

Gaussian Processes for Big Data

James Hensman

Overview

Motivation

Sparse Gaussian Processes

Stochastic Variational Inference

Examples

Overview

Motivation

Sparse Gaussian Processes

Stochastic Variational Inference

Examples

Motivation

Inference in a GP has the following demands:

$$\begin{aligned} \text{Complexity: } & O(n^3) \\ \text{Storage: } & O(n^2) \end{aligned}$$

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

$$\begin{aligned} \text{Complexity: } & O(nm^2) \\ \text{Storage: } & O(nm) \end{aligned}$$

where we get to pick m !

Still not good enough!

Big Data

- ▶ In parametric models, stochastic optimisation is used.
- ▶ This allows for application to Big Data.

This work

- ▶ Show how to use Stochastic Variational Inference in GPs
- ▶ Stochastic optimisation scheme: each step requires $O(m^3)$

Overview

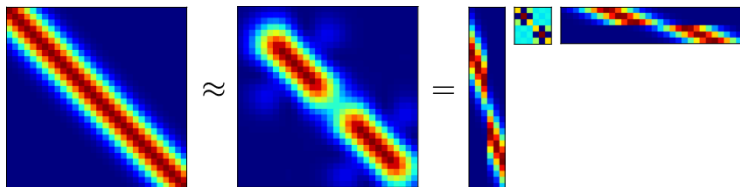
Motivation

Sparse Gaussian Processes

Stochastic Variational Inference

Examples

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.

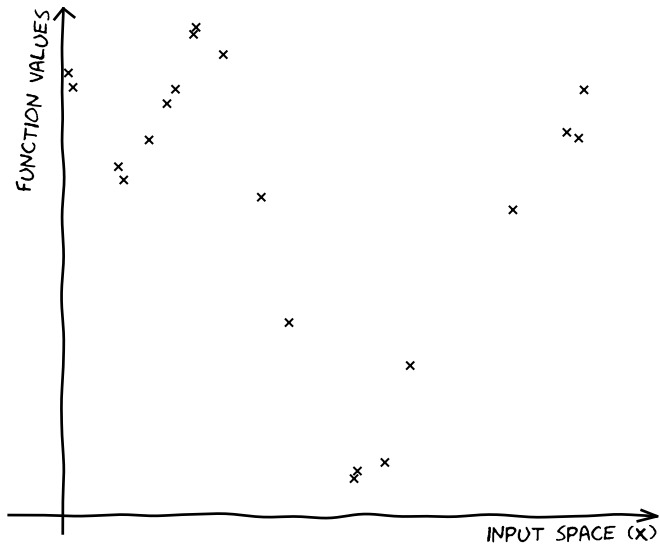
Information capture

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

- ▶ Predictions
- ▶ Marginal likelihood
- ▶ Estimating covariance parameters

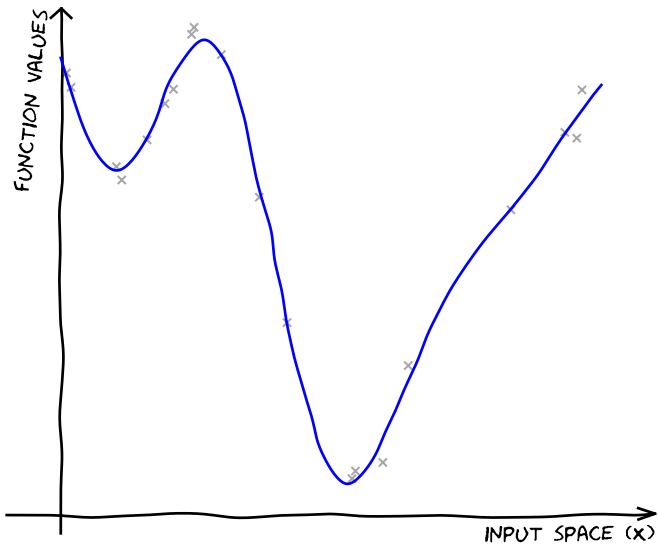
The posterior of \mathbf{f} is the central object. This means inverting \mathbf{K}_{nn} .

