

# VARIATIONAL INFERENCE

## in Gaussian process models

---

James Hensman

Approximate Inference workshop, NIPS 2015

Lancaster University

# COLLABORATORS



Alex Matthews  
Cambridge Univeristy



Nicolo Fusi  
Microsoft Research



Maurizio Filippone  
Eurecom



Rich Turner  
Cambridge Univeristy

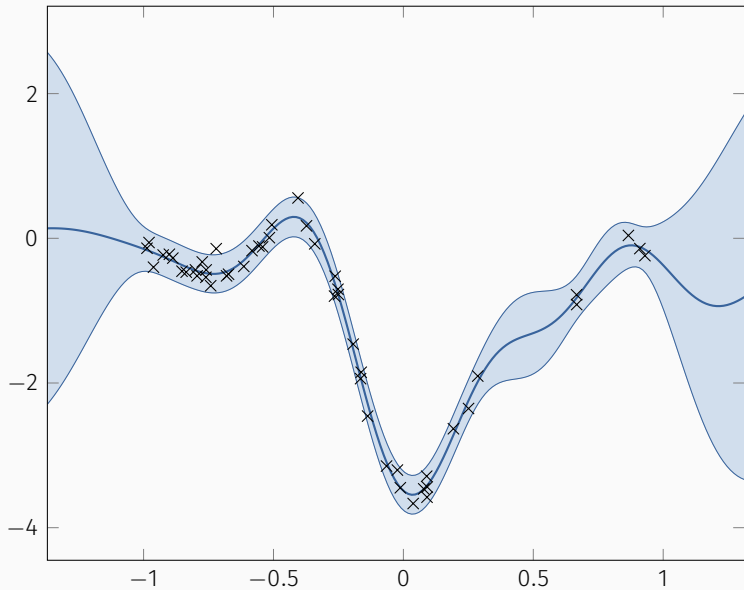


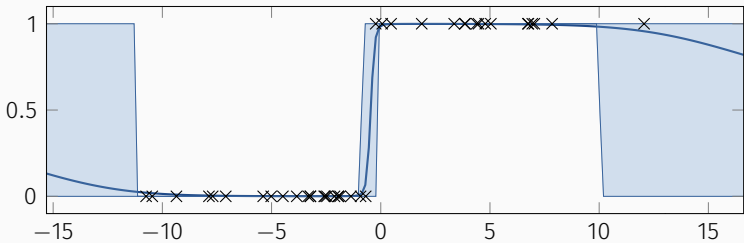
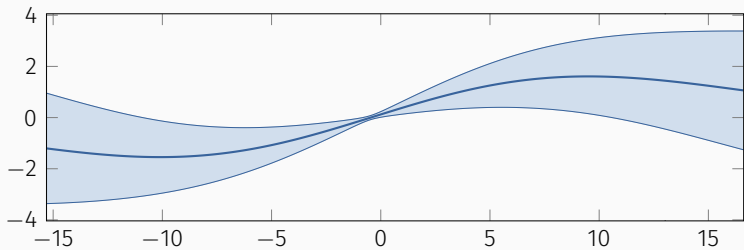
Neil D. Lawrence  
Sheffield Univeristy



Zoubin Ghahramani  
Cambridge Univeristy

# WHAT CAN GAUSSIAN PROCESSES DO?





## A unified view of variational GP approximations

- Deals with non-Gaussian posterior
- Deals with  $\mathcal{O}(n^3)$  complexity (sparse)
- The variational distribution contains a (conditionally) Gaussian process

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$$

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

with:

$$\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$$

$$y_i | f_i \sim \text{Po}(y_i | e^{f_i}) \quad \text{or} \quad \text{Bin}(y_i | \sigma(f_i)) \quad \text{or} \dots$$

- Local variational bounds (classification only) <sup>1</sup>
- Expectation Propagation <sup>2</sup>
- For classification, EP > VB <sup>3</sup>
- Variational methods need only 2N parameters <sup>4</sup>
- VB methods can be fast too! <sup>5</sup>
- VB can be applied to lots of different likelihoods <sup>6</sup>

---

<sup>1</sup>MN Gibbs, DJC MacKay - Variational Gaussian process classifiers - IEEE TNN 2000

