

SPARSE GAUSSIAN PROCESSES

GPSS 2015, Sheffield

James Hensman

September 15, 2015

Lancaster University

Motivation

A History lesson

Posteriors over functions

Posteriors over inducing points

Prediction and the KL between processes

Summary and demo

MOTIVATION

Inference in a GP has the following demands:

Complexity: $\mathcal{O}(n^3)$

Storage: $\mathcal{O}(n^2)$

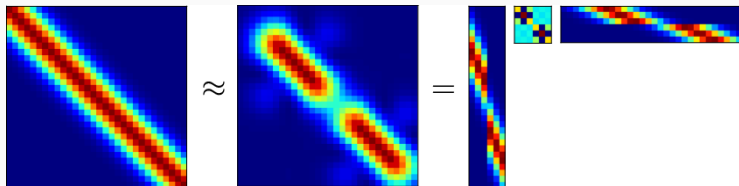
Inference in a *sparse* GP has the following demands:

Complexity: $\mathcal{O}(nm^2)$

Storage: $\mathcal{O}(nm)$

where we get to pick m !

HOW TO MAKE COMPUTATIONAL SAVINGS



$$\mathbf{K}_{nn} \approx \mathbf{Q}_{nn} = \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn}$$

Instead of inverting \mathbf{K}_{nn} , we make a low rank (or Nyström) approximation, and invert \mathbf{K}_{mm} instead.

A HISTORY LESSON

WHY ARE THEY CALLED SPARSE GPS?

Sparse (adj). From spagare, meaning “few and scattered”.

Subset of data

- Silverman 1985 (subset of regressors)
- Smola and Bartlett 2001 (greedy selection)

Pseudo-input approximations

- Snelson and Ghahramani (2005), Snelson (2007)

Variational approximations

- Titsias (2009) – derived a variational bound
- Matthews et al. (2015) – showed this minimised KL between processes