X, y



X, y

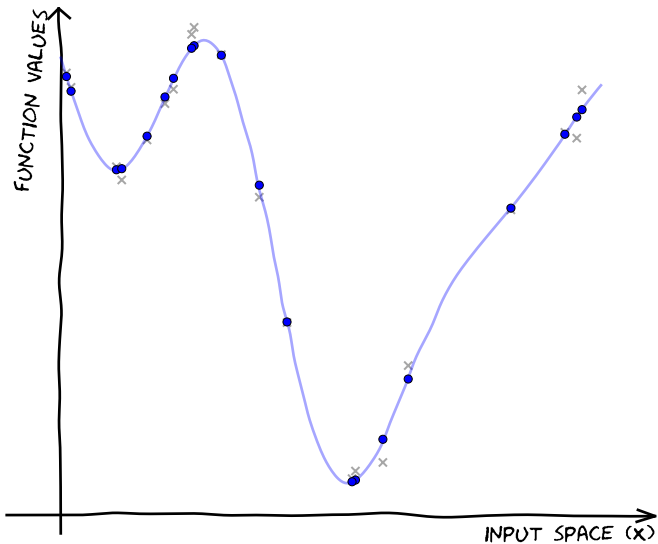
$f(x) \sim GP$



X, y

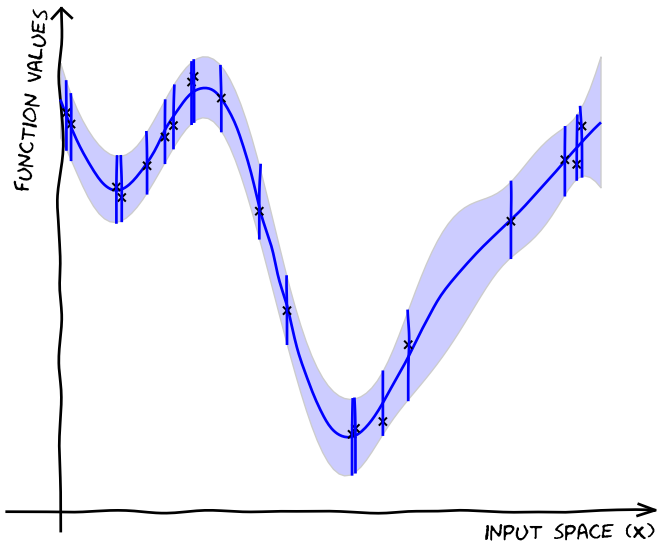
$f(x) \sim \mathcal{GP}$

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



$$\mathbf{X}, \mathbf{y}$$
$$f(\mathbf{x}) \sim \mathcal{GP}$$
$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{mn})$$