<sup>2</sup>Minka, T. P. A family of algorithms for approximate Bayesian inference. Doctoral dissertation, MIT - 2001

<sup>3</sup>H Nickisch, CE Rasmussen - Approximations for binary Gaussian process classification - JMLR 2008

<sup>4</sup>M. Opper and C. Archambeau - The variational Gaussian approximation revisited - Neural comp. 2009

<sup>5</sup>E Khan, S Mohamed, KP Murphy - Fast Bayesian inference for non-conjugate Gaussian process regression- NIPS 2012

<sup>6</sup>Nguyen and Bonilla - Automated variational inference for Gaussian process models - NIPS 2014

- Subset-of-data methods<sup>7</sup> <sup>8</sup> hence ‘sparse’.
- Pseudo-inputs introduced <sup>9</sup>
- A unifying view brings several ideas together <sup>10</sup>
- Variational approach <sup>11</sup> makes for better placement of pseudo/inducing points
- Variational approach can be optimized with SVI <sup>12</sup>

---

<sup>7</sup>AJ Smola, P Bartlett - Sparse greedy Gaussian process regression - NIPS 2001

<sup>8</sup>M Seeger, C Williams - Fast forward selection to speed up sparse Gaussian process regression - AISTATS 2003

<sup>9</sup>E Snelson, Z Ghahramani - Sparse Gaussian processes using pseudo-inputs - NIPS 2005

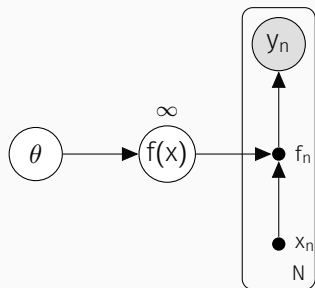
<sup>10</sup>J Quiñero-Candela, CE Rasmussen - A unifying view of sparse approximate Gaussian process regression - JMLR 2005

<sup>11</sup>M. Titsias - Variational learning of inducing variables in sparse Gaussian processes - AISTATS 2009

<sup>12</sup>J. Hensman, N. Fusi and N. Lawrence - Gaussian Processes for Big Data - UAI 2013



# A GRAPHICAL MODEL FOR GAUSSIAN PROCESSES



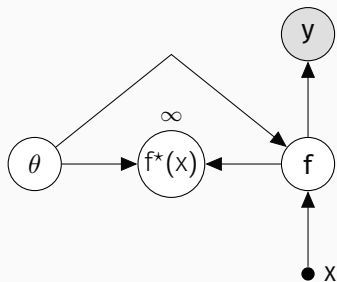
$$\theta \sim p(\theta)$$

$$f(x) \sim \mathcal{GP}(0, k(x, x'; \theta))$$

$$\mathbf{f} = [f(x_1), f(x_2) \dots f(x_n)]^T$$

$$y_n \sim p(y_n | f(x_n))$$

# A DIFFERENT GRAPHICAL MODEL FOR GAUSSIAN PROCESSES



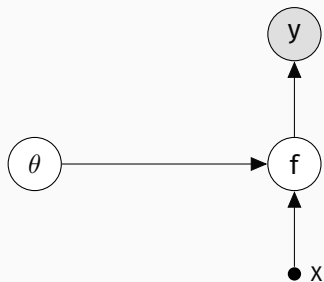
$$\theta \sim p(\theta)$$

$$\mathbf{f}|\theta \sim \mathcal{N}(0, \mathbf{K})$$

$$y_n \sim p(y_n | \mathbf{f}(x_n))$$

$$f^*(x)|\mathbf{f}, \theta \sim \mathcal{GP}(\mathbf{a}(x)^\top \mathbf{f}, b(x, x'))$$

