

# Special course on Gaussian processes: Session #4

Michael Riis Andersen

Aalto University

*michael.riis@gmail.com*

30/1-19

## 1 Computational challenges

- Computational complexity of GP regression
- Non-Gaussian likelihoods: GP classification

## 2 Approximate inference

- Variational inference: scratching the surface
- Inducing points approximations

# Computational complexity of Gaussian process regression

- The key equations for predictions (with Gaussian likelihood)

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

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

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

- Recall: If  $\mathbf{A} \in \mathbb{R}^{N \times M}$  and  $\mathbf{b} \in \mathbb{R}^M$ , then the cost of computing  $\mathbf{A}\mathbf{b}$  is  $\mathcal{O}(NM)$
- Recall: If  $\mathbf{C} \in \mathbb{R}^{N \times N}$ , then the cost of computing  $\mathbf{C}^{-1}$  is  $\mathcal{O}(N^3)$
- What is computational complexity for computing the posterior distribution for 1 test point based on a data set with  $N$  observations? What is the dominating operation?

# Computational complexity of Gaussian process regression

- The key equations for predictions (with Gaussian likelihood)

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

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

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

- Recall: If  $\mathbf{A} \in \mathbb{R}^{N \times M}$  and  $\mathbf{b} \in \mathbb{R}^M$ , then the cost of computing  $\mathbf{A}\mathbf{b}$  is  $\mathcal{O}(NM)$
- Recall: If  $\mathbf{C} \in \mathbb{R}^{N \times N}$ , then the cost of computing  $\mathbf{C}^{-1}$  is  $\mathcal{O}(N^3)$
- What is computational complexity for computing the posterior distribution for 1 test point based on a data set with  $N$  observations? What is the dominating operation?
- $\mathbf{h} = (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$  scales as  $\mathcal{O}(N^3)$

# Computational complexity of Gaussian process regression

- The key equations for predictions (with Gaussian likelihood)

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

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

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

- Recall: If  $\mathbf{A} \in \mathbb{R}^{N \times M}$  and  $\mathbf{b} \in \mathbb{R}^M$ , then the cost of computing  $\mathbf{Ab}$  is  $\mathcal{O}(NM)$
- Recall: If  $\mathbf{C} \in \mathbb{R}^{N \times N}$ , then the cost of computing  $\mathbf{C}^{-1}$  is  $\mathcal{O}(N^3)$
- What is computational complexity for computing the posterior distribution for 1 test point based on a data set with  $N$  observations? What is the dominating operation?
- $\mathbf{h} = (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$  scales as  $\mathcal{O}(N^3)$ ,  $\mu_* = \mathbf{k}_{f_* f} \mathbf{h}$  scales as  $\mathcal{O}(N)$

# Computational complexity of Gaussian process regression

- The key equations for predictions (with Gaussian likelihood)

$$\begin{aligned}p(f_* | \mathbf{y}) &= \mathcal{N}(f_* | \mu_*, \sigma_*^2) \\ \mu_* &= \mathbf{k}_{f_* f} (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \\ \sigma_*^2 &= K_{f_* f_*} - \mathbf{k}_{f_* f} (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{f_* f}^T\end{aligned}$$

- Recall: If  $\mathbf{A} \in \mathbb{R}^{N \times M}$  and  $\mathbf{b} \in \mathbb{R}^M$ , then the cost of computing  $\mathbf{Ab}$  is  $\mathcal{O}(NM)$
- Recall: If  $\mathbf{C} \in \mathbb{R}^{N \times N}$ , then the cost of computing  $\mathbf{C}^{-1}$  is  $\mathcal{O}(N^3)$
- What is computational complexity for computing the posterior distribution for 1 test point based on a data set with  $N$  observations? What is the dominating operation?
- $\mathbf{h} = (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$  scales as  $\mathcal{O}(N^3)$ ,  $\mu_* = \mathbf{k}_{f_* f} \mathbf{h}$  scales as  $\mathcal{O}(N)$
- $N \leq 1000$ : Fine,  $N \leq 10000$ : Slow, but possible,  $N > 10000$ : Prohibitively slow

# Regression vs classification

- Response variable  $\mathbf{y}$  is continuous in regression problems

$$y_n \in \mathbb{R}$$

- Response variable  $\mathbf{y}$  is discrete in classification problems

$$y_n \in \{c_1, c_2, \dots, c_K\}$$

- Classification problems

$\mathbf{X}$  = images,

$y_n \in \{\text{cat}, \text{dog}\}$

$\mathbf{X}$  = X-ray scan,

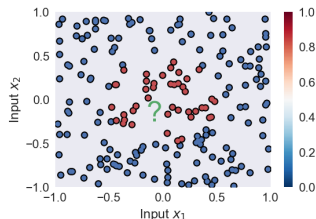
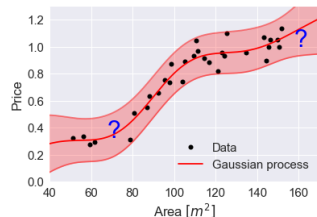
$y_n \in \{\text{tumor}, \text{no tumor}\}$

$\mathbf{X}$  = images of digits,

$y_n \in \{0, 1, 2, \dots, 9\}$

$\mathbf{X}$  = emails,

$y_n \in \{\text{spam}, \text{not spam}\}$



# Regression vs classification

- Response variable  $\mathbf{y}$  is continuous in regression problems

$$y_n \in \mathbb{R}$$

- Response variable  $\mathbf{y}$  is discrete in classification problems

$$y_n \in \{c_1, c_2, \dots, c_K\}$$

- Classification problems

$\mathbf{X}$  = images,

$y_n \in \{\text{cat}, \text{dog}\}$

$\mathbf{X}$  = X-ray scan,

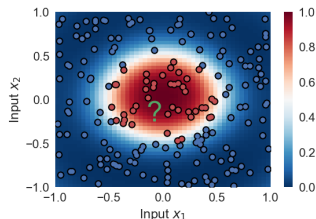
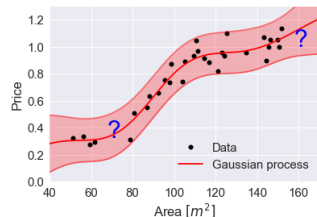
$y_n \in \{\text{tumor}, \text{no tumor}\}$

$\mathbf{X}$  = images of digits,

$y_n \in \{0, 1, 2, \dots, 9\}$

$\mathbf{X}$  = emails,

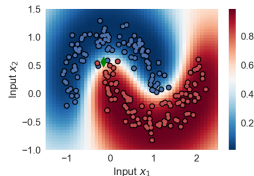
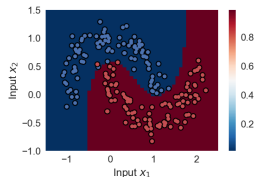
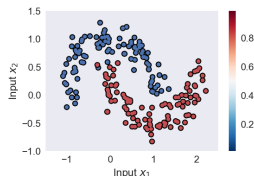
$y_n \in \{\text{spam}, \text{not spam}\}$





# Why Gaussian processes for classification?

- Complex decision boundaries
  - 1 Non-linear boundary
  - 2 Can learn complexity of decision boundary from data
- Probabilistic classification
  - 1 How would you classify the green point?
  - 2 We want to model the uncertainty



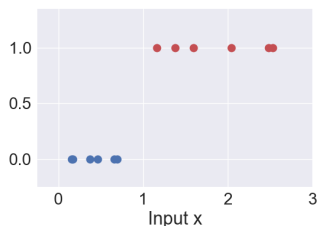
# Why don't we use regression models for classification?

- We focus on binary classification:  $y_n \in \{0, 1\}$  or  $y_n \in \{-1, 1\}$
- We are given a data set  $\{\mathbf{x}_n, y_n\}_{n=1}^N$  and we want to model

$$p(y_n = +1 | \mathbf{x}_n)$$

- What's wrong with simply using the GP regression model with labels:  $y_n \in \{0, 1\}$ :

$$p(y_n = +1 | \mathbf{x}_n) = f(\mathbf{x}_n)$$



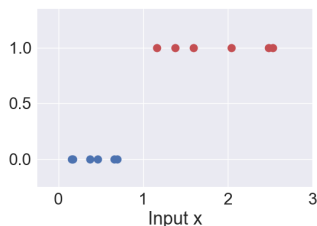
# Why don't we use regression models for classification?

- We focus on binary classification:  $y_n \in \{0, 1\}$  or  $y_n \in \{-1, 1\}$
- We are given a data set  $\{\mathbf{x}_n, y_n\}_{n=1}^N$  and we want to model

$$p(y_n = +1 | \mathbf{x}_n)$$

- What's wrong with simply using the GP regression model with labels:  $y_n \in \{0, 1\}$ :

$$p(y_n = +1 | \mathbf{x}_n) = f(\mathbf{x}_n)$$



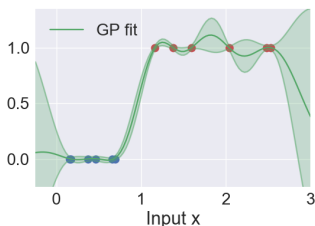
# Why don't we use regression models for classification?