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

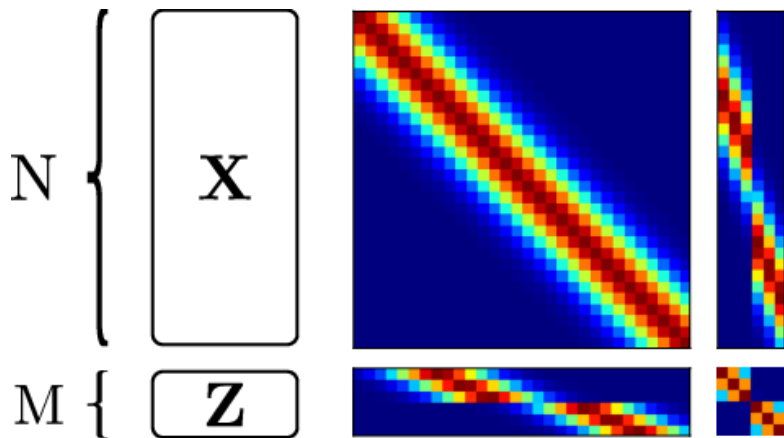


Introducing \mathbf{u}

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

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

Introducing \mathbf{u}



Introducing \mathbf{u}

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

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

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

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

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

\mathbf{X}, \mathbf{y}

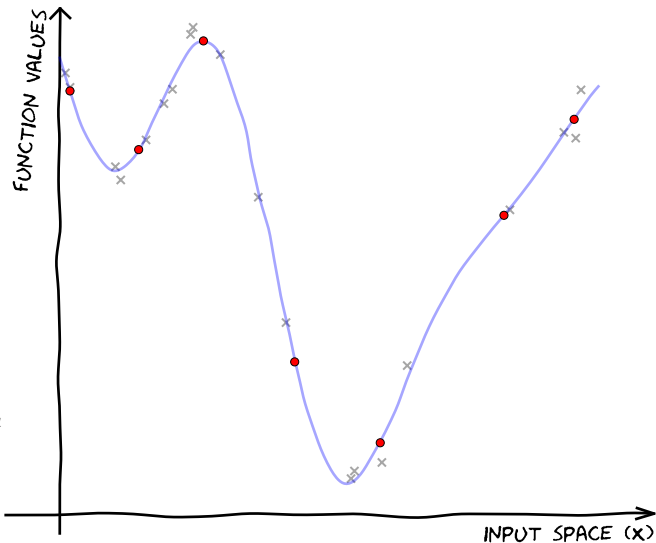
$f(\mathbf{x}) \sim \mathcal{GP}$

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

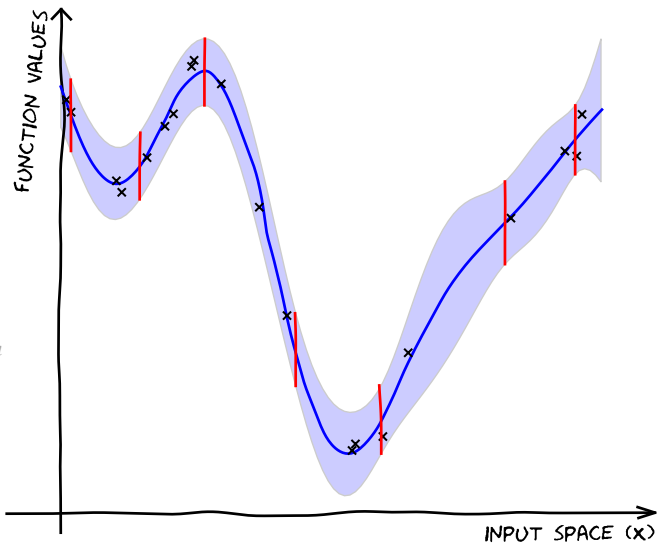
$p(\mathbf{f} | \mathbf{y}, \mathbf{X})$

\mathbf{Z}, \mathbf{u}

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



$$\mathbf{X}, \mathbf{y}$$
$$f(\mathbf{x}) \sim \mathcal{GP}$$
$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{mm})$$
$$p(\mathbf{f} | \mathbf{y}, \mathbf{X})$$
$$p(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{mm})$$
$$\tilde{p}(\mathbf{u} | \mathbf{y}, \mathbf{X})$$



The alternative posterior

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

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

The alternative posterior

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}_{mn}

Variational marginalisation of \mathbf{f}

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

Variational marginalisation of \mathbf{f}

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

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

Variational marginalisation of \mathbf{f}

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

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

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

Variational marginalisation of \mathbf{f}

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

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

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

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

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

An approximate likelihood

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

A straightforward likelihood approximation, and a penalty term

Now we can marginalise \mathbf{u}

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

- ▶ Computing the (approximate) posterior costs $O(nm^2)$
- ▶ We also get a lower bound of the marginal likelihood
- ▶ This is the standard variational sparse GP [Titsias, 2009].

Overview

Motivation

Sparse Gaussian Processes

Stochastic Variational Inference

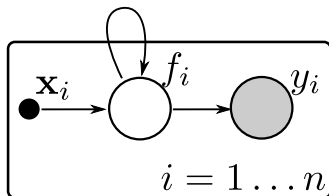
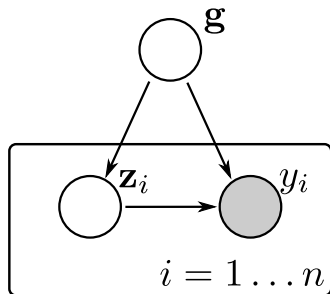
Examples

Variational Bayes

- ▶ Approximate the true posterior distribution with a simpler one.
- ▶ Usually assume factorisation in the approximation
- ▶ Iterative 'update' procedure (like EM)
- ▶ Can be seen as a coordinate-wise steepest ascent method

