

CS-E4895 Gaussian Processes

Lecture 7: Large-scale Gaussian Processes

Arno Solin

Aalto University

Monday 20.3.2023

based on slides by Zhenwen Dai, Michael Riis Andersen, Ti John, and Arno Solin

Roadmap for today

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
- 6 Variational inference for sparse GPs
- 7 Recap

Section 1

What is the computational issue?

Gaussian process regression: Operations

Learning hyperparameters: Maximize (log) marginal likelihood:

$$\theta^* = \arg \max_{\theta} \log p(\mathbf{y} \mid X, \theta) = \arg \max_{\theta} \log \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$$

Test point prediction (mean, variance):

$$p(f_* \mid \mathbf{y}) = \mathcal{N}(f_* \mid \mu_*, \sigma_*^2)$$

$$\mu_* = \mathbf{k}_{f_* f} (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\sigma_*^2 = k_{f_* f_*} - \mathbf{k}_{f_* f} (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{f_* f}^\top$$

Computational complexity of Gaussian process regression

Data set with N observations, computing posterior for 1 test point:

$$\mu_* = \mathbf{k}_{f_* f} (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\sigma_*^2 = k_{f_* f_*} - \mathbf{k}_{f_* f} (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{f_* f}^\top$$

What is computational complexity? What is the dominating operation?

- Matrix–vector multiplication (mvm):
for $\mathbf{A} \in \mathbb{R}^{N \times M}$ and $\mathbf{b} \in \mathbb{R}^M$, computing \mathbf{Ab} costs $\mathcal{O}(NM)$
- Matrix inverse: for $\mathbf{C} \in \mathbb{R}^{N \times N}$, computing \mathbf{C}^{-1} costs $\mathcal{O}(N^3)$
- $\boldsymbol{\alpha} = (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$ scales as $\mathcal{O}(N^3)$, $\mu_* = \mathbf{k}_{f_* f} \boldsymbol{\alpha}$ scales as $\mathcal{O}(N)$

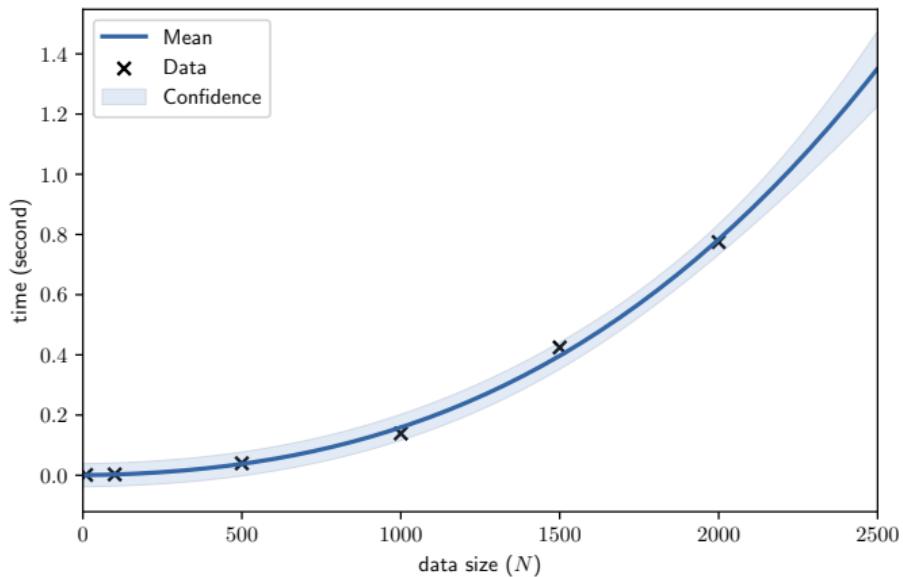
Log marginal likelihood

$$\begin{aligned}\log p(\mathbf{y} \mid \mathbf{X}) &= \log \mathcal{N}(\mathbf{y} \mid 0, \mathbf{K} + \sigma^2 \mathbf{I}) \\ &= -\frac{1}{2} \log |2\pi(\mathbf{K} + \sigma^2 \mathbf{I})| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \\ &= -\frac{N}{2} \log 2\pi - \sum_i \log \mathbf{L}_{ii} - \frac{1}{2} \|\mathbf{L}^{-1} \mathbf{y}\|^2\end{aligned}$$

- Cholesky decomposition: $\mathbf{L} = \text{chol}(\mathbf{K} + \sigma^2 \mathbf{I})$, such that $\mathbf{K} + \sigma^2 \mathbf{I} = \mathbf{L} \mathbf{L}^\top$.
- $\mathcal{O}(N^3)$ computational complexity

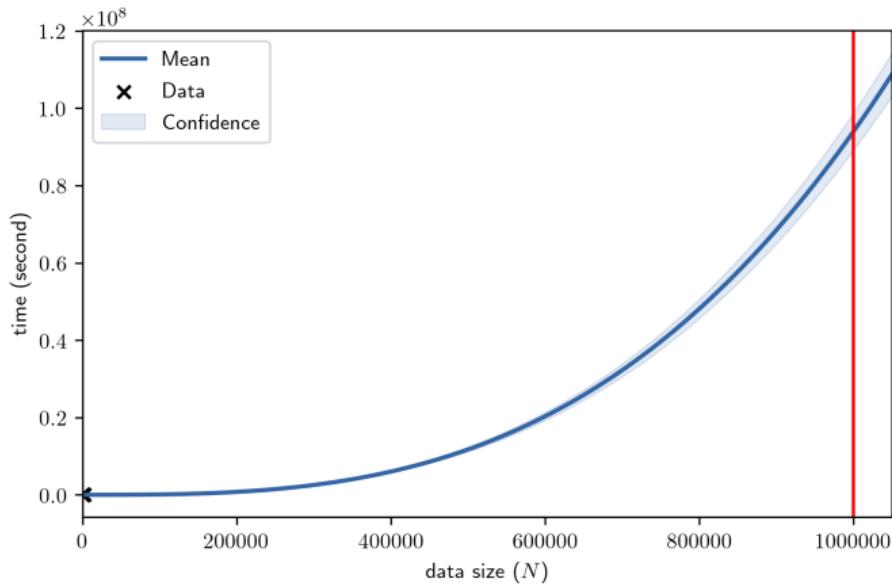
Empirical analysis of computational time

N	10	100	500	1000	1500	2000
time	1.3ms	8.5ms	28ms	0.12s	0.29s	0.76s



What if we have 1 million data points?

Predicted computational time: 9.4×10^7 seconds ≈ 2.98 years.



$N \leq 1000$: Fine, $N \leq 10000$: Slow, but possible, $N > 10000$: Prohibitively slow

Biggest bottlenecks

Line #	Time(ms)	% Time	Line Contents
2			def log_likelihood(kern, X, Y, sigma2):
3	6.0	0.0	N = X.shape[0]
4	55595.0	58.7	K = kern.K(X)
5	4369.0	4.6	Ky = K + np.eye(N)*sigma2
6	30012.0	31.7	L = np.linalg.cholesky(Ky)
7	4361.0	4.6	LinvY = dtrtrs(L, Y, lower=1)[0]
8	49.0	0.1	logL = N*np.log(2*np.pi)/-2.
9	82.0	0.1	logL += np.square(LinvY).sum()/-2.
10	208.0	0.2	logL += -np.log(np.diag(L)).sum()
11	2.0	0.0	return logL

- constructing \mathbf{K}
- computing Cholesky of \mathbf{K}

How can we address the bottleneck?

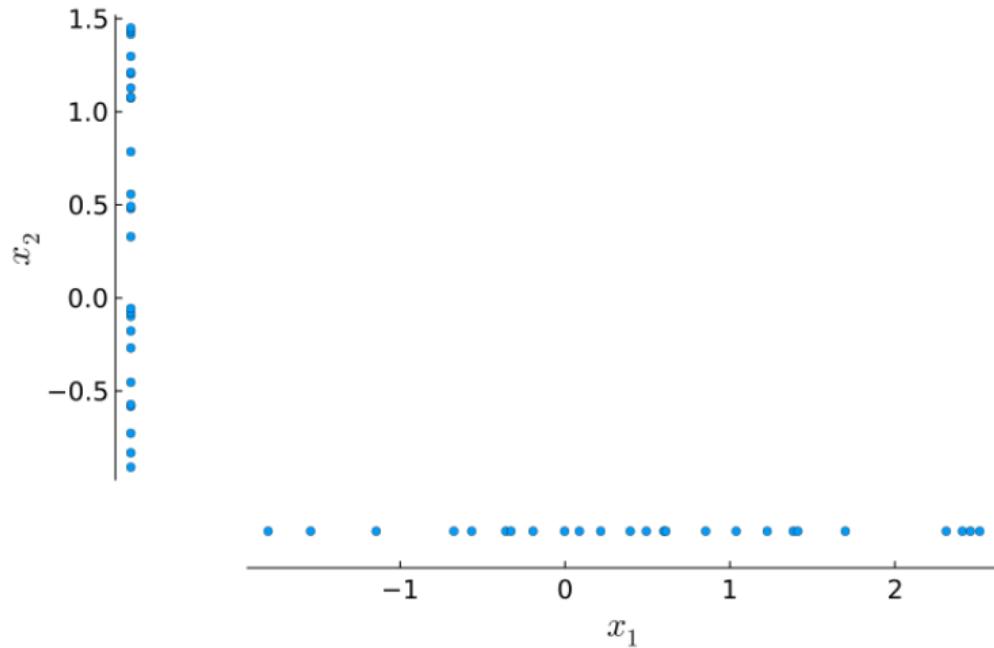
The naïve $\mathcal{O}(n^3)$ computational bottleneck ($\mathcal{O}(n^2)$ memory) can be tackled by