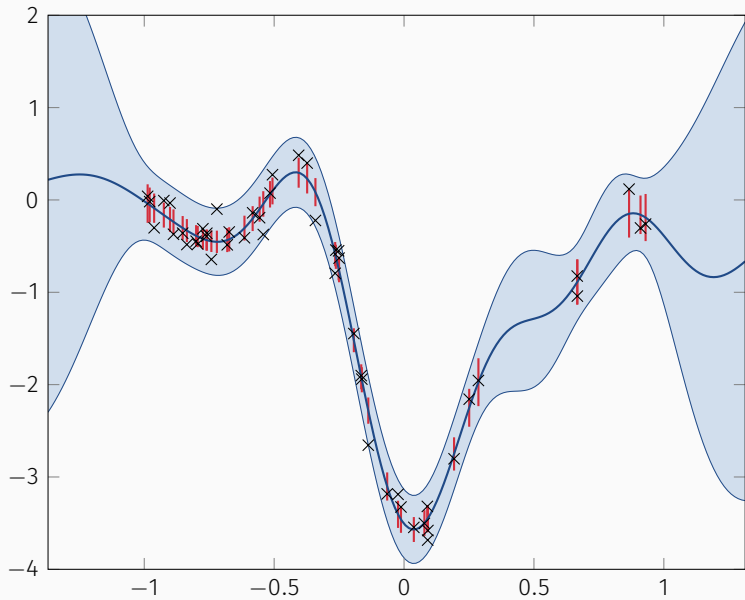
# A DIFFERENT GRAPHICAL MODEL FOR GAUSSIAN PROCESSES

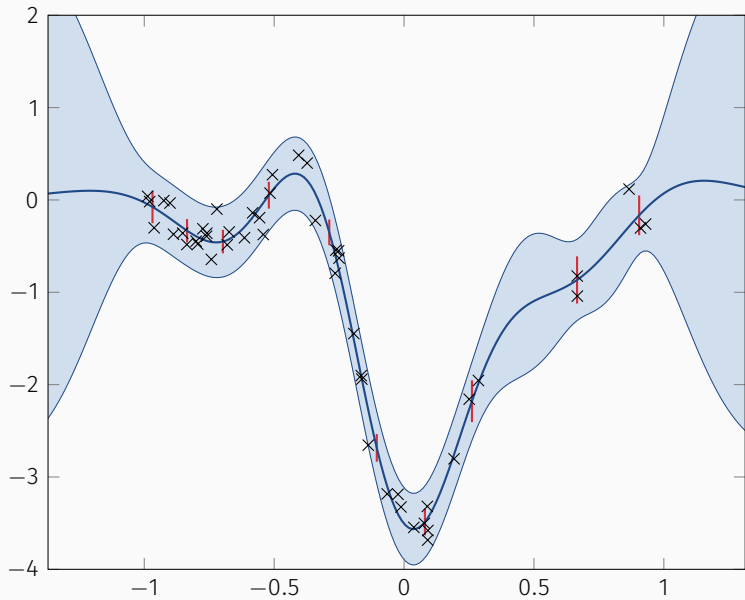


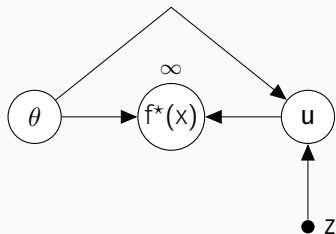
$$\theta \sim p(\theta)$$

$$f|\theta \sim \mathcal{N}(0, \mathbf{K})$$

$$y_n \sim p(y_n | f(x_n))$$







$$\theta, u \sim q(\theta, u)$$

$$f^*(x) \sim \mathcal{GP}(\mathbf{a}'(x)^\top \mathbf{u}, b'(x))$$

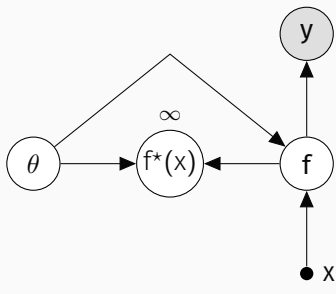
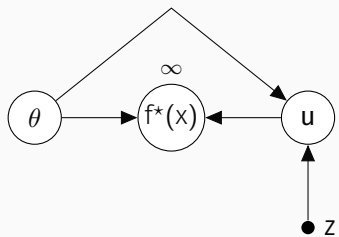
# KL DIVERGENCE BETWEEN GAUSSIAN PROCESSES?

Intuitive version:  $\mathbf{f}^*$  is a really long vector containing all points of interest.