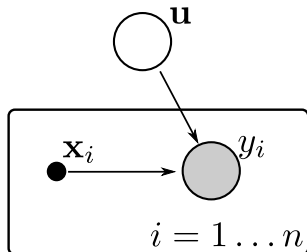
Stochastic Variational Inference

- ▶ 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}



$$\log p(\mathbf{y} | \mathbf{X}) \geq \langle \mathcal{L}_1 + \log p(\mathbf{u}) - \log q(\mathbf{u}) \rangle_{q(\mathbf{u})} \triangleq \mathcal{L}_3. \quad (1)$$

$$\begin{aligned} \mathcal{L}_3 = \sum_{i=1}^n & \left\{ \log \mathcal{N}(y_i | \mathbf{k}_{mn}^\top \mathbf{K}_{mm}^{-1} \mathbf{m}, \beta^{-1}) \right. \\ & \left. - \frac{1}{2} \beta \tilde{k}_{i,i} - \frac{1}{2} \text{tr}(\mathbf{S} \boldsymbol{\Lambda}_i) \right\} \\ & - \text{KL}(q(\mathbf{u}) \| p(\mathbf{u})) \end{aligned} \quad (2)$$

Optimisation

The variational objective \mathcal{L}_3 is a function of

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

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 $q(\mathbf{u})$.

Natural Gradients

$$\widetilde{\mathbf{g}}(\boldsymbol{\theta}) = G(\boldsymbol{\theta})^{-1} \frac{\partial \mathcal{L}_3}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{L}_3}{\partial \boldsymbol{\eta}}.$$

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

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

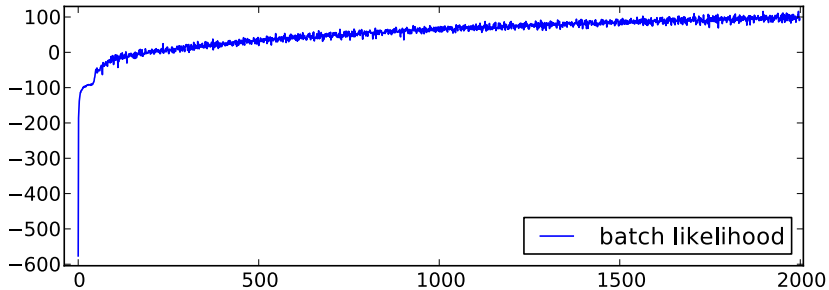
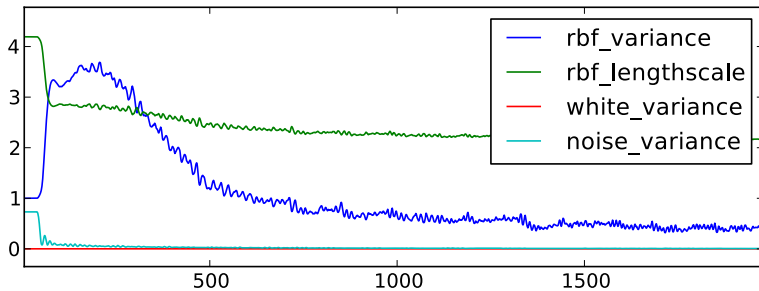
Overview