- **Exploiting structure in the data**
(data on grid, inputs are in 1D, ...)
- **Exploiting structure in the GP prior**
(GP prior is stationary, separable over input dimensions, ...)
- **Solving the linear system approximately**
(conjugate-gradient solvers)
- **Split problem into smaller chunks**
(local experts, subset of data, divide & conquer, ...)
- **Approximate the problem**
(Nyström, low-rank, inducing points, ...)
- **Approximate the problem solution**
(SVGP = sparse (and stochastic) variational methods)

Section 2

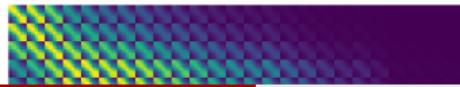
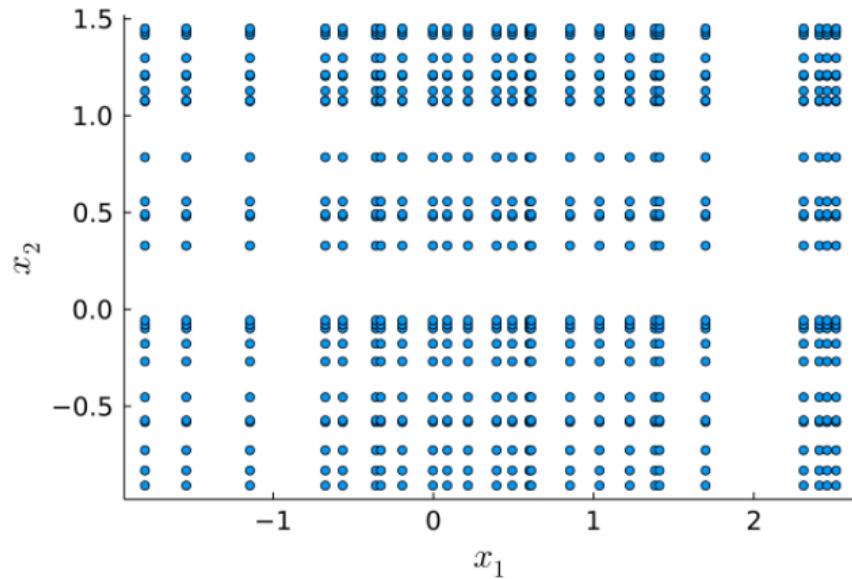
Exploiting structure in data and/or GP prior

Exploiting structure: Inputs on grid



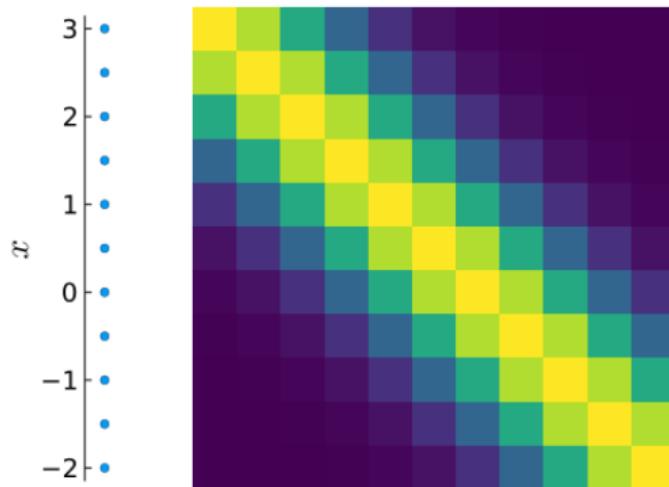
Exploiting structure: Inputs on grid

and *separable* kernel: $k(\mathbf{x}, \mathbf{x}') = \prod_{d=1}^D k_d(x_d, x'_d)$



Exploiting structure: Inputs on *regular* grid

and *stationary* kernel: $k(x, x') = k(x - x')$



Toeplitz structure: $K_{i,j} = K_{i+1,j+1} = k(x_i - x_j)$

efficient mvm via FFT: computation $\mathcal{O}(N + M \log M)$, storage $\mathcal{O}(N + M)$.

What if inputs are not exactly on grid?

Interpolate (e.g., cubically) kernel matrix between grid points: “KISS-GP”

- Can use very fine grids: very accurate
- Works well in small dimensions (up to 4)

What about higher dimensions?

More structure:

- Additive kernel
- Product kernel

Time series, graphs: sparse precision matrix, state space formulation

What about arbitrary kernels?

Section 3

Solving the linear system approximately (with Conjugate Gradients)

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
- 6 Variational inference for sparse GPs
- 7 Recap

Matrix inverse as quadratic optimization

Rewrite matrix inverse

$$\mathbf{v} = \hat{\mathbf{K}}^{-1} \mathbf{y}, \quad \hat{\mathbf{K}} = \mathbf{K} + \sigma^2 \mathbf{I}$$

as linear system:

$$\hat{\mathbf{K}}\mathbf{v} - \mathbf{y} = 0$$

Solve as quadratic optimization problem:

$$\mathbf{v}^* = \arg \min_{\mathbf{v}} \mathbf{v}^\top \hat{\mathbf{K}} \mathbf{v} - \mathbf{v}^\top \mathbf{y}$$

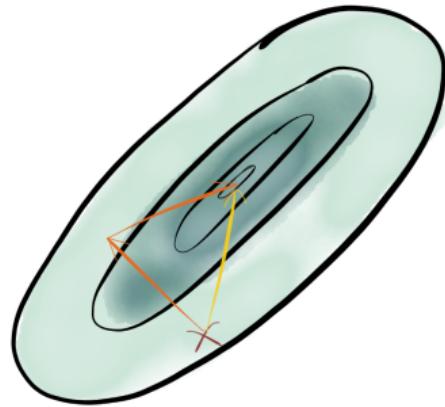
Conjugate Gradients

Efficiently solve

$$\mathbf{v}^* = \arg \min_{\mathbf{v}} \mathbf{v}^\top \hat{K} \mathbf{v} - \mathbf{v}^\top \mathbf{y}$$

using Conjugate Gradients (CG)

- Iterative method
- Each step is $\mathcal{O}(N^2)$
- Recovers exact solution after N steps $\rightarrow \mathcal{O}(N^3)$
- *Approximate* solution in much fewer steps!
("easier" problems: less steps)



Conjugate Gradient (CG)

Figure taken from Davies (2015)

Convergence and preconditioning

Condition number: ratio of largest to smallest eigenvalue,

$$\kappa(\hat{K}) = \frac{\lambda_{\max}(\hat{K})}{\lambda_{\min}(\hat{K})}$$

High condition numbers: numerically unstable, slow convergence

Improve by preconditioning: Instead of $\hat{K}\mathbf{v} - \mathbf{y} = 0$, solve

$$P^{-1}\hat{K}\mathbf{v} - P^{-1}\mathbf{y} = 0$$

If $P^{-1} = \hat{K}^{-1}$, then $\kappa(P^{-1}\hat{K}) = 1$ and we solve in one step.

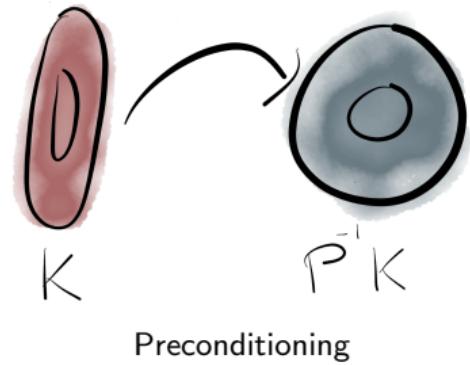
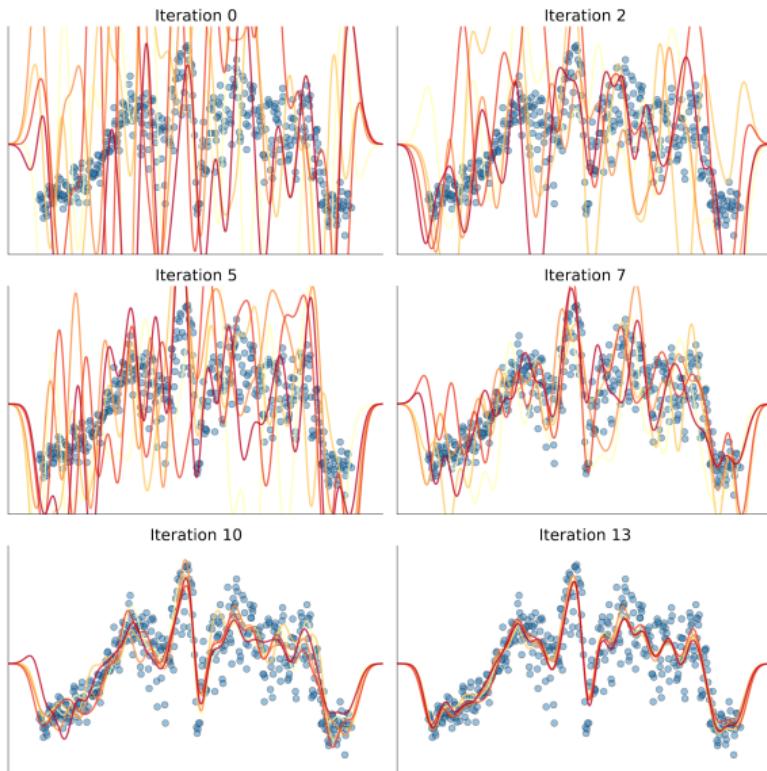


Figure taken from Davies (2015)

