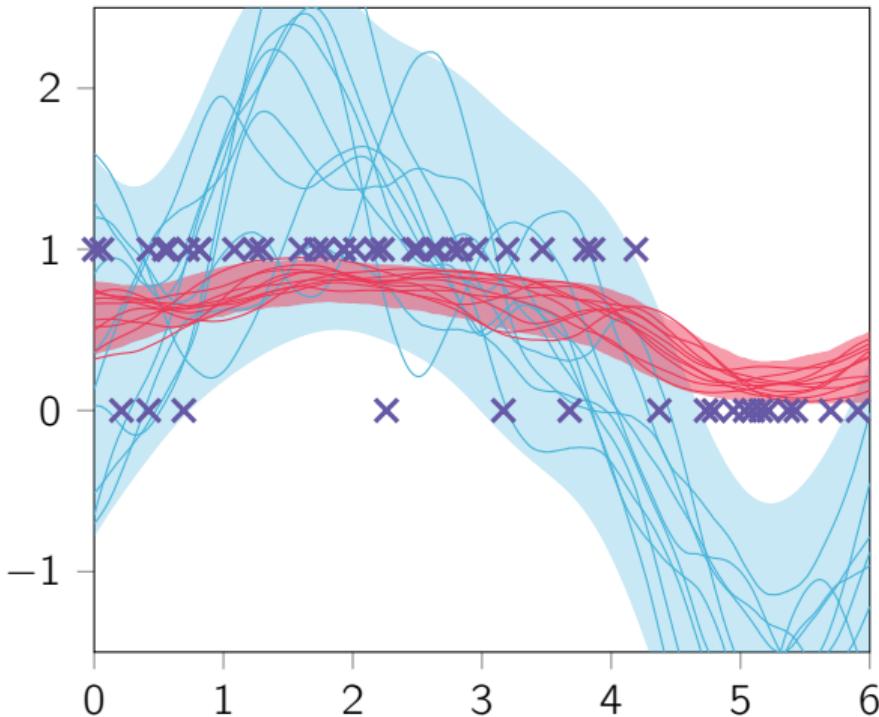












# Big Gaussian processes

A long-standing problem with GPs has been that they scale cubically with the number of data.

Recall

$$p(f(\cdot)) = \underbrace{p(f(\cdot) | \mathbf{f})}_{\text{GP prior}} \underbrace{p(\mathbf{f})}_{\text{GP conditional}} \underbrace{p(\mathbf{f})}_{\text{Gaussian dist}}$$





# Variational Bayes

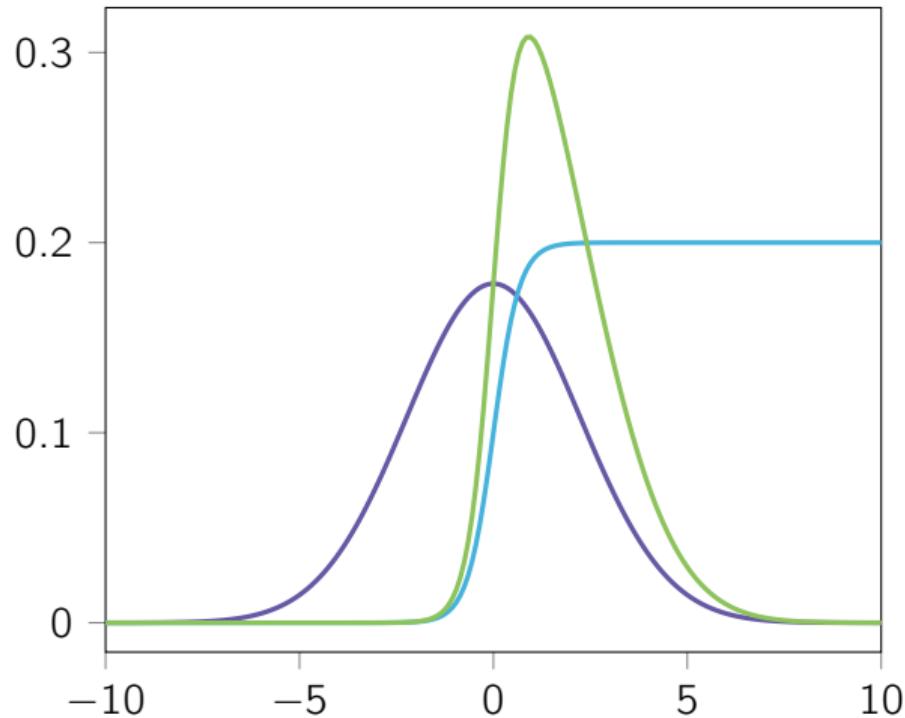


# Variational Bayes

- ▶ In Gaussian process models, we want to do Bayesian inference on functions
- ▶ We're going to have to make approximations!
- ▶ VB is our method of choice for approximation.