- We focus on binary classification:  $y_n \in \{0, 1\}$  or  $y_n \in \{-1, 1\}$
- We are given a data set  $\{\mathbf{x}_n, y_n\}_{n=1}^N$  and we want to model

$$p(y_n = +1 | \mathbf{x}_n)$$

- What's wrong with simply using the GP regression model with labels:  $y_n \in \{0, 1\}$ :

$$p(y_n = +1 | \mathbf{x}_n) = f(\mathbf{x}_n)$$



# Gaussian process classification setup (I)

- We'll use a 'squashing function'  $\phi : \mathbb{R} \rightarrow (0, 1)$  with  $y_n \in \{-1, 1\}$

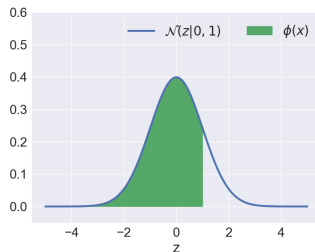
$$p(y_n | \mathbf{x}_n) = \phi(y_n \cdot f(\mathbf{x}_n)) \in (0, 1)$$

- Multiple possible choices for  $\phi(\cdot)$ , we'll use the standard normal CDF

$$\phi(x) = \int_{-\infty}^x \mathcal{N}(z|0, 1) dz$$

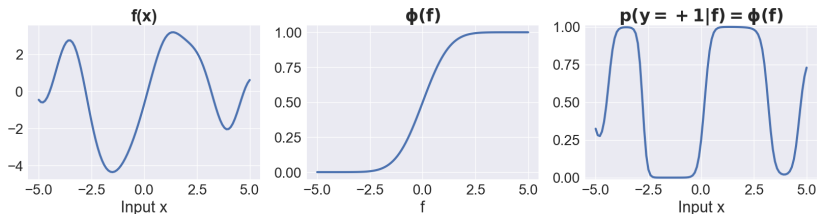
## Discuss with your neighbour

- 1 What is  $\phi(0)$ ?
- 2 What is  $\phi(-\infty)$ ?
- 3 What is  $\phi(\infty)$ ?
- 4 What is  $\phi(x) + \phi(-x)$ ?
- 5 Is  $\phi(y_n f(\mathbf{x}_n))$  normalized wrt.  $y_n$ ?

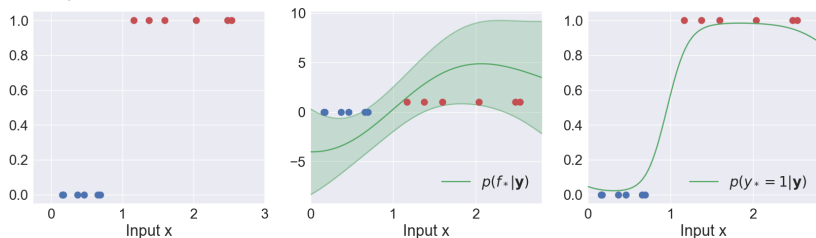


# Gaussian process classification setup (II)

- We map the unknown function  $f(\mathbf{x})$  through the squashing function



- Example re-visited



# Gaussian process classification: Inference

Three steps to compute the predictive distribution for a new test point  $\mathbf{x}_*$

$$p(\mathbf{y}, \mathbf{f}) = \prod_{n=1}^N p(y_n | f_n) p(\mathbf{f}) = \prod_{n=1}^N \phi(y_n \cdot f_n) \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K})$$

- Step 1: Compute posterior distribution of  $p(\mathbf{f} | \mathbf{y})$ :

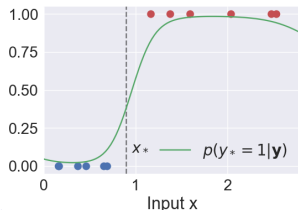
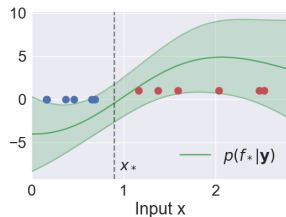
$$p(\mathbf{f} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f}) p(\mathbf{f})}{p(\mathbf{y})}$$

- Step 2: Compute posterior of  $f_*$  for new test point  $\mathbf{x}_*$ :

$$p(f_* | \mathbf{y}) = \int p(f_* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f}$$

- Step 3: Compute predictive distribution

$$p(y_* | \mathbf{y}) = \int \phi(y_* \cdot f_*) p(f_* | \mathbf{y}) df_*$$



# Gaussian process classification: Inference

Three steps to compute the predictive distribution for a new test point  $\mathbf{x}_*$

$$p(\mathbf{y}, \mathbf{f}) = \prod_{n=1}^N p(y_n | f_n) p(\mathbf{f}) = \prod_{n=1}^N \phi(y_n \cdot f_n) \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K})$$

- Step 1: Compute posterior distribution of  $p(\mathbf{f} | \mathbf{y})$ :

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

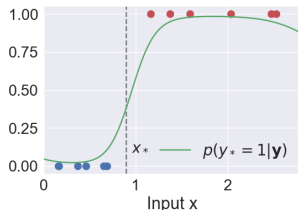
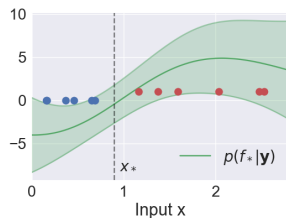
- Step 2: Compute posterior of  $f_*$  for new test point  $\mathbf{x}_*$ :

$$p(f_* | \mathbf{y}) = \int p(f_* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f}$$

- Step 3: Compute predictive distribution

$$p(y_* | \mathbf{y}) = \int \phi(y_* \cdot f_*) p(f_* | \mathbf{y}) df_*$$

- Unfortunately, these distributions are analytically intractable.





# Gaussian process classification: Inference

Three steps to compute the predictive distribution for a new test point  $\mathbf{x}_*$

$$p(\mathbf{y}, \mathbf{f}) = \prod_{n=1}^N p(y_n | f_n) p(\mathbf{f}) = \prod_{n=1}^N \phi(y_n \cdot f_n) \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K})$$

- Step 1: Compute posterior distribution of  $p(\mathbf{f} | \mathbf{y})$ :

$$p(\mathbf{f} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f}) p(\mathbf{f})}{p(\mathbf{y})} \approx q(\mathbf{f})$$

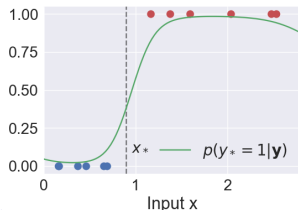
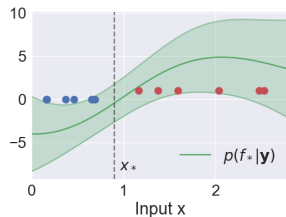
- Step 2: Compute posterior of  $f_*$  for new test point  $\mathbf{x}_*$ :

$$p(f_* | \mathbf{y}) = \int p(f_* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \approx \int p(f_* | \mathbf{f}) q(\mathbf{f}) d\mathbf{f}$$

- Step 3: Compute predictive distribution

$$p(y_* | \mathbf{y}) = \int \phi(y_* \cdot f_*) p(f_* | \mathbf{y}) df_*$$

- Unfortunately, these distributions are analytically intractable.



# Computational problems

We need to figure out what to do when

- ... likelihood is non-Gaussian?
- ... inference becomes slow due to large  $N$ ?

# Computational problems

We need to figure out what to do when

- ... likelihood is non-Gaussian?
- ... inference becomes slow due to large  $N$ ?

Variational inference

# Computational problems

We need to figure out what to do when

- ... likelihood is non-Gaussian?
- ... inference becomes slow due to large  $N$ ?

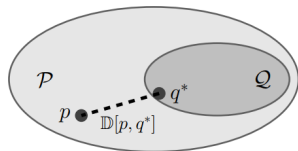
Variational inference

- General framework for approximate Bayesian inference
- Many recent application in the machine learning literature:
  - 1 GPs for big data
  - 2 GPs with non-Gaussian likelihoods
  - 3 Deep Gaussian processes
  - 4 Convolutional Gaussian processes
  - 5 Variational autoencoders (VAEs)
  - 6 ...

# Variational inference: the big picture

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

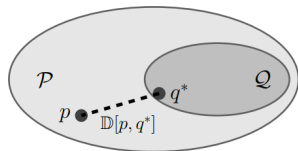
- 1 Define some "simple" family of distribution  $\mathcal{Q}$ .