CG example



Section 4

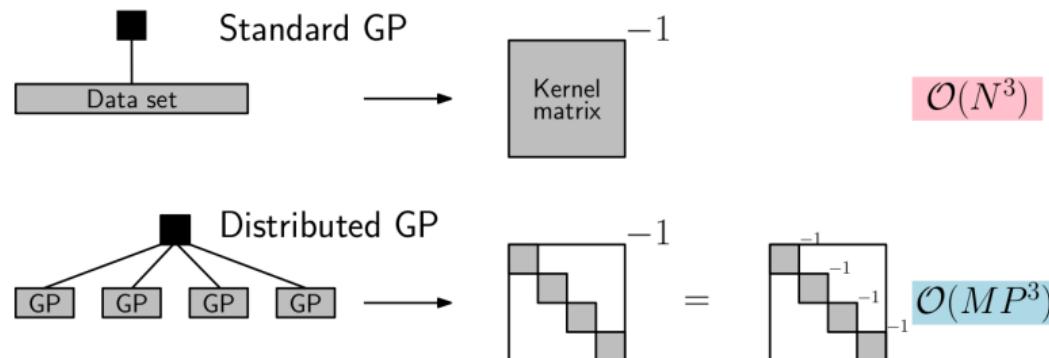
Local approximations: Split problem into smaller chunks

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
- 6 Variational inference for sparse GPs
- 7 Recap

Combining “local experts”

- Split input domain/dataset (size N) into M subsets (size $P \approx N/M$)
- Fit M independent GP models (shared kernel hyperparameters): $\mathcal{O}(P^3)$ each
- Aggregate predictions
 - + Mixture of experts
 - ✗ product of experts, e.g.

$$p(f_* | \mathbf{x}_*, \mathcal{D}) \propto \prod_{k=1}^M p_k(f_* | \mathbf{x}_*, \mathcal{D}^{(k)})$$



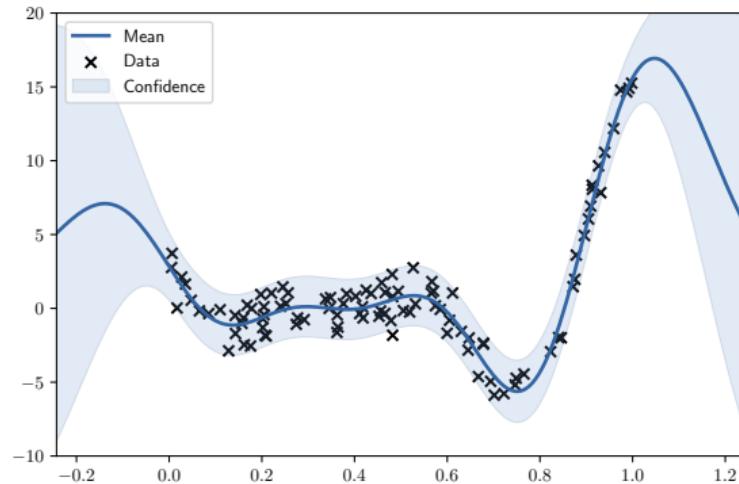
Section 5

Global approximations: Inducing points

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
 - Data redundancy
 - Nyström approximation
 - Pseudo data / inducing points
- 6 Variational inference for sparse GPs
- 7 Recap

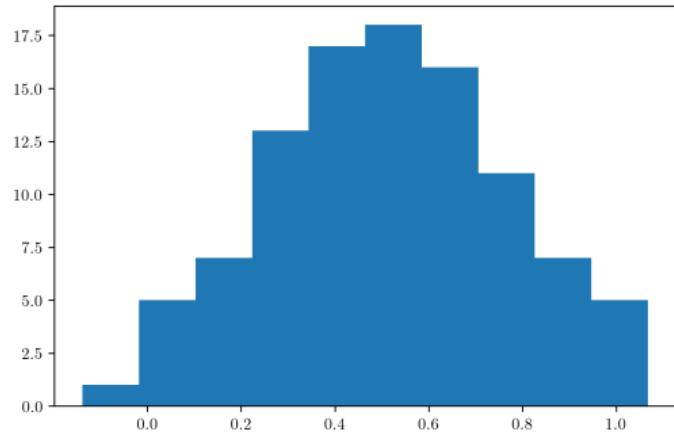
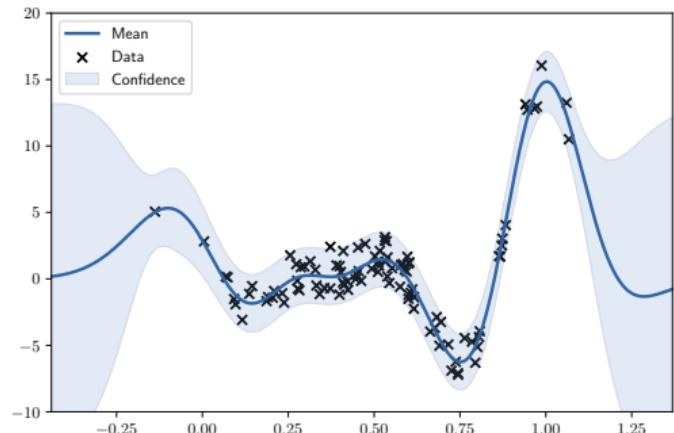
Big data (?)

- Lots of data \neq complex function
- In real world problems, we often collect a lot of data for modeling relatively simple relations.



Data subsampling?

- Real data generally not evenly distributed.
- Typically a lot of data on common cases, very few data on rare cases.



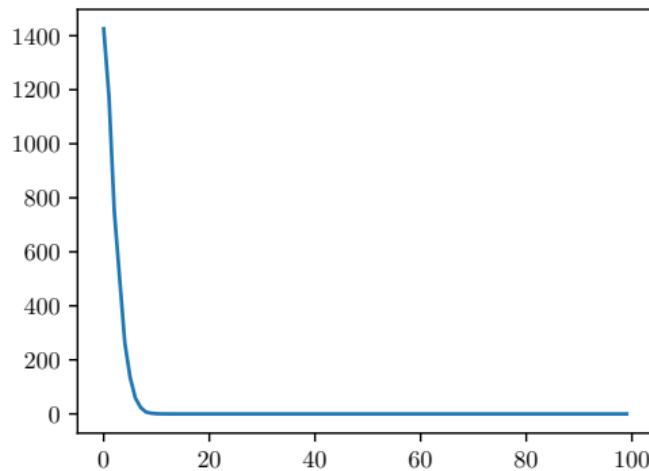
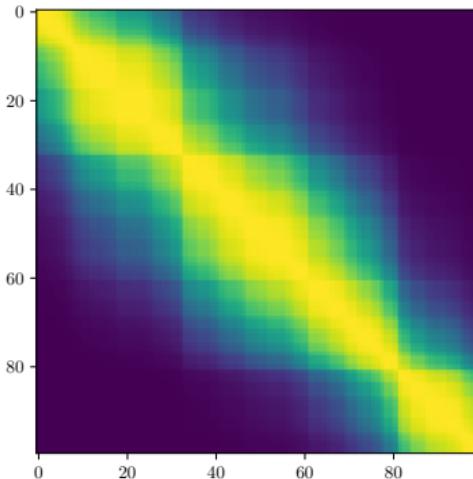
Sparse Gaussian processes

“Sparse GP” refers to various different approximations:

- Nyström approximation
- Fully independent training conditional (FITC)
- Variational sparse Gaussian process

Covariance matrix of redundant data

- With redundant data, the covariance matrix becomes low rank.



- What about low rank approximation?

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
 - Data redundancy
 - Nyström approximation
 - Pseudo data / inducing points
- 6 Variational inference for sparse GPs
- 7 Recap

Low-rank approximation

- Recall GP marginal log-likelihood:

$$\log p(\mathbf{y} \mid \mathbf{X}) = \log \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I}),$$

where \mathbf{K} is the covariance matrix computed from \mathbf{X} according to the kernel function $k(\cdot, \cdot)$ and σ^2 is the variance of the Gaussian noise distribution.

- Assume \mathbf{K} to be low rank.
- This leads to the Nyström approximation by Williams and Seeger

Approximation by subset

- Let's randomly pick a subset from the training data: $\mathbf{Z} \in \mathbb{R}^{M \times Q}$.
- Approximate the covariance matrix \mathbf{K} by $\tilde{\mathbf{K}}$.