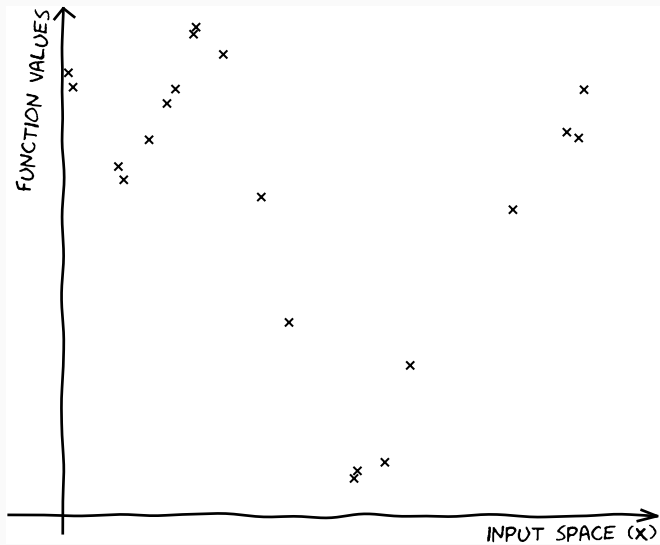
POSTERIOR OVER FUNCTIONS

Everything we want to do with a GP involves marginalising \mathbf{f}

- Predictions
- Marginal likelihood
- Estimating covariance parameters

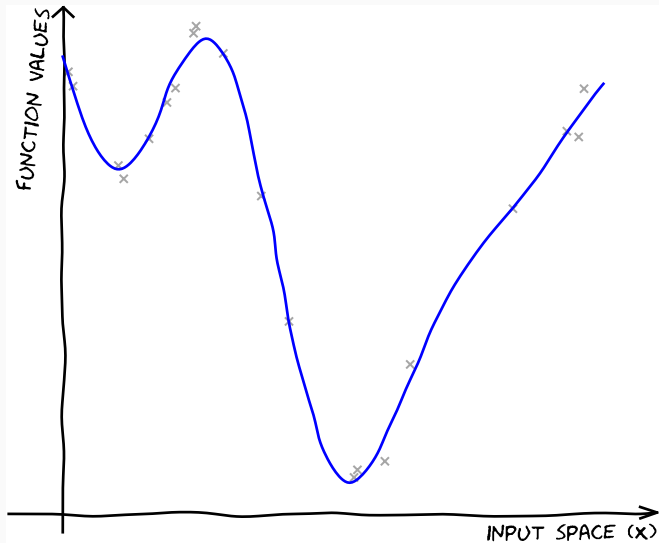
The posterior of \mathbf{f} is the central object. This means inverting \mathbf{K}_{nn} .