# Variational inference: the big picture

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

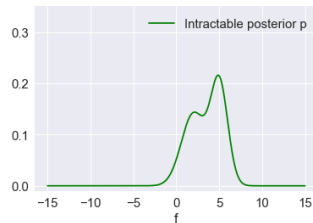
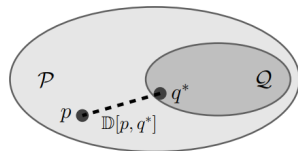
- 1 Define some "simple" family of distribution  $\mathcal{Q}$ .
- 2 Define some way to compute a "distance"  $\mathbb{D}[q, p]$  between each of the distribution  $q \in \mathcal{Q}$  and the intractable distribution  $p$



# Variational inference: the big picture

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

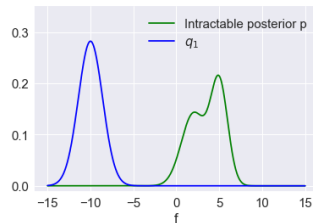
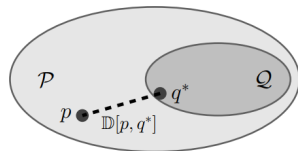
- 1 Define some "simple" family of distribution  $\mathcal{Q}$ .
- 2 Define some way to compute a "distance"  $\mathbb{D}[q, p]$  between each of the distribution  $q \in \mathcal{Q}$  and the intractable distribution  $p$



# Variational inference: the big picture

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

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

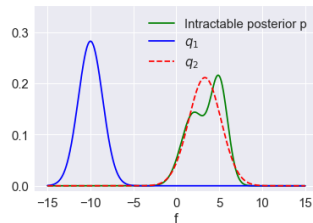
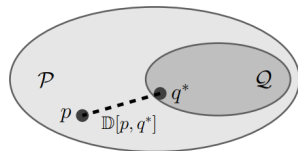




# Variational inference: the big picture

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

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

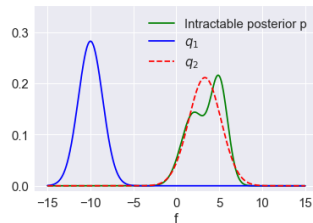
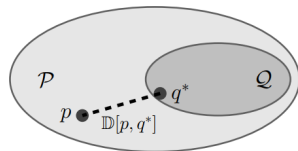


# Variational inference: the big picture

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

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

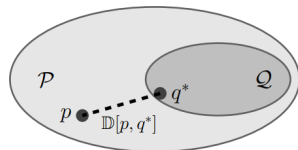
$$\mathbb{D}[q_1, p] > \mathbb{D}[q_2, p]$$



# Variational inference: the big picture

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

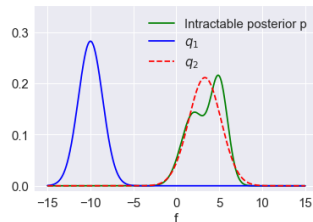
- 1 Define some "simple" family of distribution  $\mathcal{Q}$ .
- 2 Define some way to compute a "distance"  $\mathbb{D}[q, p]$  between each of the distribution  $q \in \mathcal{Q}$  and the intractable distribution  $p$



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

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

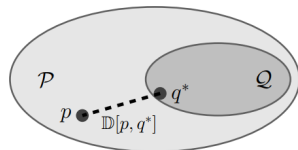
$$q^* = \arg \min_{q \in \mathcal{Q}} \mathbb{D}[q, p]$$



# Variational inference: the big picture

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

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

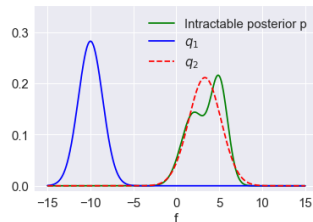


$$\mathbb{D}[q_1, p] > \mathbb{D}[q_2, p]$$

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

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

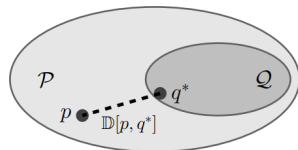
- 4 Use  $q^*$  as an approximation of  $p$



# Variational inference: the big picture

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

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

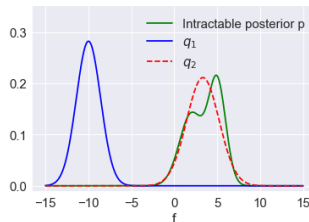


$$\mathbb{D}[q_1, p] > \mathbb{D}[q_2, p]$$

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

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

- 4 Use  $q^*$  as an approximation of  $p$



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

# Variational inference I

- We will use to the *Kullback-Leibler divergence* to "measure distances" between distributions

$$\mathbb{D}[q||p] = \int q(\mathbf{f}) \ln \frac{q(\mathbf{f})}{p(\mathbf{f})} d\mathbf{f} = \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f})} \right]$$

# Variational inference I

- We will use to the *Kullback-Leibler divergence* to "measure distances" between distributions

$$\mathbb{D}[q||p] = \int q(\mathbf{f}) \ln \frac{q(\mathbf{f})}{p(\mathbf{f})} d\mathbf{f} = \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f})} \right]$$

- Most important properties for our purpose:
  - 1 Positive definite:  $\mathbb{D}[q||p] \geq 0$
  - 2 Identity of indiscernibles:  $\mathbb{D}[q||p] = 0 \iff p = q$  (a.e.)
  - 3 Not-symmetric:  $\mathbb{D}[q||p] \neq \mathbb{D}[p||q]$

# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$



# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$

$$\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] = \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right]$$

# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \\ &= \mathbb{E}_q [\ln q(\mathbf{f}) - \ln p(\mathbf{f}|\mathbf{y})]\end{aligned}$$

# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \\ &= \mathbb{E}_q [\ln q(\mathbf{f}) - \ln p(\mathbf{f}|\mathbf{y})] \\ &= \mathbb{E}_q [\ln q(\mathbf{f})] - \mathbb{E}_q [\ln p(\mathbf{f}|\mathbf{y})]\end{aligned}$$

# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \\ &= \mathbb{E}_q [\ln q(\mathbf{f}) - \ln p(\mathbf{f}|\mathbf{y})] \\ &= \mathbb{E}_q [\ln q(\mathbf{f})] - \mathbb{E}_q [\ln p(\mathbf{f}|\mathbf{y})]\end{aligned}$$

Define the *entropy* of  $q$  as  $\mathcal{H}[q] \equiv -\mathbb{E}_q [\ln q(\mathbf{f})]$

# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \\ &= \mathbb{E}_q [\ln q(\mathbf{f}) - \ln p(\mathbf{f}|\mathbf{y})] \\ &= \mathbb{E}_q [\ln q(\mathbf{f})] - \mathbb{E}_q [\ln p(\mathbf{f}|\mathbf{y})]\end{aligned}$$

Define the *entropy* of  $q$  as  $\mathcal{H}[q] \equiv -\mathbb{E}_q [\ln q(\mathbf{f})]$

$$\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] = -\mathcal{H}[q] - \mathbb{E}_q [\ln p(\mathbf{f}|\mathbf{y})]$$

# Variational inference II

Our goal is to minimize the KL divergence between some approximation  $q \in \mathcal{Q}$  and some posterior distribution  $p(\mathbf{f}|\mathbf{y})$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= \mathbb{E}_q \left[ \ln \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \\ &= \mathbb{E}_q [\ln q(\mathbf{f}) - \ln p(\mathbf{f}|\mathbf{y})] \\ &= \mathbb{E}_q [\ln q(\mathbf{f})] - \mathbb{E}_q [\ln p(\mathbf{f}|\mathbf{y})]\end{aligned}$$

Define the *entropy* of  $q$  as  $\mathcal{H}[q] \equiv -\mathbb{E}_q [\ln q(\mathbf{f})]$

$$\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] = -\mathcal{H}[q] - \mathbb{E}_q [\ln p(\mathbf{f}|\mathbf{y})]$$

Last term depends on the exact posterior  $p(\mathbf{f}|\mathbf{y})$ , which is intractable.

# Variational inference III

Using the def. of conditional densities, we can write:  $p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y},\mathbf{f})}{p(\mathbf{y})}$

$$\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] = -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{f}|\mathbf{y})]$$

# Variational inference III

Using the def. of conditional densities, we can write:  $p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{f}|\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q\left[\ln \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}\right]\end{aligned}$$



## Variational inference III

Using the def. of conditional densities, we can write:  $p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{f}|\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q\left[\ln \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}\right] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \mathbb{E}_q[\ln p(\mathbf{y})]\end{aligned}$$

## Variational inference III

Using the def. of conditional densities, we can write:  $p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{f}|\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q\left[\ln \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}\right] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \mathbb{E}_q[\ln p(\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \ln p(\mathbf{y})\end{aligned}$$

# Variational inference III

Using the def. of conditional densities, we can write:  $p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{f}|\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q\left[\ln \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}\right] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \mathbb{E}_q[\ln p(\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \ln p(\mathbf{y})\end{aligned}$$

Let's re-arrange the terms

$$\ln p(\mathbf{y}) = \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})]$$

## Variational inference III

Using the def. of conditional densities, we can write:  $p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}$

$$\begin{aligned}\mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})] &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{f}|\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q\left[\ln \frac{p(\mathbf{y}, \mathbf{f})}{p(\mathbf{y})}\right] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \mathbb{E}_q[\ln p(\mathbf{y})] \\ &= -\mathcal{H}[q] - \mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \ln p(\mathbf{y})\end{aligned}$$

Let's re-arrange the terms

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q[\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})]$$

$\mathcal{L}[q]$  does not depend on the posterior  $p(\mathbf{f}|\mathbf{y})$ , but only on the joint density  $p(\mathbf{y}, \mathbf{f})$ .