$$p(x | y) = p(y | x)p(x)/p(y)$$



# Worked example

$$p(x) = \mathcal{N}(x | 0, \sigma_p^2)$$

$$p(y | x) = \phi(x)$$

$$p(x | y) = ?$$

$$p(y) = ?$$

$$q(x) = \mathcal{N}(x | m, \sigma_q^2)$$



$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log p(x \mid y)] dx$$



$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log p(x \mid y)] dx$$

$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log \frac{p(y \mid x)p(x)}{p(y)}] dx$$



$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log p(x \mid y)] dx$$

$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log \frac{p(y \mid x)p(x)}{p(y)}] dx$$

$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log p(y \mid x) - \log p(x)] dx + \log p(y)$$



$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log p(x \mid y)] dx$$

$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log \frac{p(y \mid x)p(x)}{p(y)}] dx$$

$$\text{KL}[q(x) \parallel p(x \mid y)] = \int q(x) [\log q(x) - \log p(y \mid x) - \log p(x)] dx + \log p(y)$$

$$\text{KL}[q(x) \parallel p(x \mid y)] = -\text{ELBO}(q(x)) + p(y)$$







# Variational Bayes with processes



## Not so different from our toy example...

$$p(x) = \mathcal{N}(x | 0, \sigma_p^2) \rightarrow p(f(\cdot)) = \mathcal{GP}(0, k(\cdot, \cdot))$$

$$p(y | x) = \phi(x) \rightarrow p(\mathbf{y} | f(\cdot)) = \prod_{i=1}^n p(y_i | f(x_i))$$

$$p(x | y) = ? \rightarrow p(f(\cdot) | \mathbf{y}) = ?$$

$$p(y) = ? \rightarrow p(\mathbf{y}) = ?$$

$$q(x) = \mathcal{N}(x | m, \sigma_q^2) \rightarrow q_{\text{m,L,Z}}(f(\cdot)) = \mathcal{GP}(\mu(\cdot), v(\cdot, \cdot))$$



# The form of $q(f(\cdot))$

Recall the GP conditional equations:

$$p(f(\cdot)) = p(f(\cdot) | \mathbf{f})p(\mathbf{f})$$

with

$$p(f(\cdot) | \mathbf{f}) = \mathcal{GP}\left(k(\cdot, X)\mathbf{K}^{-1}\mathbf{f}, k(\cdot, \cdot) - k(\cdot, X)\mathbf{K}^{-1}k(X, \cdot)\right)$$

$$p(\mathbf{f}) = \mathcal{N}(0, \mathbf{K})$$

where

$$\mathbf{f} = [f(x_i)]_{i=1}^N.$$



# The form of $q(f(\cdot))$

Now choose some points  $\mathbf{Z}$  so

$$p(f(\cdot)) = p(f(\cdot) | \mathbf{u})p(\mathbf{u})$$

with

$$\begin{aligned} p(f(\cdot) | \mathbf{u}) &= \mathcal{GP}\left(k(\cdot, \mathbf{Z})\mathbf{K}^{-1}\mathbf{u}, k(\cdot, \cdot) - k(\cdot, \mathbf{Z})\mathbf{K}^{-1}k(\mathbf{Z}, \cdot)\right) \\ p(\mathbf{u}) &= \mathcal{N}(0, \mathbf{K}) \end{aligned}$$

where

$$\mathbf{u} = [f(x_i)]_{i=1}^M.$$



Now let

$$q(f(\cdot) | \mathbf{u}) = p(f(\cdot) | \mathbf{u})$$

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{L}\mathbf{L}^\top)$$



Now let

$$q(f(\cdot) | \mathbf{u}) = p(f(\cdot) | \mathbf{u})$$

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{L}\mathbf{L}^\top)$$

so that

$$q(f(\cdot)) = \mathcal{GP}\left(k(\cdot, \mathbf{Z})\mathbf{K}^{-1}\mathbf{m}, k(\cdot, \cdot) - k(\cdot, \mathbf{Z})\mathbf{K}^{-1}[\mathbf{K} - \mathbf{L}\mathbf{L}^\top]\mathbf{K}^{-1}k(\mathbf{Z}, \cdot)\right)$$



Now let

$$q(f(\cdot) | \mathbf{u}) = p(f(\cdot) | \mathbf{u})$$

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{L}\mathbf{L}^\top)$$

so that

$$q(f(\cdot)) = \mathcal{GP}\left(k(\cdot, \mathbf{Z})\mathbf{K}^{-1}\mathbf{m}, k(\cdot, \cdot) - k(\cdot, \mathbf{Z})\mathbf{K}^{-1}[\mathbf{K} - \mathbf{L}\mathbf{L}^\top]\mathbf{K}^{-1}k(\mathbf{Z}, \cdot)\right)$$

now: optimize  $\text{KL}[q(f(\cdot)) || p(f(\cdot) | \mathbf{Y})]$  w.r.t.  $\mathbf{Z}, \mathbf{m}, \mathbf{L}$ .