Rigorous version: Matthews et al.<sup>13</sup>

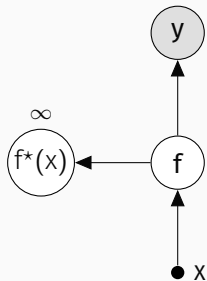
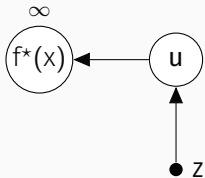
---

<sup>13</sup>On Sparse variational methods and the Kullback-Leibler divergence between stochastic processes <http://arxiv.org/abs/1504.07027>

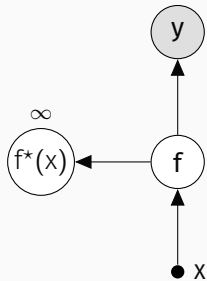
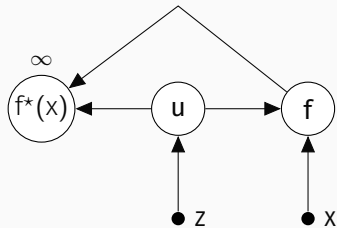




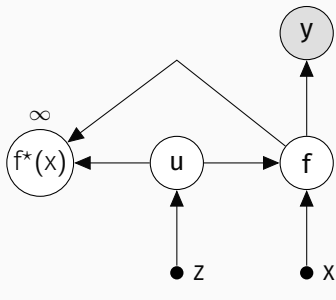
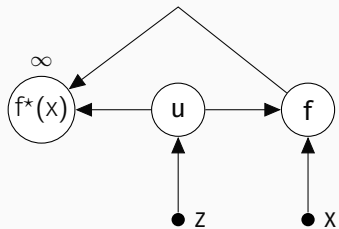
Let's ignore  $\theta$  for now

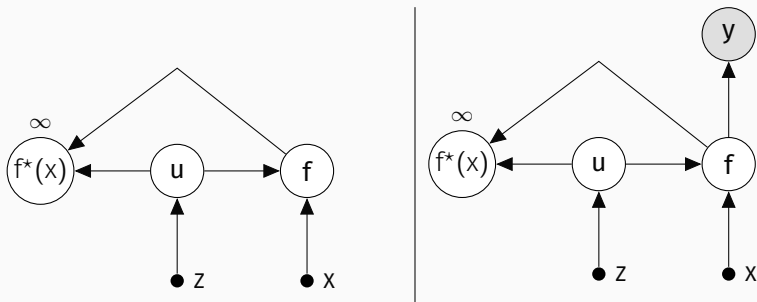


Where are the  $f$  in the approximation?

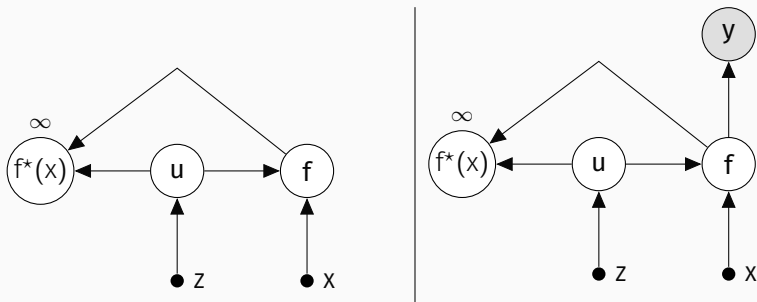


Where are the  $u$  in the model?

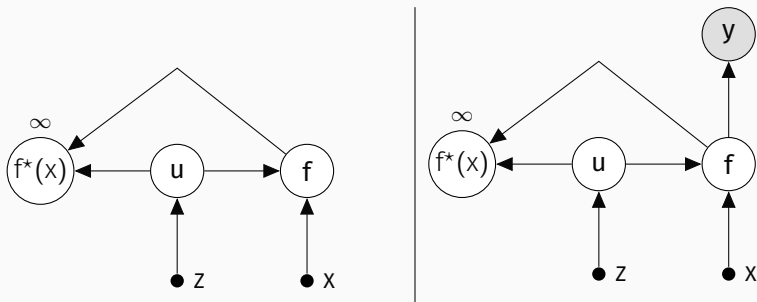




$$\text{ELBO} = \mathbb{E}_{q(f^*, f, u, \theta)} \left[ \log \frac{p(y | f)p(f | u, \theta)p(f^* | f, u, \theta)p(u | \theta)p(\theta)}{q(f | u, \theta)q(f^* | f, u, \theta)q(u | \theta)q(\theta)} \right]$$