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant
- 2  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq 0$  is non-negative



# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant
- 2  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq 0$  is non-negative
- 3  $\mathcal{L}[q]$  only depends on  $q$  and the joint density  $p(\mathbf{y}, \mathbf{f})$

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant
- 2  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq 0$  is non-negative
- 3  $\mathcal{L}[q]$  only depends on  $q$  and the joint density  $p(\mathbf{y}, \mathbf{f})$

Some consequences

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant
- 2  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq 0$  is non-negative
- 3  $\mathcal{L}[q]$  only depends on  $q$  and the joint density  $p(\mathbf{y}, \mathbf{f})$

Some consequences

- 1  $\mathcal{L}[q]$  is a *lower bound* of  $\ln p(\mathbf{y})$ . That is:  $\ln p(\mathbf{y}) \geq \mathcal{L}[q]$

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant
- 2  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq 0$  is non-negative
- 3  $\mathcal{L}[q]$  only depends on  $q$  and the joint density  $p(\mathbf{y}, \mathbf{f})$

Some consequences

- 1  $\mathcal{L}[q]$  is a *lower bound* of  $\ln p(\mathbf{y})$ . That is:  $\ln p(\mathbf{y}) \geq \mathcal{L}[q]$
- 2 Maximizing  $\mathcal{L}[q]$  is equivalent to minimizing  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$

# Variational inference IV

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

Let's make a few observations

- 1  $\ln p(\mathbf{y})$  is a constant
- 2  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq 0$  is non-negative
- 3  $\mathcal{L}[q]$  only depends on  $q$  and the joint density  $p(\mathbf{y}, \mathbf{f})$

Some consequences

- 1  $\mathcal{L}[q]$  is a *lower bound* of  $\ln p(\mathbf{y})$ . That is:  $\ln p(\mathbf{y}) \geq \mathcal{L}[q]$
- 2 Maximizing  $\mathcal{L}[q]$  is equivalent to minimizing  $\mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$

**Key take-away: we can fit the variational approx.  $q$  by optimizing  $\mathcal{L}$**

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

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

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})] + \mathcal{H}[q]}_{\mathcal{L}[q]} + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

- $\mathcal{L}[q]$  is often called the *Evidence Lower Bound* (ELBO)
- The first term in  $\mathcal{L}[q]$  can be interpreted as a data fit term and the second term can be interpreted as a regularization term

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

- $\mathcal{L}[q]$  is often called the *Evidence Lower Bound* (ELBO)
- The first term in  $\mathcal{L}[q]$  can be interpreted as a data fit term and the second term can be interpreted as a regularization term
- If we want to approximate  $p(\mathbf{f}|\mathbf{y})$ , then  $q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{V})$



$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

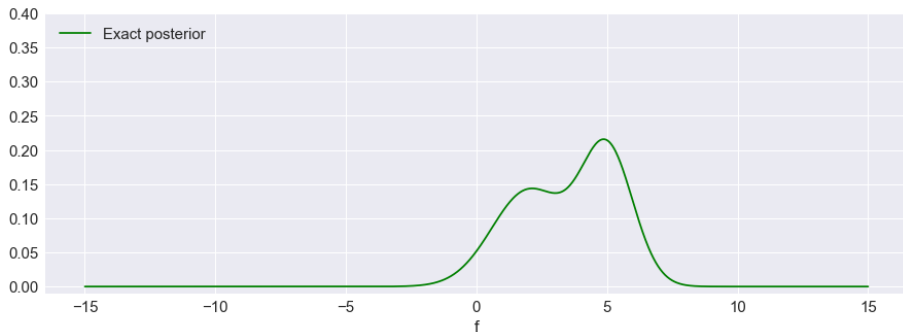
- $\mathcal{L}[q]$  is often called the *Evidence Lower Bound* (ELBO)
- The first term in  $\mathcal{L}[q]$  can be interpreted as a data fit term and the second term can be interpreted as a regularization term
- If we want to approximate  $p(\mathbf{f}|\mathbf{y})$ , then  $q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{V})$
- Define  $\boldsymbol{\lambda} = \{\mathbf{m}, \mathbf{V}\}$ , then we can write  $\mathcal{L}[q] = \mathcal{L}[\boldsymbol{\lambda}]$

$$\ln p(\mathbf{y}) = \underbrace{\mathbb{E}_q [\ln p(\mathbf{y}, \mathbf{f})]}_{\mathcal{L}[q]} + \mathcal{H}[q] + \mathbb{D}[q(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})]$$

- $\mathcal{L}[q]$  is often called the *Evidence Lower Bound* (ELBO)
- The first term in  $\mathcal{L}[q]$  can be interpreted as a data fit term and the second term can be interpreted as a regularization term
- If we want to approximate  $p(\mathbf{f}|\mathbf{y})$ , then  $q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{V})$
- Define  $\boldsymbol{\lambda} = \{\mathbf{m}, \mathbf{V}\}$ , then we can write  $\mathcal{L}[q] = \mathcal{L}[\boldsymbol{\lambda}]$
- In practice, we optimize  $\mathcal{L}[\boldsymbol{\lambda}]$  using gradient-based methods

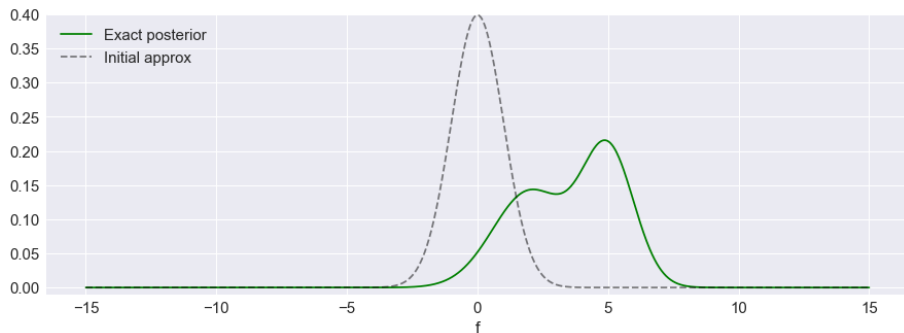
# 1D Toy example I

- Assume we have some model  $p(y, f)$  that gives rise to some intractable posterior  $p(f|y)$
- We want to approximate  $p(f|y)$  using a variational approximation
- In 1D:  $\mathcal{Q}$  is the set of univariate Gaussian, i.e.  $q_{\lambda}(x) = \mathcal{N}(x|m, v)$ , where we denote  $\lambda = \{m, v\}$
- We initialize our approximation as  $q(f) = \mathcal{N}(f|0, 1)$



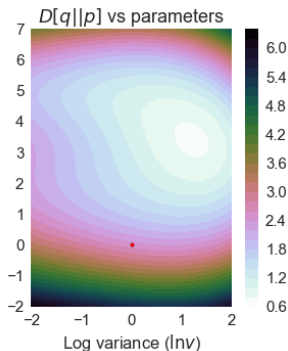
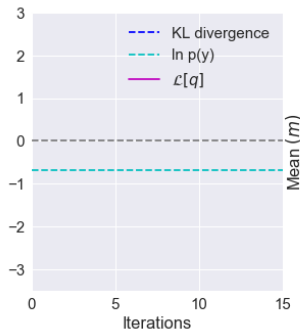
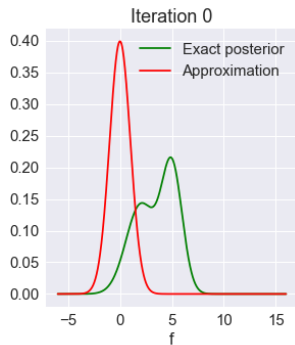
# 1D Toy example I

- Assume we have some model  $p(y, f)$  that gives rise to some intractable posterior  $p(f|y)$
- We want to approximate  $p(f|y)$  using a variational approximation
- In 1D:  $\mathcal{Q}$  is the set of univariate Gaussian, i.e.  $q_\lambda(x) = \mathcal{N}(x|m, \nu)$ , where we denote  $\lambda = \{m, \nu\}$
- We initialize our approximation as  $q(f) = \mathcal{N}(f|0, 1)$



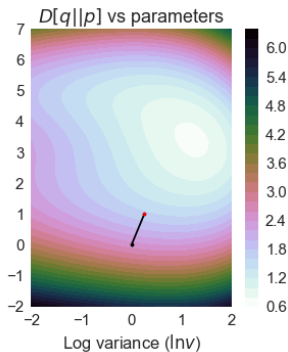
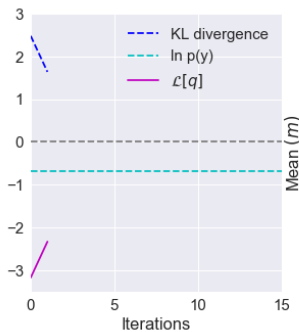
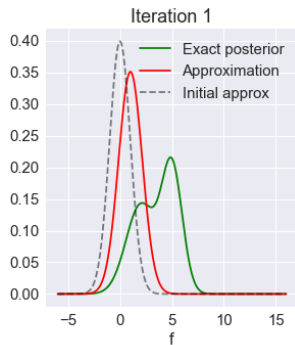
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