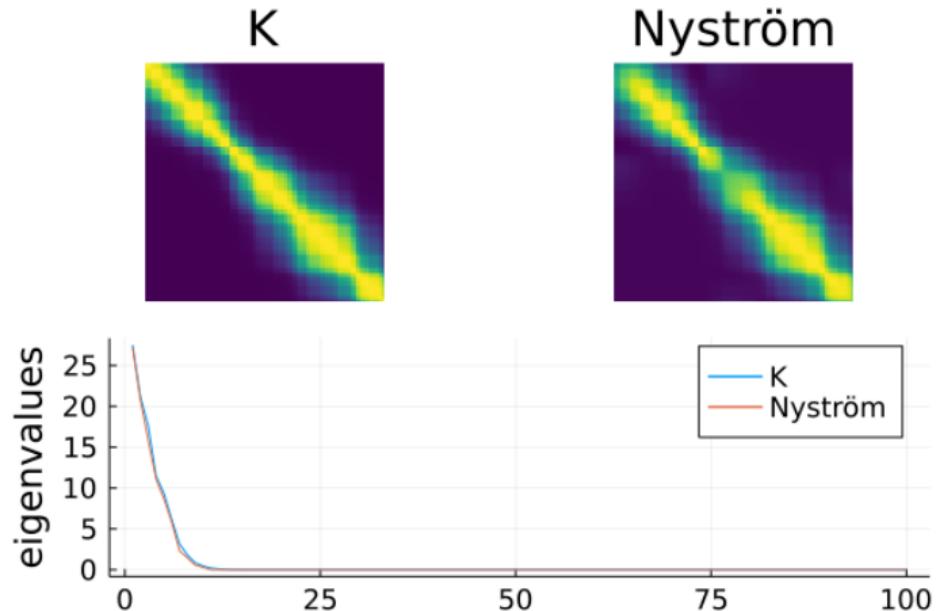
$$\tilde{\mathbf{K}} = \mathbf{K}_z \mathbf{K}_{zz}^{-1} \mathbf{K}_z^\top, \text{ where } \mathbf{K}_z = \mathbf{K}(\mathbf{X}, \mathbf{Z}) \text{ and } \mathbf{K}_{zz} = \mathbf{K}(\mathbf{Z}, \mathbf{Z}).$$

- Note that $\tilde{\mathbf{K}} \in \mathbb{R}^{N \times N}$, $\mathbf{K}_z \in \mathbb{R}^{N \times M}$ and $\mathbf{K}_{zz} \in \mathbb{R}^{M \times M}$.
- The log-likelihood is approximated by

$$\log p(\mathbf{y} \mid \mathbf{X}, \theta) \approx \log \mathcal{N} \left(\mathbf{y} \mid \mathbf{0}, \mathbf{K}_z \mathbf{K}_{zz}^{-1} \mathbf{K}_z^\top + \sigma^2 \mathbf{I} \right).$$

Nyström approximation example

Nyström approximation using 10 random data points:



Efficient computation using Woodbury formula

Naive formulation: still $\mathcal{O}(N^3)$!

$$\tilde{\mathcal{L}} = -\frac{1}{2} \log |2\pi(\tilde{\mathbf{K}} + \sigma^2 \mathbf{I})| - \frac{1}{2} \mathbf{y}^\top (\tilde{\mathbf{K}} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

Apply the Woodbury matrix identity:

$$\begin{aligned} (UC & V & +A)^{-1} &= A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1} \\ (\mathbf{K}_z \mathbf{K}_{zz}^{-1} \mathbf{K}_z^\top + \sigma^2 \mathbf{I})^{-1} &= \sigma^{-2} \mathbf{I} - \sigma^{-4} \mathbf{K}_z (\mathbf{K}_{zz} + \sigma^{-2} \mathbf{K}_z^\top \mathbf{K}_z)^{-1} \mathbf{K}_z^\top \end{aligned}$$

- Note that $(\mathbf{K}_{zz} + \sigma^{-2} \mathbf{K}_z^\top \mathbf{K}_z) \in \mathbb{R}^{M \times M}$.
- The computational complexity reduces to $O(NM^2)$.

Nyström approximation

- Approximation directly on the covariance matrix
- Randomly selected subset of data points
- Approximation becomes exact if the whole data set is taken, *i.e.*, $\mathbf{K}\mathbf{K}^{-1}\mathbf{K}^\top = \mathbf{K}$.

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
 - Data redundancy
 - Nyström approximation
 - Pseudo data / inducing points
- 6 Variational inference for sparse GPs
- 7 Recap

Inducing point methods

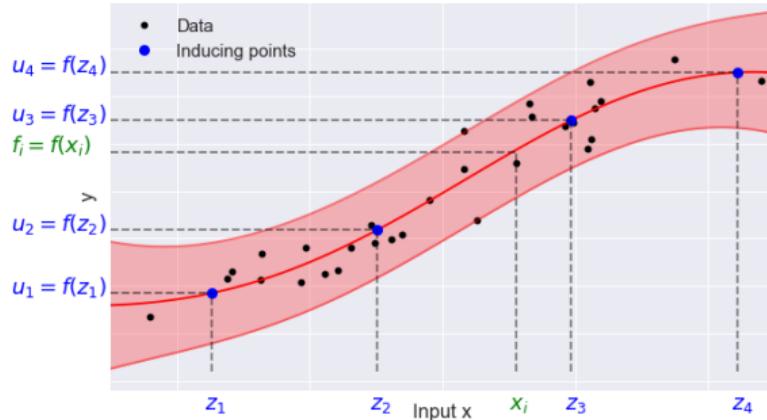
Main idea: “represent” the information from the full dataset using a smaller “virtual” dataset

- Recall our GP model:

$$p(\mathbf{y}, \mathbf{f}) = p(\mathbf{y} \mid \mathbf{f})p(\mathbf{f}), \quad \text{where } \mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)]$$

- We will now introduce a set of *inducing points* $\{\mathbf{z}_m\}_{m=1}^M$
- They live in the same space as the input points, i.e. $\mathbf{x}_i, \mathbf{z}_j \in \mathbb{R}^D$
- Let u_m denote the value of the function f evaluated at each \mathbf{z}_m , i.e. $u_m = f(\mathbf{z}_m)$
- ... and $\mathbf{u} = [f(\mathbf{z}_1), f(\mathbf{z}_2), \dots, f(\mathbf{z}_M)]$

Inducing point methods



- Goal: choose set of inducing points such that it contains same information as full dataset
- Remember: Both $u_j = f(z_j)$ and $f_i = f(x_i)$ are random variables
- How can we learn \mathbf{u} ?
Next step: Formulate joint model $p(\mathbf{y}, \mathbf{f}, \mathbf{u})$

Section 6

Variational inference for sparse GPs

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
- 6 Variational inference for sparse GPs
 - Recap: Variational inference
 - Variational inference for inducing variables
 - Collapsed bound for Gaussian likelihood
 - Mini-batch learning
 - Mini-batch learning for GPs
 - Numerical considerations

Variational inference: The big picture

Recipe for approximating intractable distribution $p \in \mathcal{P}$

- ① Define some “simple” family of distributions \mathcal{Q} .
- ② Define some way to compute a “distance” $\mathbb{D}[p, q]$ between intractable distribution p and each distribution $q \in \mathcal{Q}$

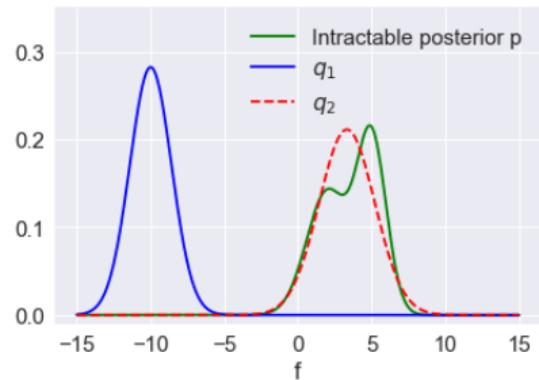
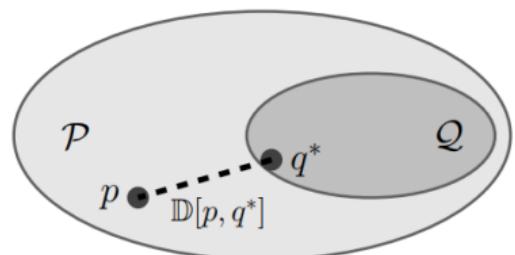
$$\mathbb{D}[p, q_1] > \mathbb{D}[p, q_2]$$

- ③ Search for $q \in \mathcal{Q}$ such that $\mathbb{D}[p, q]$ is minimized

$$q^* = \arg \min_{q \in \mathcal{Q}} \mathbb{D}[p, q]$$

- ④ Use q^* as an approximation of p

Here we will always choose \mathcal{Q} to be the set of multivariate Gaussian distributions.



How to “measure distances” between distributions?

Here: *Kullback-Leibler divergence*

$$\mathbb{D}[p, q] := \text{KL}[q \parallel p] = \int q(\mathbf{f}) \log \frac{q(\mathbf{f})}{p(\mathbf{f})} d\mathbf{f} = \mathbb{E}_q \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f})} \right]$$

Important properties:

- ① Non-symmetric: $\text{KL}[q \parallel p] \neq \text{KL}[p \parallel q]$
- ② Positive: $\text{KL} \geq 0$ (Gibbs' inequality)
- ③ Minimum: $\text{KL}[q \parallel p] = 0 \iff q \equiv p$.