$$\text{ELBO} = \mathbb{E}_{q(\cancel{f^*}, f, u, \theta)} \left[ \log \frac{p(y | f)p(f | u, \theta)\cancel{p(f^* | f, u, \theta)}p(u | \theta)p(\theta)}{q(f | u, \theta)\cancel{q(f^* | f, u, \theta)}q(u | \theta)q(\theta)} \right]$$



$$\text{ELBO} = \mathbb{E}_{q(\cancel{f^*}, f, u, \theta)} \left[ \log \frac{p(y | f) p(\cancel{f} | u, \theta) p(\cancel{f^*} | f, u, \theta) p(u | \theta) p(\theta)}{q(\cancel{f} | u, \theta) q(\cancel{f^*} | f, u, \theta) q(u | \theta) q(\theta)} \right]$$

## Strategy 1: Gaussian<sup>14</sup>

Let  $q(\mathbf{u}, \theta) = \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{L}\mathbf{L}^\top) \delta(\theta - \hat{\theta})$

Optimize wrt  $\mathbf{m}, \mathbf{L}, \hat{\theta}$  (and  $\mathbf{Z}$ !)

## Strategy 2: Free-form<sup>15</sup>

Given the limited size of  $\mathbf{Z}$  (and thus  $\mathbf{u}$ ), write down the optimal, intractable, form for  $q(\mathbf{u}, \theta)$ , and sample from it using HMC.

---

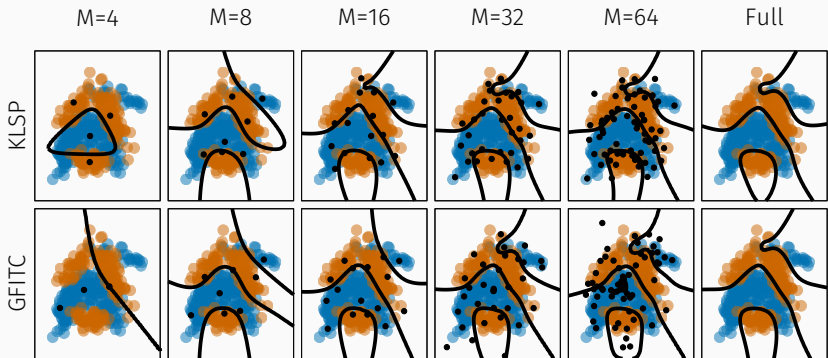
<sup>14</sup> Hensman, A Matthews, Z Ghahramani - Scalable Variational Gaussian Process Classification - AISTATS 2015

<sup>15</sup> Hensman, AGG Matthews, M Filippone - MCMC for Variationally Sparse Gaussian Processes - NIPS 2015

The objective function (which minimizes the KL between the q-process and the p-process) is

$$\mathcal{L} = \sum_i \mathbb{E}_{q(f_i)}[\log p(\mathbf{y}_i | f_i)] - \text{KL}[q(u) || p(u)]$$







Left: three k-means centers used to initialize the inducing point positions. Center: the positions of the same inducing points after optimization. Right: difference.