X, y



X, y

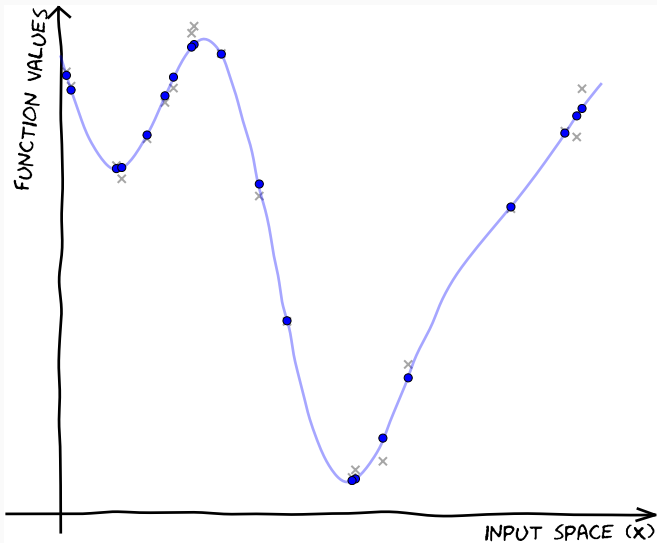
$f(x) \sim \mathcal{GP}$



X, y

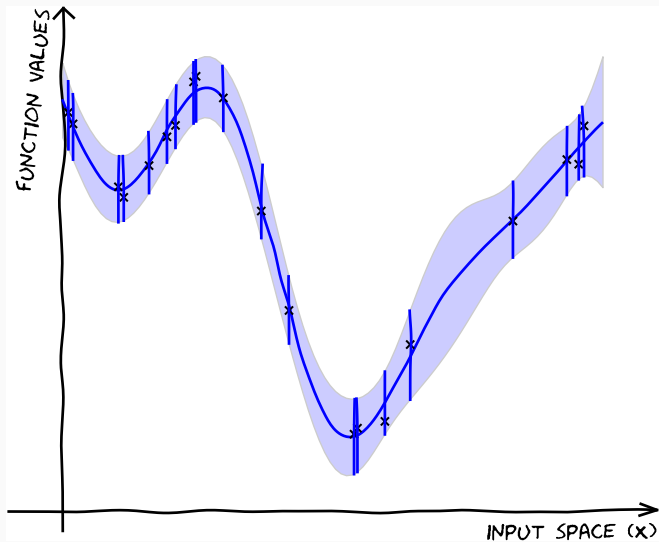
$f(x) \sim \mathcal{GP}$

$p(f) = \mathcal{N}(0, K_{nn})$



$$X, y$$
$$f(x) \sim \mathcal{GP}$$
$$p(f) = \mathcal{N}(0, K_{nn})$$

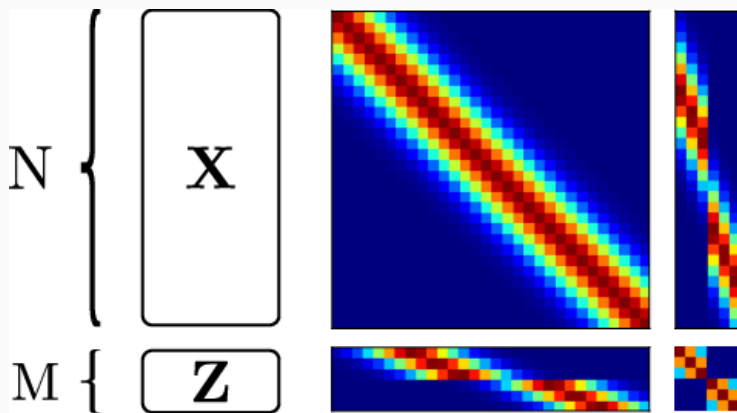
$$p(f | y, X)$$
$$p(f^* | f, X, x^*)$$



POSTERIOR OVER INDUCING POINTS

Take an extra M points on the function, $\mathbf{u} = \mathbf{f}(\mathbf{Z})$.

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



Take and extra M points on the function, $\mathbf{u} = \mathbf{f}(\mathbf{Z})$.

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

$$p(\mathbf{y} | \mathbf{f}) = \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{f} | \mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{u}, \tilde{\mathbf{K}})$$

$$p(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{mm})$$

X, y

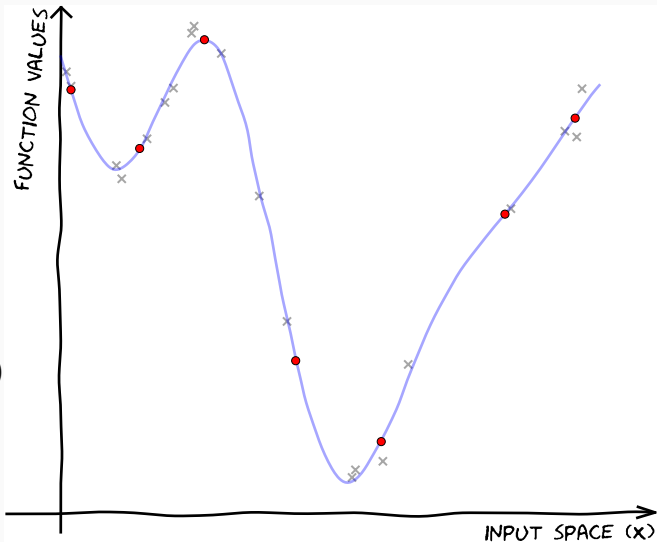
$f(x) \sim \mathcal{GP}$

$p(f) = \mathcal{N}(0, K_{nn})$

$p(f|y, X)$

Z, u

$p(u) = \mathcal{N}(0, K_{mm})$



X, y

$f(x) \sim \mathcal{GP}$

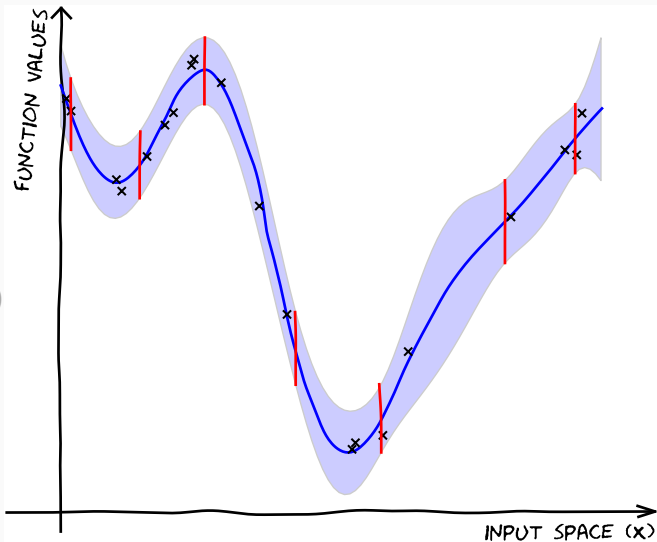
$p(f) = \mathcal{N}(0, K_{nn})$

$p(f|y, X)$

$p(u) = \mathcal{N}(0, K_{mm})$

$\tilde{p}(u|y, X)$

$p(f^* | u, Z, x^*)$



Instead of doing

$$p(\mathbf{f} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{X})}{\int p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{X})d\mathbf{f}}$$