Variational inference: Minimizing $\text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f} \mid \mathbf{y})]$ by bounding

$$\begin{aligned}\text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f} \mid \mathbf{y})] &= \text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f})] - \int q(\mathbf{f}) [\log p(\mathbf{y} \mid \mathbf{f})] d\mathbf{f} + \log p(\mathbf{y}) \\ \log p(\mathbf{y}) &= \underbrace{\int q(\mathbf{f}) [\log p(\mathbf{y} \mid \mathbf{f})] d\mathbf{f} - \text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f})]}_{\mathcal{L}[q]} + \underbrace{\text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f} \mid \mathbf{y})]}_{\geq 0} \\ &\geq \int q(\mathbf{f}) [\log p(\mathbf{y} \mid \mathbf{f})] d\mathbf{f} - \text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f})] = \mathcal{L}[q]\end{aligned}$$

- $\log p(\mathbf{y})$ is a constant
- $\mathcal{L}[q]$ does not depend on $p(\mathbf{f} \mid \mathbf{y})$
- $\mathcal{L}[q] \leq \log p(\mathbf{y})$, so $\mathcal{L}[q]$ is a *lower bound* on the marginal log likelihood
- Maximizing $\mathcal{L}[q]$ is equivalent to minimizing $\text{KL}[q(\mathbf{f}) \parallel p(\mathbf{f} \mid \mathbf{y})]$

Key take-away: we can fit variational approximation q by optimizing \mathcal{L}

Variational inference: ELBO

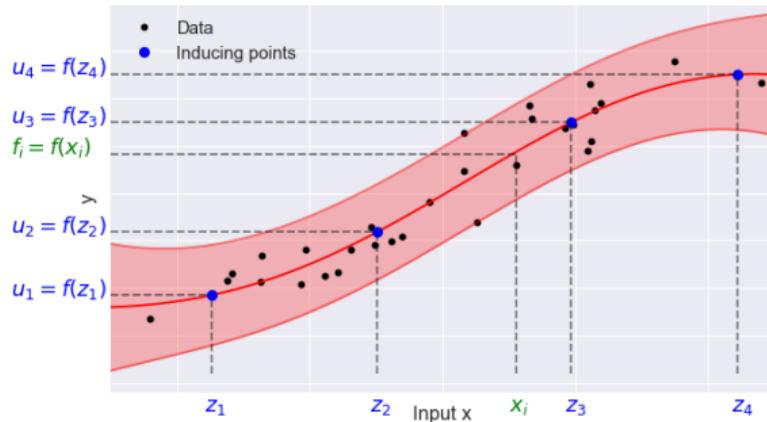
$$\log p(\mathbf{y}) \geq \mathcal{L}[q] = \underbrace{\int q(\mathbf{f}) [\log p(\mathbf{y} | \mathbf{f})] d\mathbf{f}}_{\text{data fit}} - \underbrace{\text{KL}[q(\mathbf{f}) \| p(\mathbf{f})]}_{\text{regularization}}$$

$\mathcal{L}[q]$ often called the *Evidence Lower Bound* (ELBO)

- To approximate $p(\mathbf{f} | \mathbf{y})$, use $q(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{m}, \mathbf{S})$
- Define $\boldsymbol{\lambda} = \{\mathbf{m}, \mathbf{S}\}$, then we can write $\mathcal{L}[q] = \mathcal{L}[\boldsymbol{\lambda}]$
- In practice, we optimize $\mathcal{L}[\boldsymbol{\lambda}]$ using gradient-based methods

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
- 6 Variational inference for sparse GPs
 - Recap: Variational inference
 - Variational inference for inducing variables
 - Collapsed bound for Gaussian likelihood
 - Mini-batch learning
 - Mini-batch learning for GPs
 - Numerical considerations

Joint model over function values and inducing variables



- Goal: choose set of inducing points such that it contains same information as full dataset
- Remember: Both $u_j = f(z_j)$ and $f_i = f(x_i)$ are random variables
- Before: $p(\mathbf{y}, \mathbf{f}) = p(\mathbf{y} | \mathbf{f})p(\mathbf{f})$.
Now: joint model $p(\mathbf{y}, \mathbf{f}, \mathbf{u}) = p(\mathbf{y} | \mathbf{f})p(\mathbf{f}, \mathbf{u})$

Inducing point methods: the joint model

- The augmented model

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

- We recover the original model by marginalizing over \mathbf{u} :

$$p(\mathbf{y}, \mathbf{f}) = \int p(\mathbf{y} \mid \mathbf{f})p(\mathbf{f}, \mathbf{u}) \, d\mathbf{u} = p(\mathbf{y} \mid \mathbf{f}) \int p(\mathbf{f}, \mathbf{u}) \, d\mathbf{u} = p(\mathbf{y} \mid \mathbf{f})p(\mathbf{f})$$

- Let's decompose the “augmented” model as follows

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

Inducing point methods: the joint model II

- The joint model

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

- Using Gaussian conditional densities (lecture #1):

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

$$p(\mathbf{f} \mid \mathbf{u}) = \mathcal{N}\left(\mathbf{f} \mid \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}\mathbf{u}, \tilde{\mathbf{K}}\right), \quad \tilde{\mathbf{K}} = \mathbf{K}_{nn} - \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}\mathbf{K}_{mn}$$

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

- Covariance of inducing points: $[\mathbf{K}_{mm}]_{ij} = k(\mathbf{z}_i, \mathbf{z}_j)$
- Cross-covariance between inducing points and training: $[\mathbf{K}_{mn}]_{ij} = k(\mathbf{z}_i, \mathbf{x}_j)$
- Covariance of training points: $[\mathbf{K}_{nn}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$

Variational Sparse Gaussian Process

- Typical variational lower bound of a marginal likelihood:

$$\begin{aligned}\log p(\mathbf{y} \mid \mathbf{X}) &= \log \int_{\mathbf{f}, \mathbf{u}} p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) p(\mathbf{u} \mid \mathbf{Z}) \\ &\geq \int_{\mathbf{f}, \mathbf{u}} q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) p(\mathbf{u} \mid \mathbf{Z})}{q(\mathbf{f}, \mathbf{u})} \equiv \mathcal{L}\end{aligned}$$

- Let's define a *special* variational posterior (make use of structure in our problem):

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}), \quad \text{where } q(\mathbf{u}) = \mathcal{N}(\mathbf{u} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

- Plug it into the lower bound:

$$\begin{aligned}\mathcal{L} &= \int_{\mathbf{f}, \mathbf{u}} p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}) \log \frac{p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) p(\mathbf{u} \mid \mathbf{Z})}{p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u})} \\ &= \langle \log p(\mathbf{y} \mid \mathbf{f}) \rangle_{p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u})} - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u} \mid \mathbf{Z})]\end{aligned}$$

Likelihood term

- Remember: $p(\mathbf{y} \mid \mathbf{f}) = \prod_{i=1}^N p(y_i \mid f_i)$
- Let's have a closer look at the first term

$$\begin{aligned}\mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\log p(\mathbf{y} \mid \mathbf{f})] &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[\log \prod_{i=1}^N p(y_i \mid f_i) \right] = \sum_{i=1}^N \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\log p(y_i \mid f_i)] \\ &= \sum_{i=1}^N \iint q(\mathbf{u}, \mathbf{f}) \log p(y_i \mid f_i) \mathrm{d}\mathbf{u} \mathrm{d}\mathbf{f} \\ &= \sum_{i=1}^N \iint p(f_i \mid \mathbf{u}) \mathcal{N}(\mathbf{u} \mid \mathbf{m}, \mathbf{S}) \log p(y_i \mid f_i) \mathrm{d}\mathbf{u} \mathrm{d}f_i \\ &= \sum_{i=1}^N \iint p(f_i \mid \mathbf{u}) \mathcal{N}(\mathbf{u} \mid \mathbf{m}, \mathbf{S}) \mathrm{d}\mathbf{u} \log p(y_i \mid f_i) \mathrm{d}f_i\end{aligned}$$

Decomposing the likelihood term

- Let's define the univariate distribution

$$q(f_i) \equiv \int p(f_i | \mathbf{u}) \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}) d\mathbf{u} = \mathcal{N}\left(f_i | \mathbf{k}_{im} \mathbf{K}_{mm}^{-1} \mathbf{m}, \tilde{K}_{ii} + \mathbf{k}_{im} \mathbf{K}_{mm}^{-1} \mathbf{S} \mathbf{K}_{mm}^{-1} \mathbf{k}_{mi}\right)$$

- Variational expectations term:

$$\begin{aligned}\mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\log p(\mathbf{y} | \mathbf{f})] &= \sum_{i=1}^N \iint p(f_i | \mathbf{u}) \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}) d\mathbf{u} \log p(y_i | f_i) df_i \\ &= \sum_{i=1}^N \int q(f_i) \log p(y_i | f_i) df_i\end{aligned}$$

- As in the previous lecture:
 - Decomposes into a sum over 1D integrals
 - Solved analytically for some likelihoods
 - Fast to approximate using numerical integration for others

Sparse variational posterior for **Gaussian** likelihood

- Variational posterior:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z})q(\mathbf{u}), \quad \text{where } q(\mathbf{u}) = \mathcal{N}(\mathbf{u} \mid \mu, \Sigma).$$

- Lower bound:

$$\begin{aligned}\mathcal{L} &= \int_{\mathbf{f}, \mathbf{u}} p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z})q(\mathbf{u}) \log \frac{p(\mathbf{y} \mid \mathbf{f})p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z})p(\mathbf{u} \mid \mathbf{Z})}{p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z})q(\mathbf{u})} \\ &= \langle \log p(\mathbf{y} \mid \mathbf{f}) \rangle_{p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z})q(\mathbf{u})} - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u} \mid \mathbf{Z})] \\ &= \langle \log \mathcal{N}(\mathbf{y} \mid \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u}, \sigma^2\mathbf{I}) \rangle_{q(\mathbf{u})} - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u} \mid \mathbf{Z})]\end{aligned}$$

- There is no inversion of any big covariance matrices in the first term:

$$-\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \left\langle (\mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u} - \mathbf{y})^\top (\mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u} - \mathbf{y}) \right\rangle_{q(\mathbf{u})}$$

- The overall complexity of the lower bound is $O(NM^2)$.