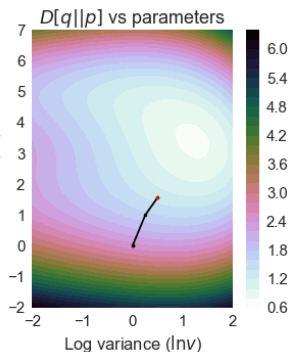
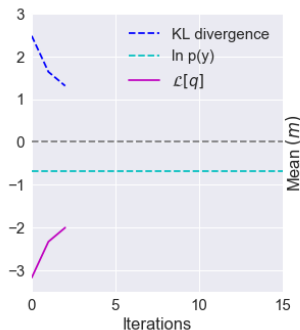
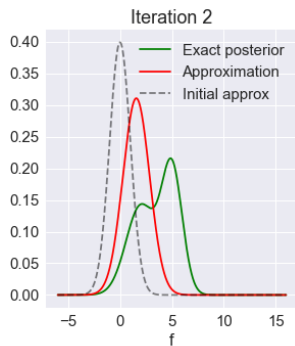
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



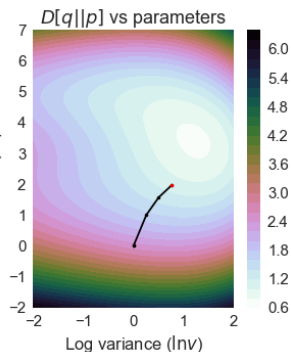
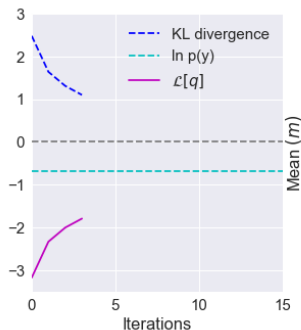
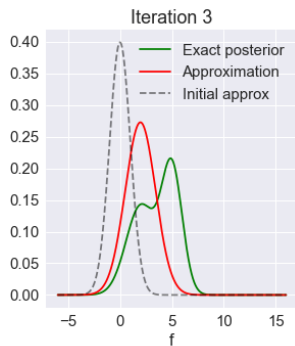
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



# 1D Toy example II

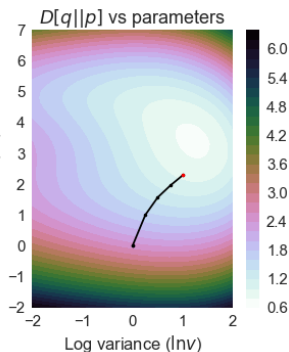
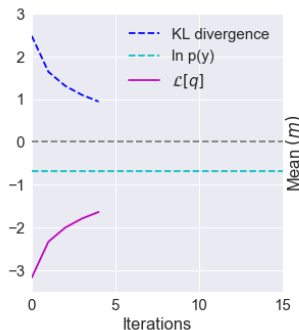
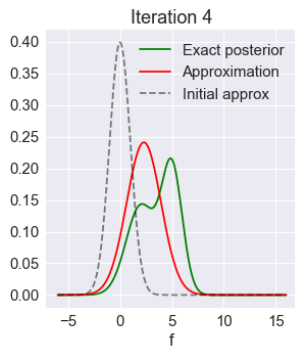
- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$





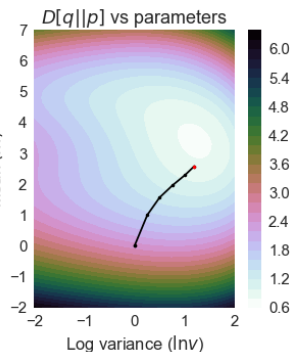
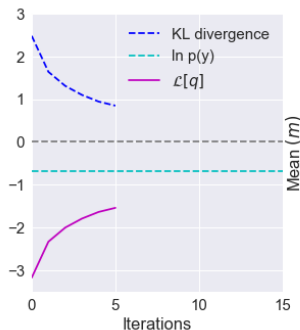
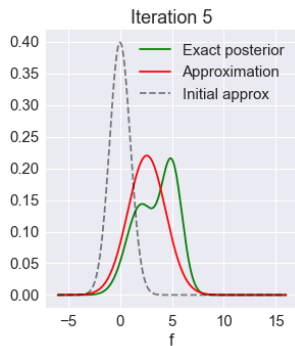
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



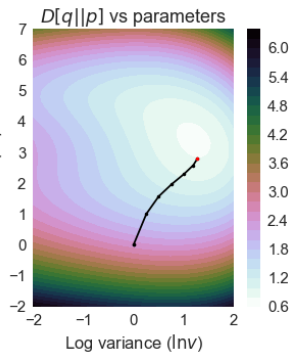
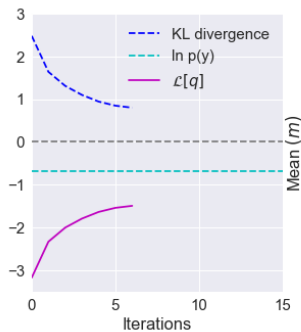
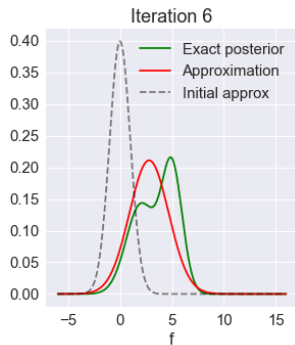
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



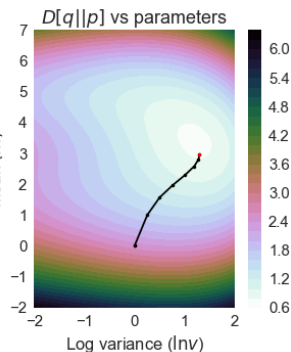
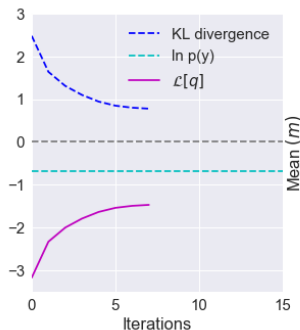
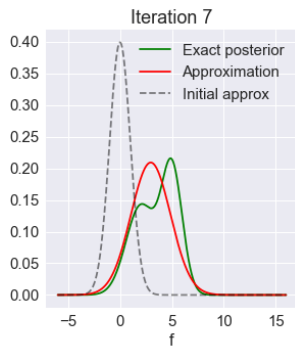
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



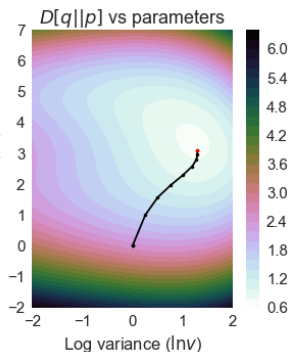
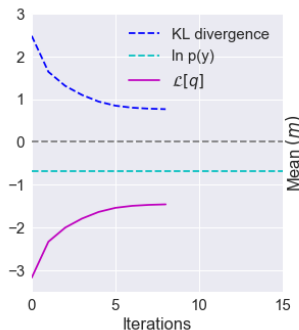
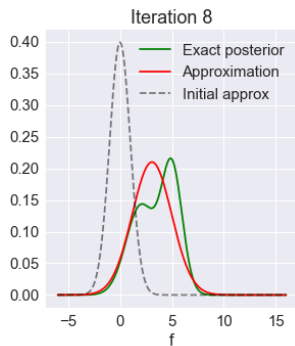
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



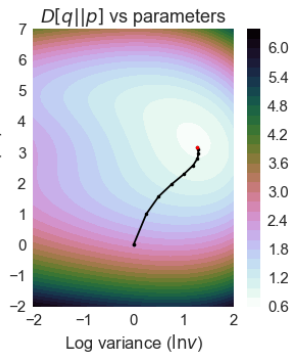
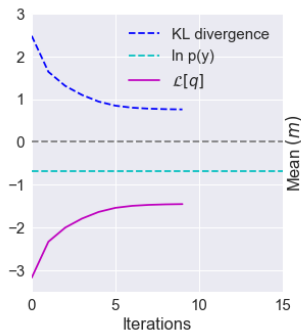
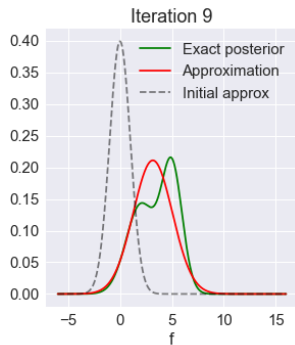
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



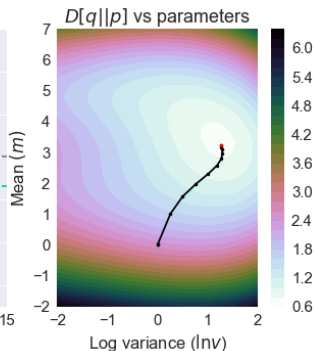
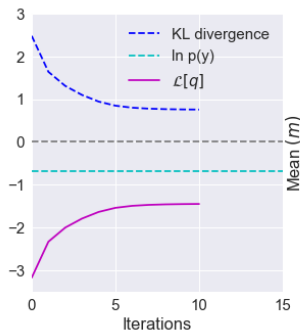
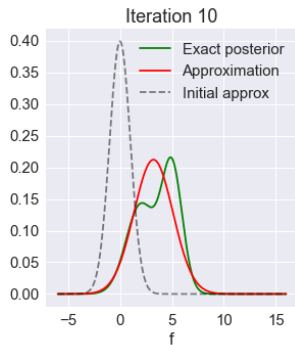
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