# What happens when we adjust **m**?



# What happens when we adjust **L**?



# What happens when we adjust **Z**?





# The ELBO for GP models

$$\text{ELBO} = \int q(f(\cdot)) \left[ \log p(\mathbf{y} | f(\cdot)) + \log \frac{p(f(\cdot))}{q(f(\cdot))} \right] df(\cdot)$$

- ▶ We can compute the term  $p(f(\cdot))/q(f(\cdot))$  because of our careful construction of  $q$ .
- ▶ We can compute the integral of  $p(\mathbf{y} | f(\cdot))$  since we only need the marginals of  $q(f(\cdot))$ .
- ▶ We can minibatch if needed!



# Questions

?



- ▶ What about parameters of the kernel?
- ▶ How many **Z** should I pick?
- ▶ How reasonable is the Gaussian assumption?
- ▶ What if the likelihood does not factorize?
- ▶ What if the likelihood is Gaussian?





Bonus content: coping with image inputs



---

# Convolutional Gaussian Processes

---

**Mark van der Wilk**  
Department of Engineering  
University of Cambridge, UK  
mv310@cam.ac.uk

**Carl Edward Rasmussen**  
Department of Engineering  
University of Cambridge, UK  
cer54@cam.ac.uk

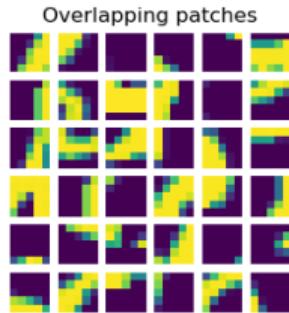
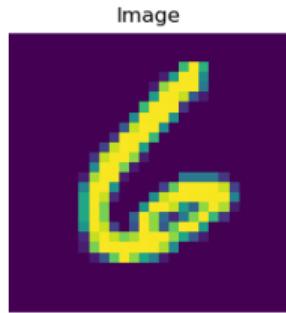
**James Hensman**  
prowler.io  
Cambridge, UK  
james@prowler.io

## Abstract

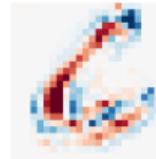
We present a practical way of introducing convolutional structure into Gaussian processes, making them more suited to high-dimensional inputs like images. The main contribution of our work is the construction of an inter-domain inducing point approximation that is well-tailored to the convolutional kernel. This allows us to gain the generalisation benefit of a convolutional kernel, together with fast but accurate posterior inference. We investigate several variations of the convolutional kernel, and apply it to MNIST and CIFAR-10, where we obtain significant improvements over existing Gaussian process models. We also show how the marginal likelihood can be used to find an optimal weighting between convolutional and RBF kernels to further improve performance. This illustration of the usefulness of the marginal likelihood may help automate discovering architectures in larger models.



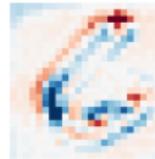
# Convolutional Gaussian Processes



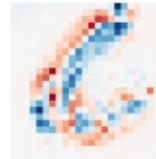
0,  $f(x) = 85.4$



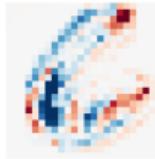
1,  $f(x) = -70.0$



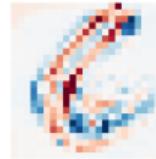
2,  $f(x) = -33.8$



$f(x) = -143.0$



4,  $f(x) = -1.4$



5,  $f(x) = -7.7$



6,  $f(x) = 206.6$



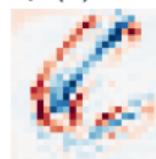
$$f(\mathbf{x}) = \sum_p w_p g\left(\mathbf{x}^{[p]}\right)$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_p \sum_{p'} k_g\left(\mathbf{x}^{[p]}, \mathbf{x}'^{[p']}\right)$$

$f(x) = -110.0$



8,  $f(x) = 8.7$



9,  $f(x) = -66.8$



$$k(\mathbf{x}, \mathbf{z}') = \sum_p k_g\left(\mathbf{x}^{[p]}, \mathbf{z}'\right)$$





# PROWLER.io

AI for Autonomous  
Decision Making