Tighten the Bound

- Find the optimal parameters of $q(\mathbf{u})$:

$$(\mathbf{m}^*, \mathbf{S}^*) = \arg \max_{\mathbf{m}, \mathbf{S}} \mathcal{L}(\mathbf{m}, \mathbf{S}).$$

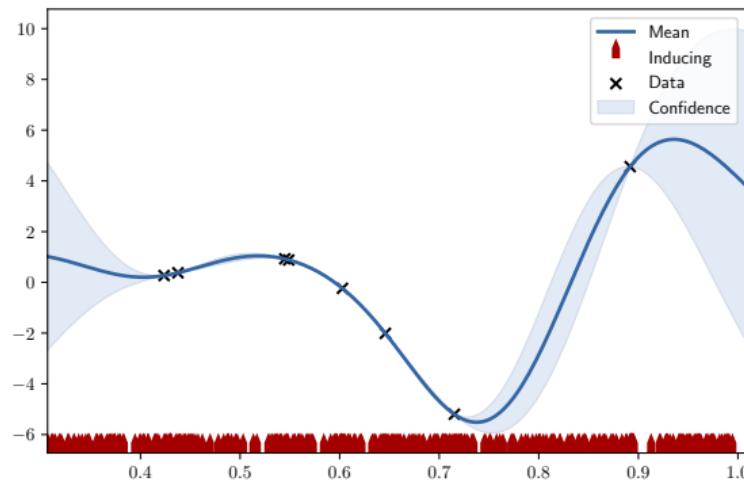
- Make the bound as tight as possible by plugging in \mathbf{m}^* and \mathbf{S}^* :

$$\mathcal{L} = \log \mathcal{N} \left(\mathbf{y} \mid \mathbf{0}, \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \mathbf{K}_{fu}^\top + \sigma^2 \mathbf{I} \right) - \frac{1}{2\sigma^2} \text{tr} \left(\mathbf{K}_{ff} - \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \mathbf{K}_{fu}^\top \right).$$

- The 1st term is the same as in the Nyström approximation.
- The overall complexity of the lower bound remains $O(NM^2)$.

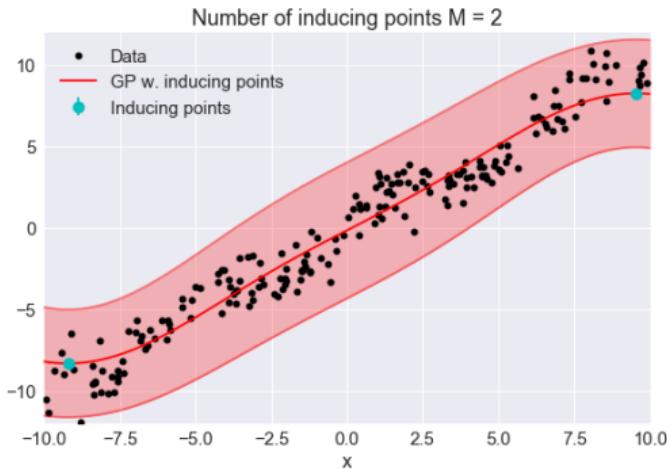
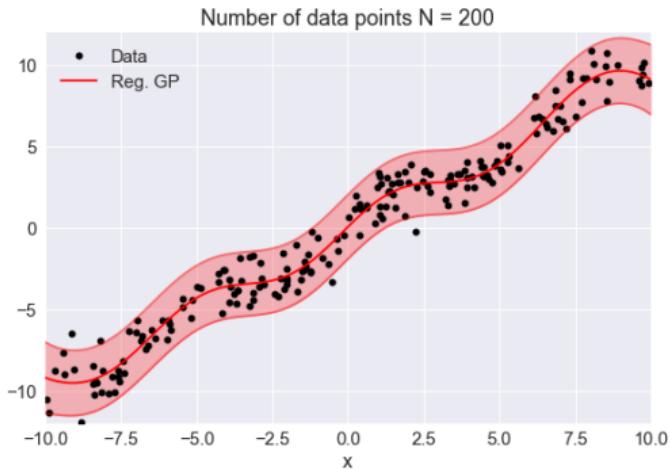
Variational sparse GP

- Note that \mathcal{L} is not a valid log-pdf, $\int_{\mathbf{y}} \exp(\mathcal{L}(\mathbf{y})) \leq 1$, due to the trace term.
- As inducing points are variational parameters, optimizing the inducing inputs \mathbf{Z} always leads to a better bound.
- The model does not “overfit” with too many inducing points.



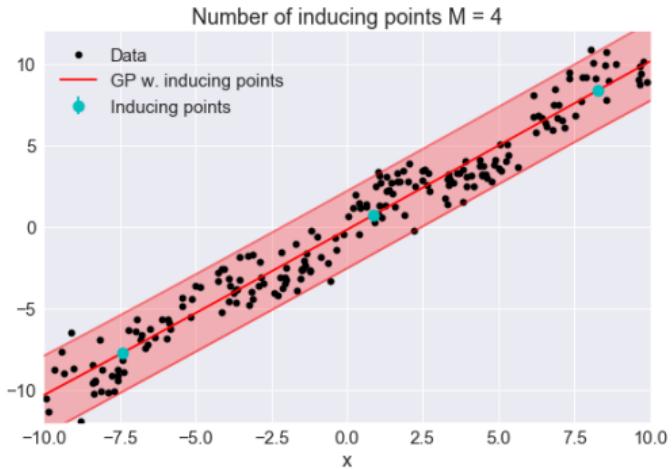
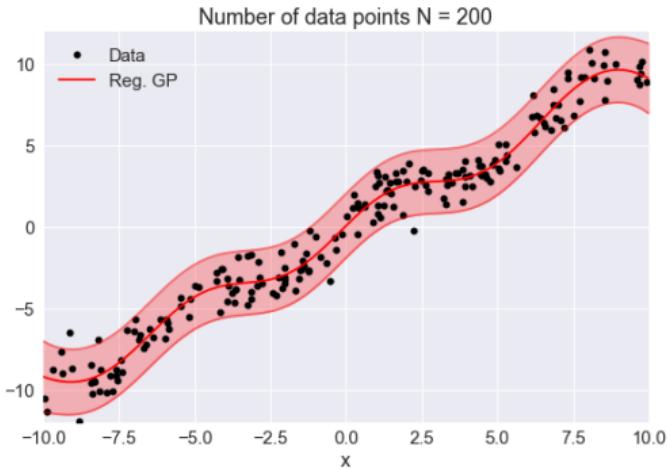
Example: Number of inducing points

We can think of number of inducing points as parameter that trades off speed for accuracy



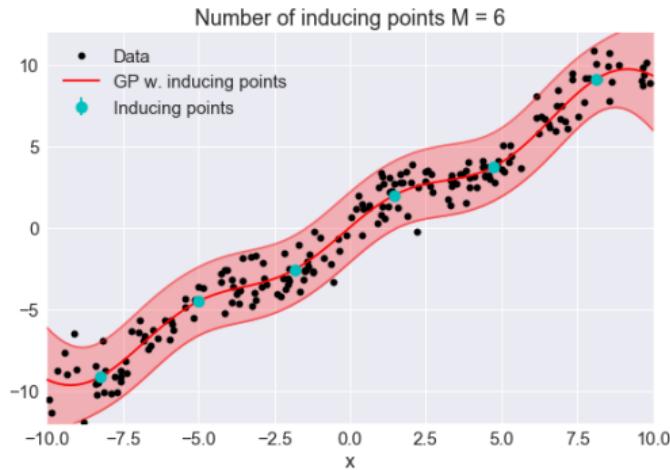
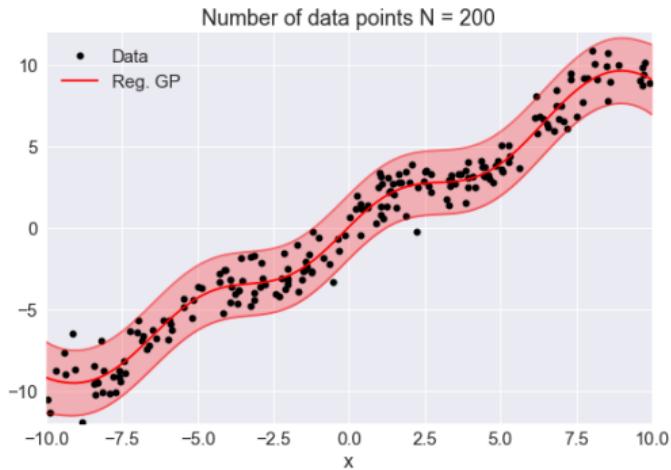
Example: Number of inducing points

We can think of number of inducing points as parameter that trades off speed for accuracy



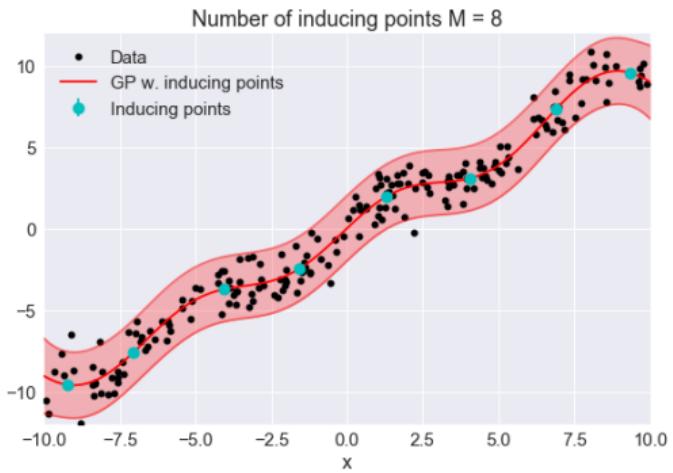
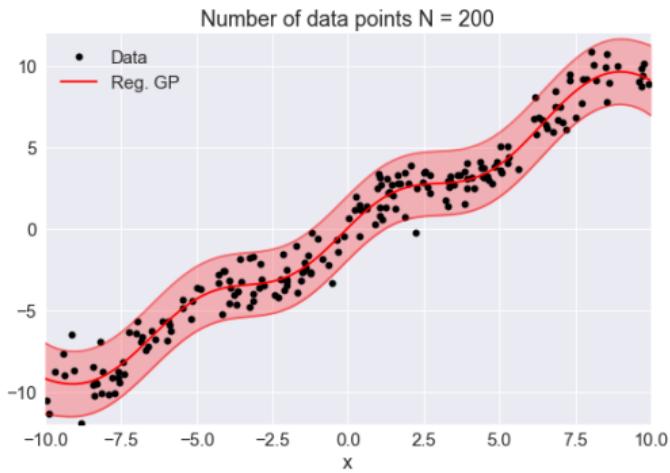
Example: Number of inducing points