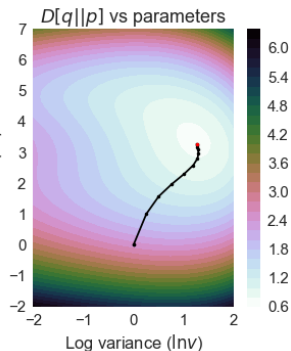
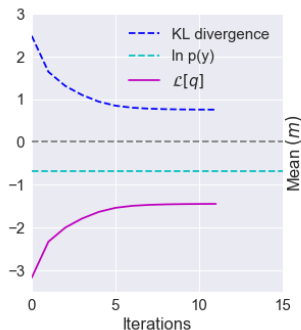
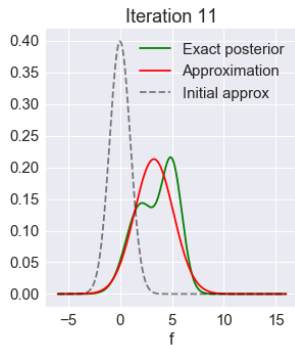
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



# 1D Toy example II

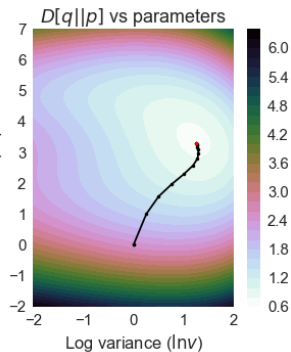
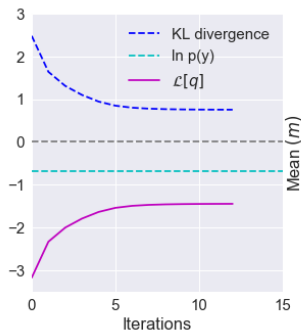
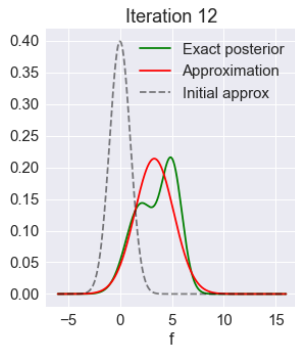
- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$





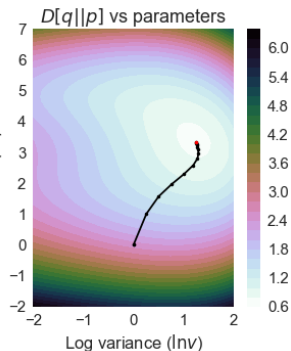
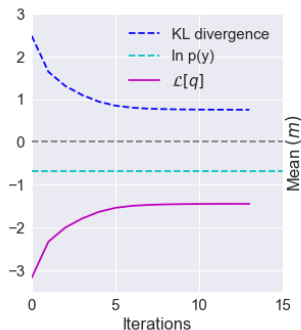
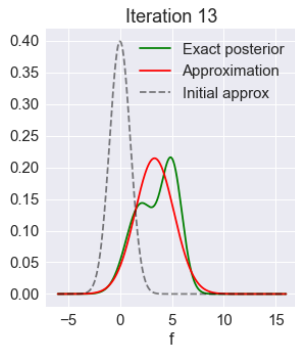
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



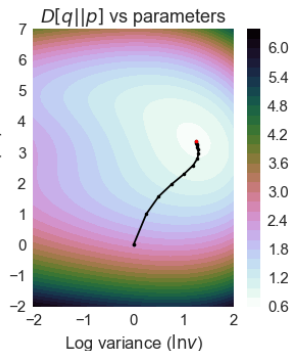
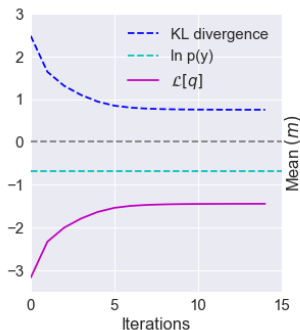
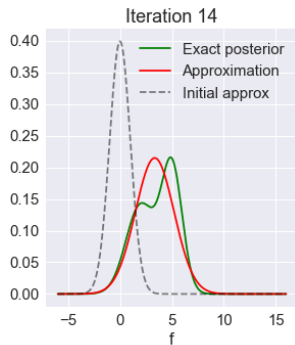
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



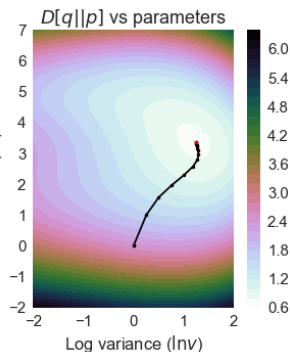
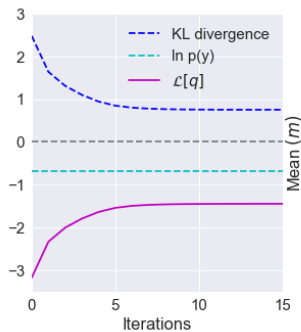
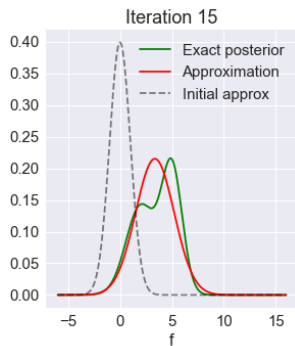
# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



# 1D Toy example II

- Gradient ascent:  $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L} [\lambda]$
- $\ln p(\mathbf{y}) = \mathcal{L} [\lambda] + \mathbb{D} [q_{\lambda}(\mathbf{f}) || p(\mathbf{f}|\mathbf{y})] \geq \mathcal{L} [\lambda]$



- Let's see how we can use combine the ideas from variational inference with inducing points methods to solve the two computational problems:
  - 1 The computational complexity of GPs is  $\mathcal{O}(N^3)$
  - 2 How to handle non-Gaussian likelihoods

# Solution: Inducing point methods

- The main idea is to "represent" the information from the full dataset using a smaller "virtual" dataset

# Solution: Inducing point methods

- The main idea is to "represent" the information from the full dataset using a smaller "virtual" dataset
- Recall our GP model:

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

# Solution: Inducing point methods

- The main idea is to "represent" the information from the full dataset using a smaller "virtual" dataset
- Recall our GP model:

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

- 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.  $\mathbf{x}_i, \mathbf{z}_j \in \mathbb{R}^D$



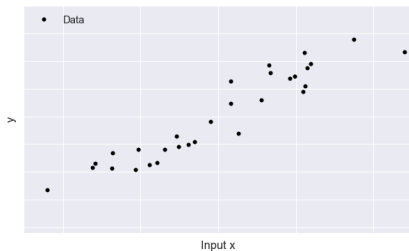
# Solution: Inducing point methods

- The main idea is to "represent" the information from the full dataset using a smaller "virtual" dataset
- Recall our GP model:

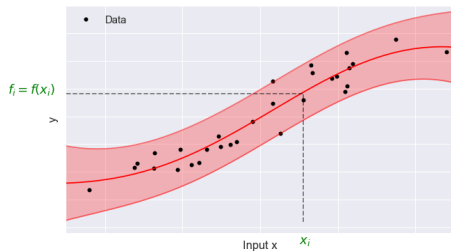
$$p(\mathbf{y}, \mathbf{f}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}), \quad \text{where } \mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)]$$

- 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  $\mathbf{u} = [f(\mathbf{z}_1), f(\mathbf{z}_2), \dots, f(\mathbf{z}_M)]$

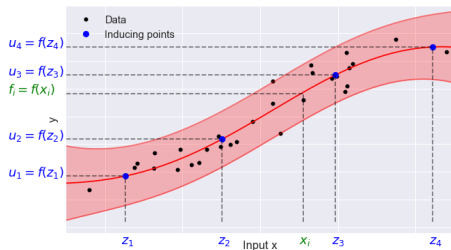
# Inducing point methods



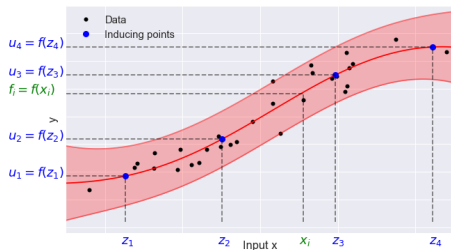
# Inducing point methods



# Inducing point methods

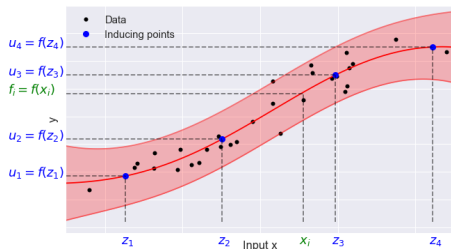


# Inducing point methods



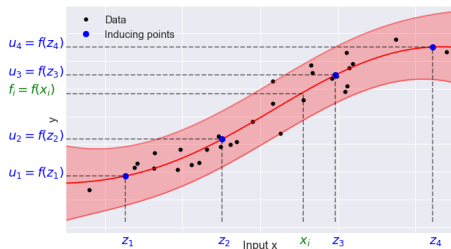
- Goal: choose the set of inducing points such that it contains the same information as the full dataset