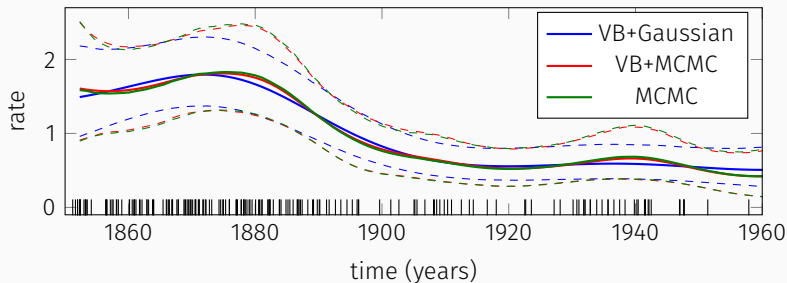
Data:  $N=60,000$ ,  $D=784$

Accuracy: 98.04%

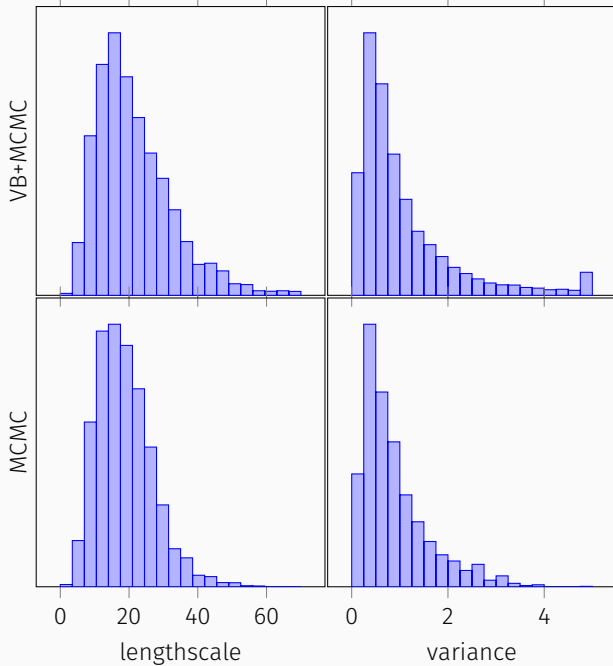
The 'perfect' distribution  $\hat{q}(\mathbf{u}, \theta)$  which minimises the KL divergence (with no further restrictions) is

$$\log \hat{q}(\mathbf{u}, \theta) = \mathbb{E}_{p(\mathbf{f} | \mathbf{u})}[\log p(\mathbf{y} | \mathbf{f})] + \log p(\mathbf{u}, \theta) + \text{const.}$$

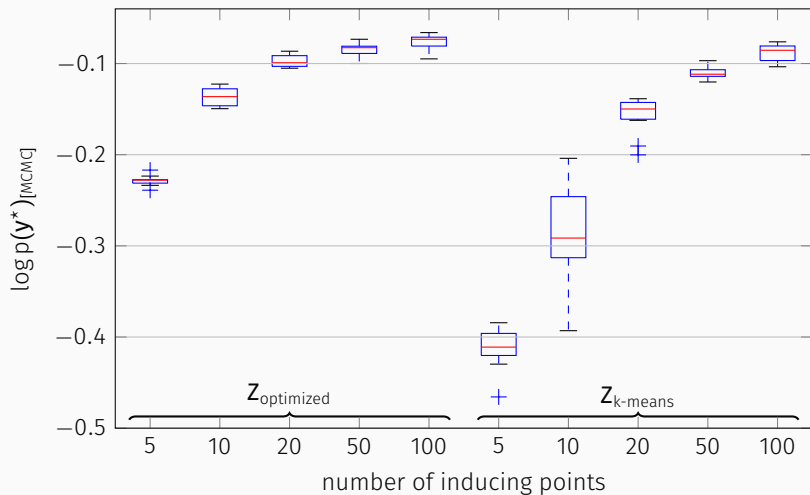
Sampling  $\hat{q}$  costs  $\mathcal{O}(NM^2)$ .



The posterior of the rates for the coal mining disaster data.



# THE EFFECT OF INDUCING POINTS SELECTION



- Exact inference (Gaussian likelihood,  $\mathbf{Z} = \mathbf{X}$ )
- Subset-of-data methods (e.g. IVM <sup>16</sup>)
- Inter-domain approximations <sup>17</sup>
- Black box likelihoods <sup>18</sup>
- Log Gaussian Cox processes <sup>19</sup>

---

<sup>16</sup>Lawrence, Seeger and Herbrich - The Informative Vector Machine - NIPS 2003

<sup>17</sup>Alvarez, Rosasco and Lawrence - Kernels for vector valued functions, a review - Foundations and Trends in ML 2011

<sup>18</sup>Dezfouli and Bonilla - Gaussian Process Models with Black-Box Likelihoods - NIPS 2015

<sup>19</sup>Lloyd et al - Variational Inference for Gaussian Process Modulated Poisson Processes - ICML 2015

THANKS FOR LISTENING