Motivation

Sparse Gaussian Processes

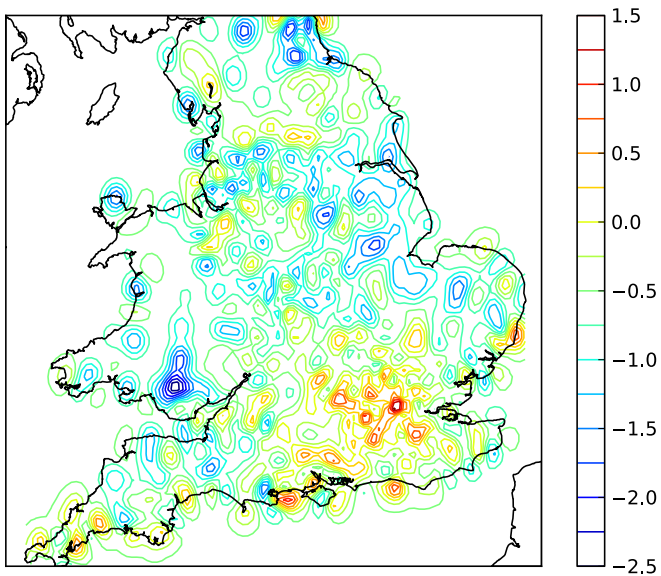
Stochastic Variational Inference

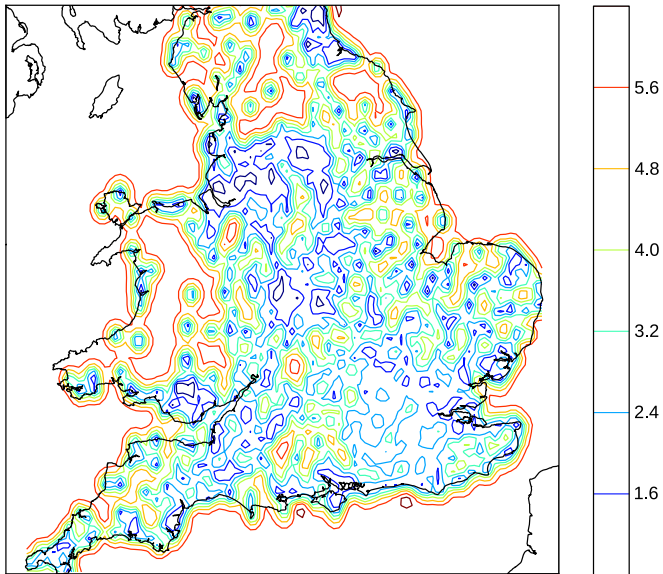
Examples



UK apartment prices

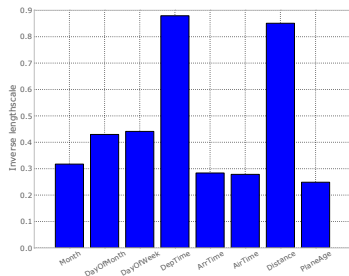
- ▶ Monthly price paid data for February to October 2012 (England and Wales)
- ▶ from <http://data.gov.uk/dataset/land-registry-monthly-price-paid-data/>
- ▶ 75,000 entries
- ▶ Cross referenced against a postcode database to get latitude and longitude
- ▶ Regressed the normalised logarithm of the apartment prices

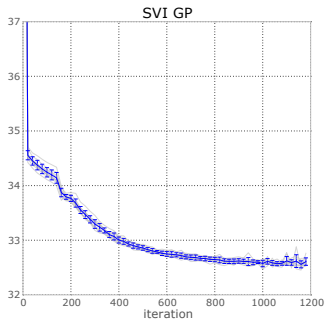
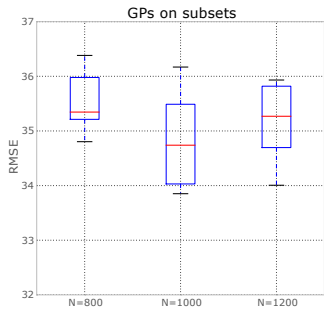


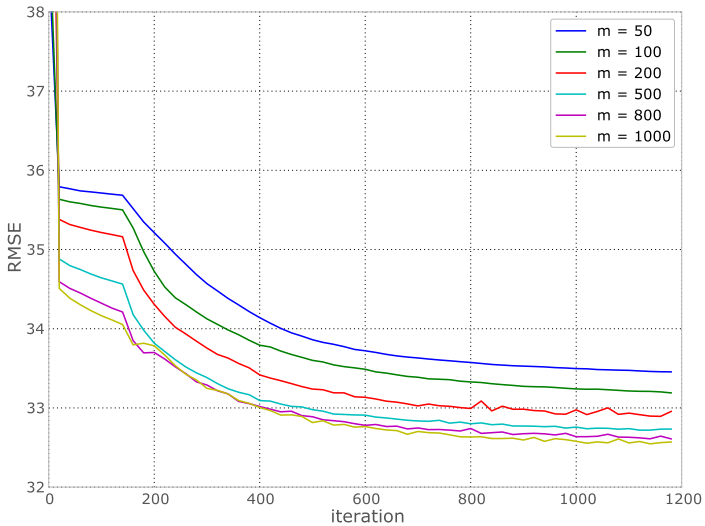


Airline data

- ▶ Flight delays for every commercial flight in the USA from January to April 2008.
- ▶ Average delay was 30 minutes.
- ▶ We randomly selected 800,000 datapoints (we have limited memory!)
- ▶ 700,000 train, 100,000 test







Download the code!

`github.com/SheffieldML/GPy`

Cite our paper!

Hensman, Fusi and Lawrence,
Gaussian Processes for Big Data
Proceedings of UAI 2013

Michalis K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In David van Dyk and Max Welling, editors, *Proceedings of the Twelfth International Workshop on Artificial Intelligence and Statistics*, volume 5, pages 567–574, Clearwater Beach, FL, 16-18 April 2009. JMLR W&CP 5.