We'll do

$$p(\mathbf{u} | \mathbf{y}, \mathbf{Z}) = \frac{p(\mathbf{y} | \mathbf{u})p(\mathbf{u} | \mathbf{Z})}{\int p(\mathbf{y} | \mathbf{u})p(\mathbf{u} | \mathbf{Z})d\mathbf{u}}$$

Instead of doing

$$p(\mathbf{f} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{X})}{\int p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{X})d\mathbf{f}}$$

We'll do

$$p(\mathbf{u} | \mathbf{y}, \mathbf{Z}) = \frac{p(\mathbf{y} | \mathbf{u})p(\mathbf{u} | \mathbf{Z})}{\int p(\mathbf{y} | \mathbf{u})p(\mathbf{u} | \mathbf{Z})d\mathbf{u}}$$

but $p(\mathbf{y} | \mathbf{u})$ involves inverting \mathbf{K}_{nn}

$$p(\mathbf{y} | \mathbf{u}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$p(\mathbf{y} | \mathbf{u}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln p(\mathbf{y} | \mathbf{f}) + \ln \frac{p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$p(\mathbf{y} | \mathbf{u}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln p(\mathbf{y} | \mathbf{f}) + \ln \frac{p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \mathbb{E}_{p(\mathbf{f} | \mathbf{u})} [\ln p(\mathbf{y} | \mathbf{f})] + \mathbb{E}_{p(\mathbf{f} | \mathbf{u})} \left[\ln \frac{p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})} \right]$$

$$p(\mathbf{y} | \mathbf{u}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln p(\mathbf{y} | \mathbf{f}) + \ln \frac{p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \mathbb{E}_{p(\mathbf{f} | \mathbf{u})} [\ln p(\mathbf{y} | \mathbf{f})] + \mathbb{E}_{p(\mathbf{f} | \mathbf{u})} \left[\ln \frac{p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y}, \mathbf{u})} \right]$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \tilde{p}(\mathbf{y} | \mathbf{u}) + \text{KL}[p(\mathbf{f} | \mathbf{u}) || p(\mathbf{f} | \mathbf{y}, \mathbf{u})]$$

No inversion of \mathbf{K}_{nn} required

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{u}, \mathbf{X}) d\mathbf{f}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{u}, \mathbf{X}) d\mathbf{f}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \mathbb{E}_{p(\mathbf{f} | \mathbf{u}, \mathbf{X})} [p(\mathbf{y} | \mathbf{f})]$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{u}, \mathbf{X}) d\mathbf{f}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \mathbb{E}_{p(\mathbf{f} | \mathbf{u}, \mathbf{X})} [p(\mathbf{y} | \mathbf{f})]$$

$$\ln p(\mathbf{y} | \mathbf{u}) \geq \mathbb{E}_{p(\mathbf{f} | \mathbf{u}, \mathbf{X})} [\ln p(\mathbf{y} | \mathbf{f})] \triangleq \ln \tilde{p}(\mathbf{y} | \mathbf{u})$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{u}, \mathbf{X}) d\mathbf{f}$$

$$\ln p(\mathbf{y} | \mathbf{u}) = \ln \mathbb{E}_{p(\mathbf{f} | \mathbf{u}, \mathbf{X})} [p(\mathbf{y} | \mathbf{f})]$$

$$\ln p(\mathbf{y} | \mathbf{u}) \geq \mathbb{E}_{p(\mathbf{f} | \mathbf{u}, \mathbf{X})} [\ln p(\mathbf{y} | \mathbf{f})] \triangleq \ln \tilde{p}(\mathbf{y} | \mathbf{u})$$