# Inducing point methods



- Goal: choose the set of inducing points such that it contains the same information as the full dataset
- Remember: Both  $u_j = f(z_j)$  and  $f_i = f(x_i)$  are random variables

# Inducing point methods



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

# Inducing point methods: the joint model

- The augmented model

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

- Let's decompose the "augmented" model as follows

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

- We can get back to the original model by marginalizing over  $\mathbf{u}$

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



# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$
- Let's write down the KL divergence between  $q(\mathbf{f}, \mathbf{u})$  and  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

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

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$
- Let's write down the KL divergence between  $q(\mathbf{f}, \mathbf{u})$  and  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

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

- As before, we use Bayes rule and do some algebra:

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

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$
- Let's write down the KL divergence between  $q(\mathbf{f}, \mathbf{u})$  and  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

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

- As before, we use Bayes rule and do some algebra:

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

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$
- Let's write down the KL divergence between  $q(\mathbf{f}, \mathbf{u})$  and  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

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

- As before, we use Bayes rule and do some algebra:

$$\begin{aligned} \mathbb{D}[q||p] &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[ \ln \frac{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \ln p(\mathbf{y}) \\ &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[ \ln \frac{q(\mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})} \right] + \ln p(\mathbf{y}) \\ &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] + \ln p(\mathbf{y}) \end{aligned}$$

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$
- Let's write down the KL divergence between  $q(\mathbf{f}, \mathbf{u})$  and  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

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

- As before, we use Bayes rule and do some algebra:

$$\begin{aligned} \mathbb{D}[q||p] &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[ \ln \frac{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \ln p(\mathbf{y}) \\ &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[ \ln \frac{q(\mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})} \right] + \ln p(\mathbf{y}) \\ &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] + \ln p(\mathbf{y}) \end{aligned}$$

- Re-arranging yields

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

# Setting up the approximation

- The idea is now to derive a variational approximation for the posterior  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$
- We choose  $\mathcal{Q}$  be the set of all distributions of the form  $q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$ , where  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$
- Let's write down the KL divergence between  $q(\mathbf{f}, \mathbf{u})$  and  $p(\mathbf{f}, \mathbf{u}|\mathbf{y})$

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

- As before, we use Bayes rule and do some algebra:

$$\begin{aligned} \mathbb{D}[q||p] &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[ \ln \frac{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \ln p(\mathbf{y}) \\ &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} \left[ \ln \frac{q(\mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})} \right] + \ln p(\mathbf{y}) \\ &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] + \ln p(\mathbf{y}) \end{aligned}$$

- Re-arranging yields

$$\begin{aligned} \ln p(\mathbf{y}) &= \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] + \mathbb{D}[q||p] \\ &\geq \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] \equiv \mathcal{L}_3 \end{aligned}$$



# The inducing points approximation

- **Take-away #1:** We can now tractably optimize the lower bound wrt.  $\mathbf{m}$ ,  $\mathbf{S}$ , and even  $\mathbf{z}$

$$\ln p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] \equiv \mathcal{L}_3$$

# The inducing points approximation

- **Take-away #1:** We can now tractably optimize the lower bound wrt.  $\mathbf{m}$ ,  $\mathbf{S}$ , and even  $\mathbf{z}$

$$\ln p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] \equiv \mathcal{L}_3$$

- We will now show that the first decomposes in a very convenient way

# The inducing points approximation

- **Take-away #1:** We can now tractably optimize the lower bound wrt.  $\mathbf{m}$ ,  $\mathbf{S}$ , and even  $\mathbf{z}$

$$\ln p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] \equiv \mathcal{L}_3$$

- We will now show that the first decomposes in a very convenient way
- Remember:  $p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N p(y_i|f_i)$

# The inducing points approximation

- **Take-away #1:** We can now tractably optimize the lower bound wrt.  $\mathbf{m}$ ,  $\mathbf{S}$ , and even  $\mathbf{z}$

$$\ln p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\ln q(\mathbf{u})] \equiv \mathcal{L}_3$$

- We will now show that the first decomposes in a very convenient way
- Remember:  $p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N p(y_i|f_i)$
- Let's have a closer look at the first term

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

# Decomposing the likelihood term

- Let's define the univariate distribution

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

- then we can write

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

# Decomposing the likelihood term

- Let's define the univariate distribution

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

- then we can write

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

- Thus, the "likelihood term" decomposes into a sum over 1D integrals

# Decomposing the likelihood term

- Let's define the univariate distribution

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

- then we can write

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

- Thus, the "likelihood term" decomposes into a sum over 1D integrals
- Can be solved analytically for Gaussian likelihoods and some classification likelihoods

# Decomposing the likelihood term

- Let's define the univariate distribution

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

- then we can write

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

- Thus, the "likelihood term" decomposes into a sum over 1D integrals
- Can be solved analytically for Gaussian likelihoods and some classification likelihoods
- But it is fast to approximate 1D integrals using numerical integration for other likelihoods
- Take away #2:** We can tractably optimize the bound even with non-Gaussian likelihoods



# The resulting bound

- Substituting back into  $\mathcal{L}_3$

$$\ln p(\mathbf{y}) \geq \mathcal{L}_3 = \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i + \mathbb{E}_{q(\mathbf{u})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u})} [\ln q(\mathbf{u})]$$

- We want to optimize  $\mathcal{L}_3$  wrt.  $\lambda = \{\mathbf{m}, \mathbf{S}, \mathbf{z}\}$  using gradient-based methods

$$\nabla_{\lambda} \mathcal{L}_3 = \nabla_{\lambda} \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i + \nabla_{\lambda} \mathbb{E}_{q(\mathbf{u})} [\ln p(\mathbf{u})] - \nabla_{\lambda} \mathbb{E}_{q(\mathbf{u})} [\ln q(\mathbf{u})]$$

# The resulting bound

- Substituting back into  $\mathcal{L}_3$

$$\ln p(\mathbf{y}) \geq \mathcal{L}_3 = \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i + \mathbb{E}_{q(\mathbf{u})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u})} [\ln q(\mathbf{u})]$$

- We want to optimize  $\mathcal{L}_3$  wrt.  $\lambda = \{\mathbf{m}, \mathbf{S}, \mathbf{z}\}$  using gradient-based methods

$$\nabla_{\lambda} \mathcal{L}_3 = \nabla_{\lambda} \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i + \nabla_{\lambda} \mathbb{E}_{q(\mathbf{u})} [\ln p(\mathbf{u})] - \nabla_{\lambda} \mathbb{E}_{q(\mathbf{u})} [\ln q(\mathbf{u})]$$

- We can approximate the gradient as follows (mini-batching)

$$\nabla_{\lambda} \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i \approx \frac{N}{|S|} \sum_{i \in S} \nabla_{\lambda} \int q(f_i) \ln p(y_i|f_i) df_i$$

# The resulting bound

- Substituting back into  $\mathcal{L}_3$

$$\ln p(\mathbf{y}) \geq \mathcal{L}_3 = \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i + \mathbb{E}_{q(\mathbf{u})} [\ln p(\mathbf{u})] - \mathbb{E}_{q(\mathbf{u})} [\ln q(\mathbf{u})]$$

- We want to optimize  $\mathcal{L}_3$  wrt.  $\lambda = \{\mathbf{m}, \mathbf{S}, \mathbf{z}\}$  using gradient-based methods

$$\nabla_{\lambda} \mathcal{L}_3 = \nabla_{\lambda} \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i + \nabla_{\lambda} \mathbb{E}_{q(\mathbf{u})} [\ln p(\mathbf{u})] - \nabla_{\lambda} \mathbb{E}_{q(\mathbf{u})} [\ln q(\mathbf{u})]$$

- We can approximate the gradient as follows (mini-batching)

$$\nabla_{\lambda} \sum_{i=1}^N \int q(f_i) \ln p(y_i|f_i) df_i \approx \frac{N}{|S|} \sum_{i \in S} \nabla_{\lambda} \int q(f_i) \ln p(y_i|f_i) df_i$$

- **Take away #3:** Because it decomposes as a sum over the data points, the bound becomes amenable to stochastic gradient descent (mini-batching) and hence, we can scale the method to really really large datasets!

# Example from the paper

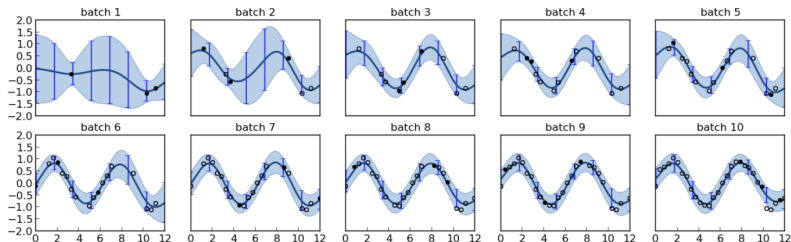


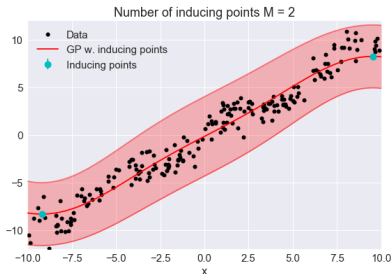
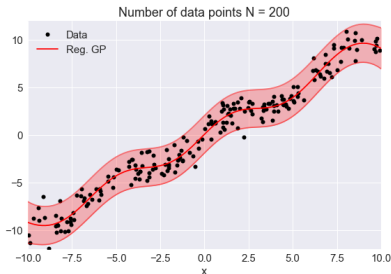
Figure 2: Stochastic variational inference on a trivial GP regression problem. Each pane shows the posterior of the GP after a batch of data, marked as solid points. Previously seen (and discarded) data are marked as empty points, the distribution  $q(\mathbf{u})$  is represented by vertical errorbars.