We can think of number of inducing points as parameter that trades off speed for accuracy



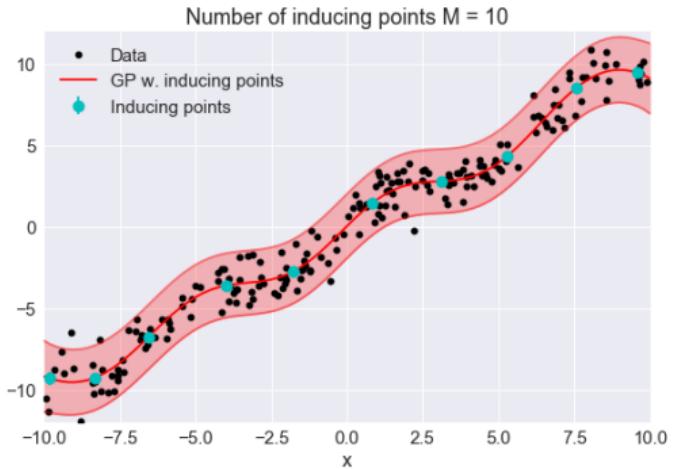
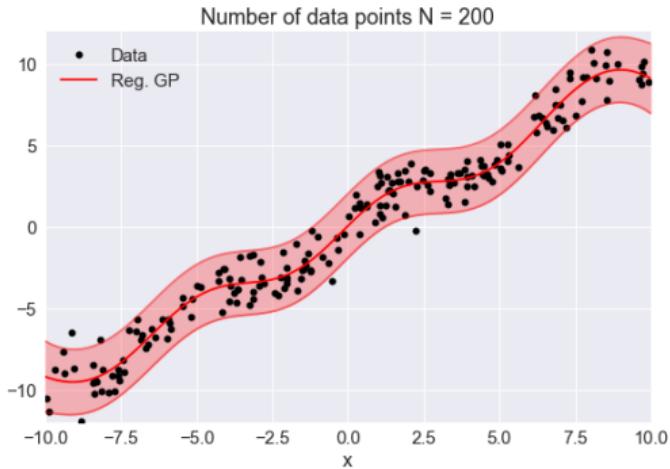
Example: Number of inducing points

We can think of number of inducing points as parameter that trades off speed for accuracy



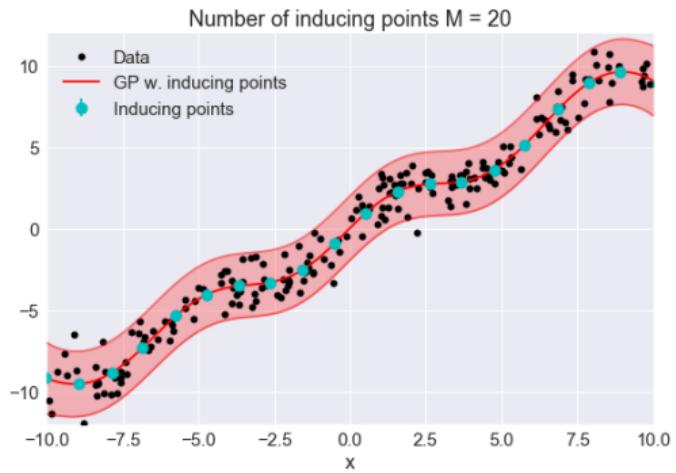
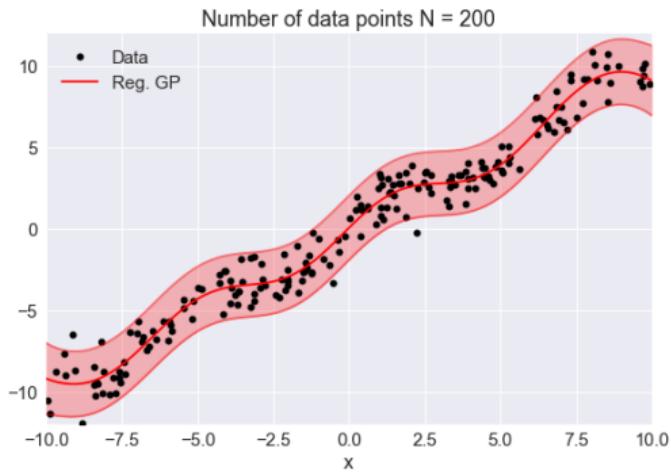
Example: Number of inducing points

We can think of number of inducing points as parameter that trades off speed for accuracy



Example: Number of inducing points

We can think of number of inducing points as parameter that trades off speed for accuracy



Limitations of sparse GPs

Variational sparse GP has computational complexity $O(NM^2)$.

The computation becomes infeasible under two scenarios:

- The number of data points N is very high, e.g., millions of data points.
- The function is very complex, which requires tens of thousands of inducing points.

- 1 What is the computational issue?
- 2 Exploiting structure in data and/or GP prior
- 3 Solving the linear system approximately (with Conjugate Gradients)
- 4 Local approximations: Split problem into smaller chunks
- 5 Global approximations: Inducing points
- 6 Variational inference for sparse GPs
 - Recap: Variational inference
 - Variational inference for inducing variables
 - Collapsed bound for Gaussian likelihood
 - Mini-batch learning
 - Mini-batch learning for GPs
 - Numerical considerations

Mini-batch learning (1)

- Mini-batch learning allows DNNs to be trained on millions of data points.
- Given a set of inputs and labels, $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$, $(\mathbf{x}_i, y_i) \sim p(\mathbf{x}, y)$, the true loss function is defined as

$$c_{\text{true}} = \int l(f_\theta(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy \approx \frac{1}{N} \sum_{i=1}^N l(f_\theta(\mathbf{x}_i), y_i) = c,$$

where $f_\theta(\cdot)$ is DNN and $l(\cdot, \cdot)$ is the loss function.

- Gradient descent (GD) updates the parameters by

$$\theta_{t+1} = \theta_t - \eta \frac{dc}{d\theta}.$$

Mini-batch learning (2)

- Mini-batch learning approximates the loss by subsampling the data,

$$c_{\text{MB}} = \frac{1}{B} \sum_{\mathbf{x}_i, y_i \sim \tilde{p}(\mathbf{x}, y)} l(f_\theta(\mathbf{x}_i), y_i).$$

- Stochastic gradient descent (SGD) updates the parameters by

$$\theta_{t+1} = \theta_t - \eta \frac{dc_{\text{MB}}}{d\theta}.$$

- Can mini-batch learning be applied to GPs as well?

Mini-batch Learning for GPs

- Mini-batch learning relies on the objective being an expectation w.r.t. the data, i.e.,
 $\langle l(f_\theta(\mathbf{x}), y) \rangle_{p(\mathbf{x}, y)}$.
- The log-marginal likelihood of GP:

$$\log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$$

- The variational lower bound of sparse GP:

$$\log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \mathbf{K}_{fu}^\top + \sigma^2 \mathbf{I}) - \frac{1}{2\sigma^2} \text{tr}(\mathbf{K}_{ff} - \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \mathbf{K}_{fu}^\top)$$

- Not an expectation w.r.t. the data!

“Uncollapsed” lower bound

- Let's revisit original (“uncollapsed”) variational lower bound of sparse GP:

$$\mathcal{L} = \langle \log p(\mathbf{y} \mid \mathbf{f}) \rangle_{p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z})q(\mathbf{u})} - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u})]$$

- The 2nd term, $\text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u})]$, does not depend on the data.

“Uncollapsed” lower bound

- In the 1st term, as $p(\mathbf{y} \mid \mathbf{f}) = \mathcal{N}(\mathbf{y} \mid \mathbf{f}, \sigma^2 \mathbf{I})$,

$$\log p(\mathbf{y} \mid \mathbf{f}) = \sum_{n=1}^N \log \mathcal{N}(y_n \mid f_n, \sigma^2)$$

- Denote $q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}) = \int p(\mathbf{f} \mid \mathbf{u}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}) d\mathbf{u}$.

$$\begin{aligned}\langle \log p(\mathbf{y} \mid \mathbf{f}) \rangle_{q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z})} &= \left\langle \sum_{n=1}^N \log \mathcal{N}(y_n \mid f_n, \sigma^2) \right\rangle_{q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z})} \\ &= \sum_{n=1}^N \langle \log \mathcal{N}(y_n \mid f_n, \sigma^2) \rangle_{q(f_n \mid \mathbf{x}_n, \mathbf{Z})}\end{aligned}$$

Stochastic Variational GP (SVGP)