No inversion of \mathbf{K}_{nn} required

$$\tilde{p}(\mathbf{y} | \mathbf{u}) = \prod_{i=1}^n \mathcal{N}(y_i | \mathbf{k}_{mn}^\top \mathbf{K}_{mm}^{-1} \mathbf{u}, \sigma^2) \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{k}_{nn} - \mathbf{k}_{mn}^\top \mathbf{K}_{mm}^{-1} \mathbf{k}_{mn}) \right\}$$

A straightforward likelihood approximation, and a penalty term

$$\tilde{p}(\mathbf{u} | \mathbf{y}, \mathbf{Z}) = \frac{\tilde{p}(\mathbf{y} | \mathbf{u})p(\mathbf{u} | \mathbf{Z})}{\int \tilde{p}(\mathbf{y} | \mathbf{u})p(\mathbf{u} | \mathbf{Z})d\mathbf{u}}$$

- Computing the (approximate) posterior costs $\mathcal{O}(nm^2)$
- We also get a lower bound of the marginal likelihood
- This is the standard variational sparse GP as Titsias 2009
- looks like a low rank approximation.

$$\begin{aligned}\tilde{p}(\mathbf{y}) &= \int \tilde{p}(\mathbf{y} | \mathbf{u}) p(\mathbf{u} | \mathbf{Z}) d\mathbf{u} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} + \sigma^2 \mathbf{I}) \exp \sum_i \left\{ -\frac{1}{2\sigma^2} (\mathbf{k}_{nn} - \mathbf{k}_{mn}^T \mathbf{K}_{mm}^{-1} \mathbf{k}_{mn}) \right\}\end{aligned}$$

The variational objective $\ln \tilde{p}(\mathbf{y})$ is a function of

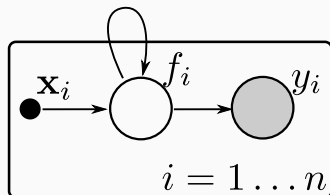
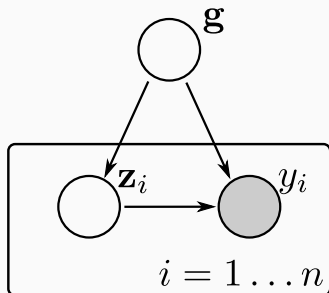
- the parameters of the covariance function $\boldsymbol{\theta}$
- the inducing inputs, \mathbf{Z}

Strategy: jointly optimize $\boldsymbol{\theta}$ and \mathbf{Z} .

DISTRIBUTED COMPUTATION AND STOCHASTIC OPTIMIZATION

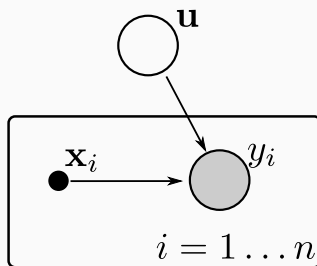
STOCHASTIC VARIATIONAL INFERENCE

- Combine the ideas of stochastic optimisation with Variational inference
- example: apply Latent Dirichlet allocation to project Gutenberg
- Can apply variational techniques to Big Data
- How could this work in GPs?



MAINTAIN THE FACTORISATION!

- The variational marginalisation of \mathbf{f} introduced factorisation across the datapoints (conditioned on \mathbf{u})
- Marginalising \mathbf{u} re-introduced dependencies between the data
- Solution: a variational treatment of \mathbf{u}



$$\log p(\mathbf{y} | \mathbf{X}) \geq \langle \log \tilde{p}(\mathbf{y} | \mathbf{u}) + \log p(\mathbf{u}) - \log q(\mathbf{u}) \rangle_{q(\mathbf{u})} \triangleq \mathcal{L}. \quad (1)$$

$$\begin{aligned} \mathcal{L} = \sum_{i=1}^n \left\{ \right. & \log \mathcal{N}(y_i | \mathbf{k}_{mn}^T \mathbf{K}_{mm}^{-1} \mathbf{m}, \beta^{-1}) \\ & - \frac{1}{2} \beta \tilde{k}_{i,i} - \frac{1}{2} \text{tr}(\mathbf{S} \boldsymbol{\Lambda}_i) \left. \right\} \\ & - \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) \end{aligned} \quad (2)$$