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

# Inducing points method summary

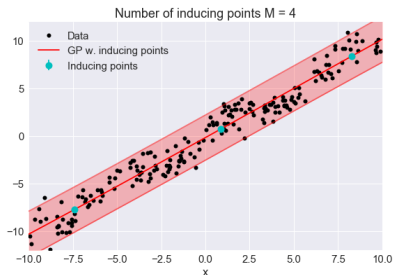
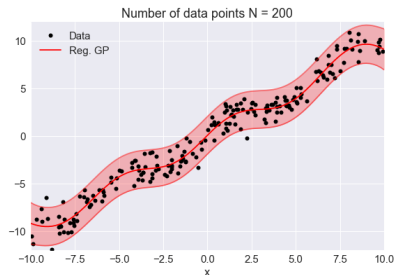
- The inducing point approximation allows us to
  - ... scale Gaussian processes to big data
  - ... use non-Gaussian likelihoods
- It reduces the computational complexity from  $\mathcal{O}(N^3)$  to  $\mathcal{O}(M^3)$ , where  $M \ll N$
- It's implemented in most GP toolboxes, e.g. GPy (numpy) and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

# Example: Number of inducing points



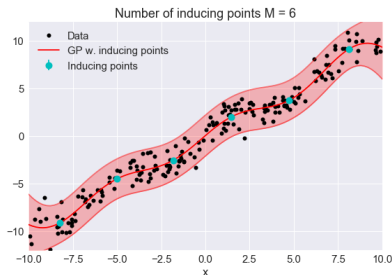
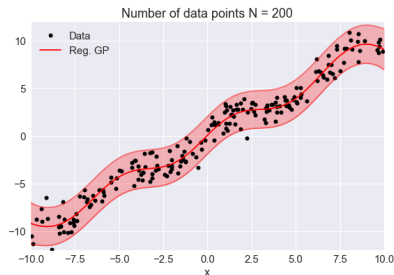
- We can think of the number of inducing points as a parameter that trades off speed for accuracy

# Example: Number of inducing points



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

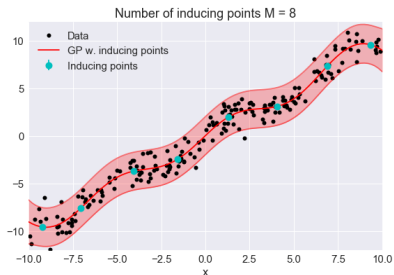
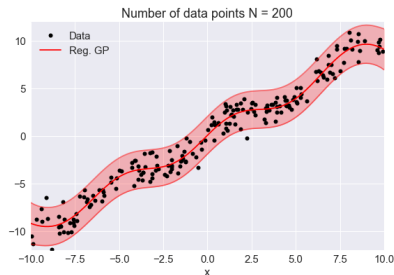
# Example: Number of inducing points



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

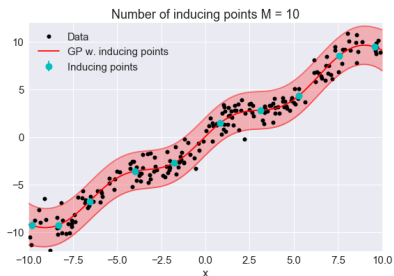
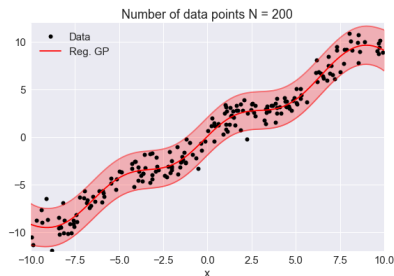


# Example: Number of inducing points



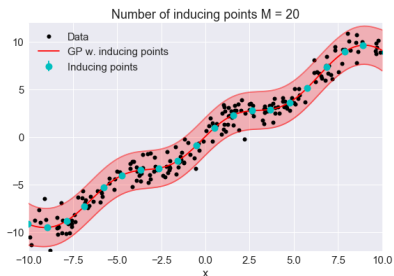
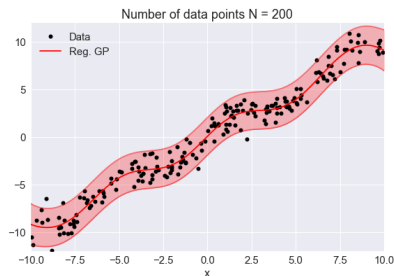
- We can think of the number of inducing points as a parameter that trades off speed for accuracy

# Example: Number of inducing points



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

# Example: Number of inducing points



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

# Gaussian process classification: Inference

Three steps to compute the predictive distribution for a new test point  $\mathbf{x}_*$

$$p(\mathbf{y}, \mathbf{f}) = \prod_{n=1}^N p(y_n | f_n) p(\mathbf{f}) = \prod_{n=1}^N \phi(y_n \cdot f_n) \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K})$$

- Step 1: Compute posterior distribution of  $p(\mathbf{f} | \mathbf{y})$ :

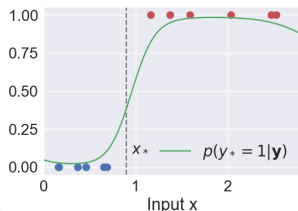
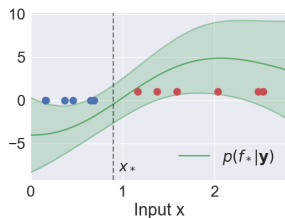
$$p(\mathbf{f} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f}) p(\mathbf{f})}{p(\mathbf{y})} \approx q(\mathbf{f})$$

- Step 2: Compute posterior of  $f_*$  for new test point  $\mathbf{x}_*$ :

$$p(f_* | \mathbf{y}) = \int p(f_* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \approx \int p(f_* | \mathbf{f}) q(\mathbf{f}) d\mathbf{f}$$

- Step 3: Compute predictive distribution

$$p(y_* | \mathbf{y}) = \int \phi(y_* \cdot f_*) p(f_* | \mathbf{y}) df_*$$



# Predictive distribution

- Using the (approximate) posterior  $q(f_*)$ , we can compute  $p(y_*|\mathbf{y})$

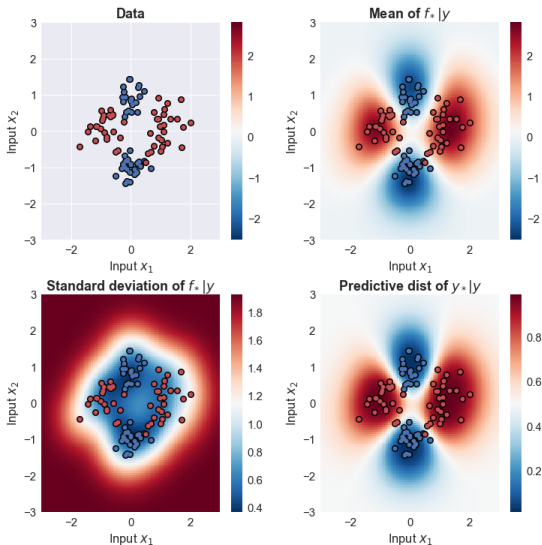
$$\begin{aligned} p(y_* = 1|\mathbf{y}) &= \int p(y_*|f_*)p(f_*|\mathbf{y})df_* \\ &= \int \phi(y_* \cdot f_*) p(f_*|\mathbf{y})df_* \\ &\approx \int \phi(y_* \cdot f_*) q(f_*) df_* \\ &= \int \phi(y_* \cdot f_*) \mathcal{N}(f_*|\mu_*, \sigma_*^2) df_* \\ &= \phi\left(\frac{\mu_*}{\sqrt{1 + \sigma_*^2}}\right) \end{aligned}$$

## Discuss with your neighbor

- What can we say about the predictive distributions for  $y_*$  when  $\mu_*$  is positive? or negative?
- How does the uncertainty of the posterior distribution of  $f_*$  influence the predictions for  $y_*$ ? What happens as  $\sigma_*^2$  approaches  $\infty$ ?

# Gaussian process classification example

- Non-linear classification problem
- $N = 100$  data points
- Squared exponential kernel
- Hyperparameters are chosen by optimizing  $\mathcal{L}_3$



# End of today's lecture

- This will be my last lecture
- Markus Heinonen, Arno Solin and Aki Vehtari will handle the rest of the course
- Next time: Markus Heinonen will give a lecture about spectral kernels
- In two weeks: Arno Solin will give a lecture about spatio-temporal modelling
- Now time for questions and assignment #2