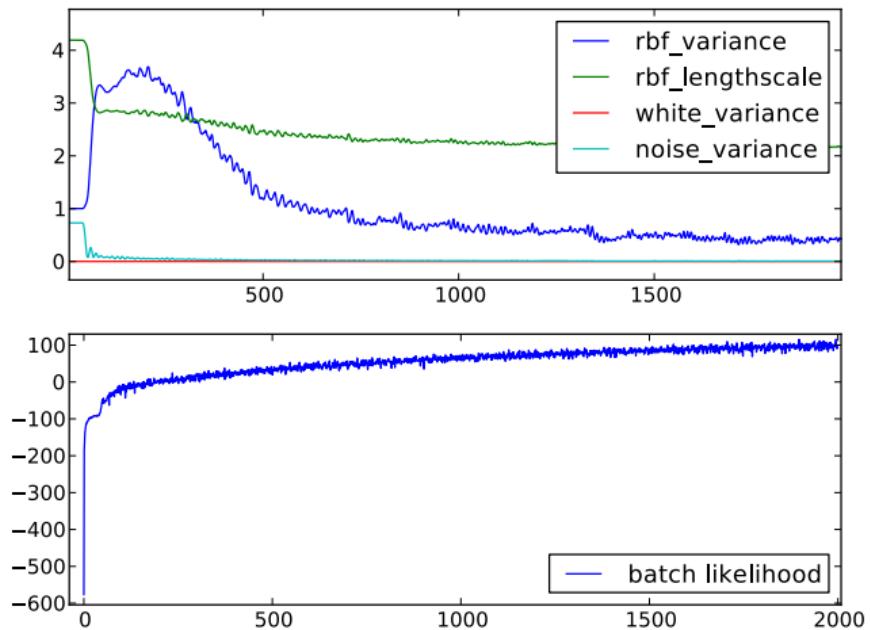
- The resulting lower bound can be written as the sum over the data,

$$\begin{aligned}\mathcal{L} &= \sum_{n=1}^N \langle \log \mathcal{N}(y_n | f_n, \sigma^2) \rangle_{q(f_n | \mathbf{x}_n, \mathbf{z})} - \text{KL}[q(\mathbf{u}) \| p(\mathbf{u})] \\ &\approx \frac{N}{B} \sum_{\mathbf{x}_i, y_i \sim \tilde{p}(\mathbf{x}, y)} \langle \log \mathcal{N}(y_i | f_i, \sigma^2) \rangle_{q(f_i | \mathbf{x}_i, \mathbf{z})} - \text{KL}[q(\mathbf{u}) \| p(\mathbf{u})] = \mathcal{L}_{\text{MB}}\end{aligned}$$

- This allows us to do mini-batch learning with SGD,

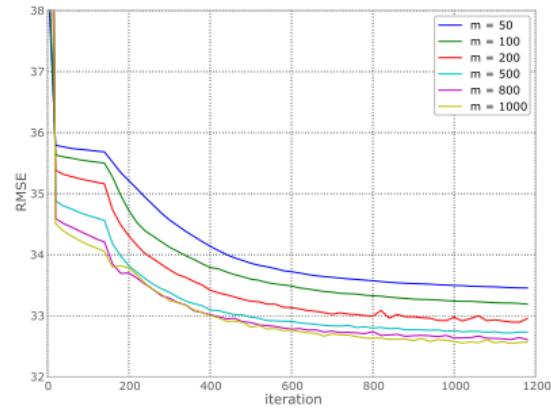
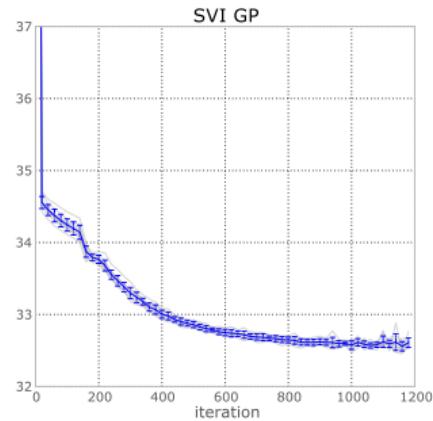
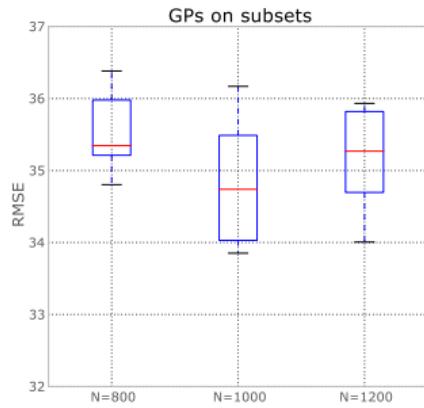
$$\theta_{t+1} = \theta_t - \eta \frac{d\mathcal{L}_{\text{MB}}}{d\theta}.$$

Example: 2D synthetic data



Example: Airline delay data

Flight delays for every commercial flight in the USA from January to April 2008. 700,000 train, 100,000 test



Example from the paper

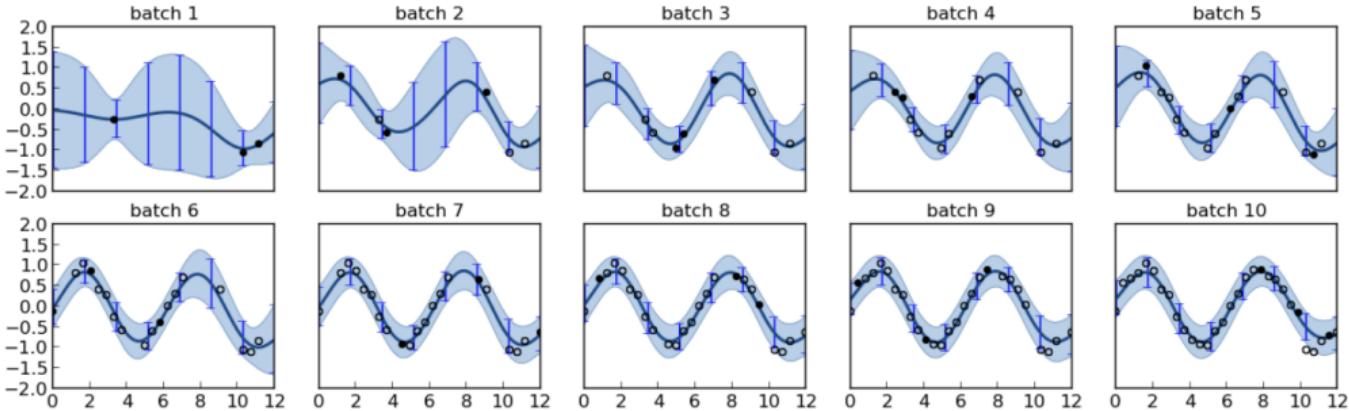


Figure 2: Stochastic variational inference on a trivial GP regression problem. Each pane shows the posterior of the GP after a batch of data, marked as solid points. Previously seen (and discarded) data are marked as empty points, the distribution $q(\mathbf{u})$ is represented by vertical errorbars.

(from Hensman et al: Gaussian processes for big data)

Avoid direct inverses

- Kernel matrix may have high condition number (particularly under squared exponential kernel)
 - direct inverse numerically unstable
- Instead: Cholesky and (triangular) solve
- Sparse inducing points help

Jitter

- Kernel matrix is positive definite ... or is it?
- Might not be *numerically* positive definite (near-zero but negative eigenvalues)
 - Cholesky failures...
- **Solution:** add small diagonal *jitter* (or *nugget*) to kernel matrix (e.g., $10^{-6}\mathbf{I}$)
- Corresponds to “observing” function values under (very small) noise to make up for finite machine precision

Whitening

- Prior $p(\mathbf{u})$ strongly correlates the elements of $q(\mathbf{u}) = \mathcal{N}(\mathbf{m}, \mathbf{S})$
→ Optimization over \mathbf{m} and \mathbf{S} is challenging
- **Whitening:** reparameterization
 - Define $\mathbf{u} = \mathbf{L}\mathbf{v}$ (+prior mean)
where $\mathbf{L} = \text{chol}(\mathbf{K})$ or $\mathbf{K} = \mathbf{L}\mathbf{L}^\top$
 - Prior $p(\mathbf{v}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, completely uncorrelated
 - Optimize over $q(\mathbf{v}) = \mathcal{N}(\mathbf{m}', \mathbf{S}')$ instead
 - Approximate posterior $q(\mathbf{u}) = \mathcal{N}(\mathbf{L}\mathbf{m}', \mathbf{L}\mathbf{S}'\mathbf{L}^\top)$
- (similar idea to preconditioning in CG)

How do we initialize all these parameters?

- Kernel hyperparameters: as for exact GP regression
see <https://tinyurl.com/guide2gp> (scale your data!)
- Approximate posterior $q(\mathbf{u})$: at the prior $\mathcal{N}(\mathbf{0}, \mathbf{K})$
whitening: initialize $q(\mathbf{v}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Inducing point locations Z :
 - Low dimensions: grid
 - (Randomly shuffled!) subset of data
 - Cluster centers (e.g., k-means)
 - Determinantal point process
- Random restarts

Inducing points method summary

- The inducing point approximation allows us to
 - ... scale Gaussian processes to big data
 - ... use non-Gaussian likelihoods
- It reduces the computational complexity from $\mathcal{O}(N^3)$ to $\mathcal{O}(M^3)$, where $M \ll N$
- It's implemented in most GP toolboxes, e.g. GPy (numpy) and gpflow (tensorflow)

Section 7

Recap

Recap on scaling to large data

The naïve $\mathcal{O}(n^3)$ computational bottleneck ($\mathcal{O}(n^2)$ memory) can be tackled by

- **Exploiting structure in the data**
(data on grid, inputs are in 1D, ...)
- **Exploiting structure in the GP prior**
(GP prior is stationary, separable over input dimensions, ...)
- **Solving the linear system approximately**
(conjugate-gradient solvers)
- **Split problem into smaller chunks**
(local experts, subset of data, divide & conquer, ...)
- **Approximate the problem**
(Nyström, low-rank, inducing points, ...)
- **Approximate the problem solution**
(SVGP = sparse (and stochastic) variational methods)

End of today's lecture

- Computational bottlenecks in exact GP regression (most methods typically apply to general likelihoods)
- The limitations we discussed on the first lecture are actually not that severe
- We will look at one more approach for speeding up GPs later in the course