The variational objective \mathcal{L} is a function of

- the parameters of the covariance function
- the parameters of $q(\mathbf{u})$
- the inducing inputs, \mathbf{Z}

Original strategy: set \mathbf{Z} . Take the data in small minibatches, take stochastic gradient steps in the covariance function parameters, stochastic natural gradient steps in the parameters of $q(\mathbf{u})$.

New strategy: optimize everything jointly with AdaDelta.

$$\tilde{\mathbf{g}}(\boldsymbol{\theta}) = \mathbf{G}(\boldsymbol{\theta})^{-1} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{L}}{\partial \boldsymbol{\eta}}.$$

$$\begin{aligned} \boldsymbol{\theta}_{2(t+1)} &= -\frac{1}{2} \mathbf{S}_{(t+1)} \\ &= -\frac{1}{2} \mathbf{S}_{(t)} + \ell \left(-\frac{1}{2} \boldsymbol{\Lambda} + \frac{1}{2} \mathbf{S}_{(t)} \right), \end{aligned}$$

$$\begin{aligned} \boldsymbol{\theta}_{1(t+1)} &= \mathbf{S}_{(t+1)} \mathbf{m}_{(t+1)} \\ &= \mathbf{S}_{(t)} \mathbf{m}_{(t)} + \ell \left(\beta \mathbf{K}_{mm} \mathbf{K}_{mn} \mathbf{y} - \mathbf{S}_{(t)} \mathbf{m}_{(t)} \right), \end{aligned}$$

PREDICTION AND THE KL BETWEEN PROCESSES

We have minimised the KL divergence

$$\text{KL}[\tilde{p}(\mathbf{u})p(\mathbf{f}|\mathbf{u})||p(\mathbf{f}, \mathbf{u} | \mathbf{y})]$$

We have minimised the KL divergence

$$\text{KL}[\tilde{p}(\mathbf{u})p(\mathbf{f}|\mathbf{u})||p(\mathbf{f}, \mathbf{u} | \mathbf{y})]$$

but this turns out to be equivalent to

$$\text{KL}[p(\mathbf{f}^* | \mathbf{u})\tilde{p}(\mathbf{u})||p(\mathbf{f}^* | \mathbf{y})p(\mathbf{f} | \mathbf{y})]$$

To predict, just compute the required quantity of the variational stochastic process.

$$p(\mathbf{f}^* | \mathbf{y}) \approx \int p(\mathbf{f}^* | \mathbf{u}) \tilde{p}(\mathbf{u} | \mathbf{y}) d\mathbf{u}$$

SUMMARY AND DEMO

- I have guided you through the variational sparse GP method (for regression).

- I have guided you through the variational sparse GP method (for regression).
- Great framework for extensibility
 - Non Gaussian likelihoods
 - Multiple outputs
 - Stochastic optimization
 - ...

- I have guided you through the variational sparse GP method (for regression).
- Great framework for extensibility
 - Non Gaussian likelihoods
 - Multiple outputs
 - Stochastic optimization
 - ...
- Move away from thinking of a model approximation: separate model and inference

- I have guided you through the variational sparse GP method (for regression).
- Great framework for extensibility
 - Non Gaussian likelihoods
 - Multiple outputs
 - Stochastic optimization
 - ...
- Move away from thinking of a model approximation: separate model and inference
- Work of many authors (406,000 scholar hits!). Apologies for the 405,995 omissions