

# Non-Gaussian likelihoods for Gaussian Processes

Alan Saul





# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

Expectation propagation

Comparisons

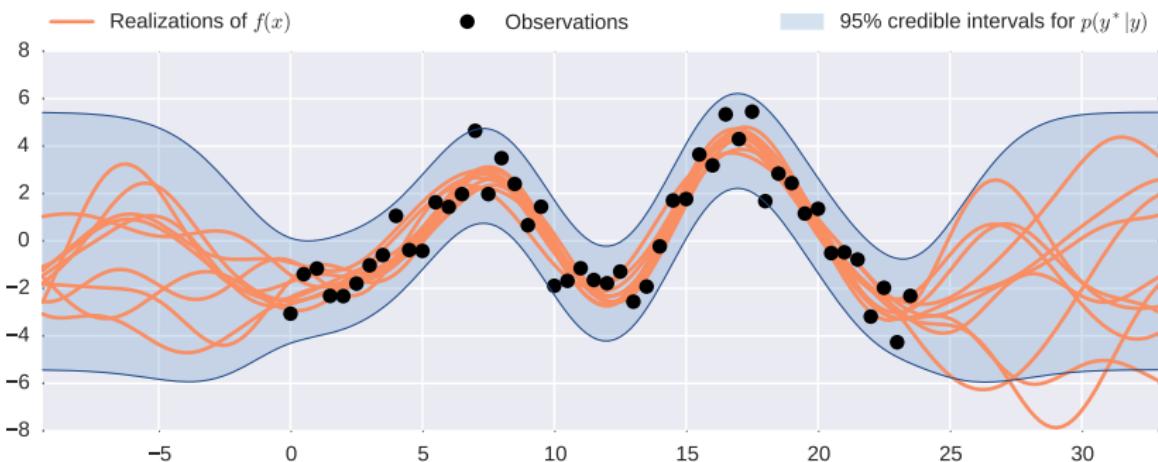
# GP regression - recap so far



Model the observations as a distorted version of the process  
 $\mathbf{f}_i = f(\mathbf{x}_i)$ :

$$\mathbf{y}_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma^2)$$

$f$  is a non-linear function, in our case we assume it is latent, and is assigned a Gaussian process prior.



# GP regression setting



So far we have assumed that the latent values,  $\mathbf{f}$ , have been corrupted by Gaussian noise. Everything remains analytically tractable.

Gaussian Prior:  $\mathbf{f} \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_{\mathbf{ff}}) = p(\mathbf{f})$

Gaussian likelihood:  $\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I}) = \prod_{i=1}^n p(\mathbf{y}_i | \mathbf{f}_i)$

Gaussian posterior:  $p(\mathbf{f} | \mathbf{y}) \propto \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K}_{\mathbf{ff}})$

# Outline



## Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

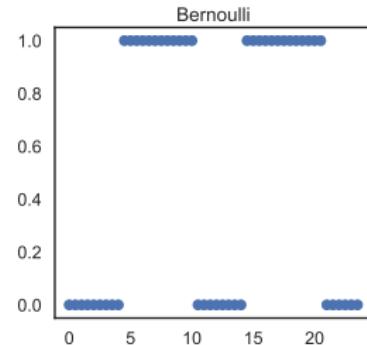
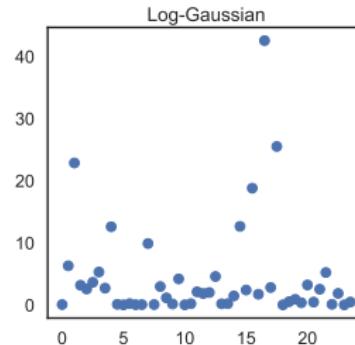
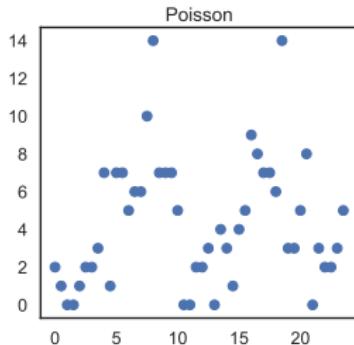
Expectation propagation

Comparisons



# Motivation

- ▶ You have been given some data you wish to model.
- ▶ You believe that the observations are connected through some underlying unknown function.
- ▶ You know from your understanding of the data generation process, that the observations are not Gaussian.
- ▶ You still want to learn, as best as possible, what is the unknown function being used, and make predictions.



# Likelihood



- ▶  $p(\mathbf{y}|\mathbf{f})$  is the probability that we would see some random variables,  $\mathbf{y}$ , if we knew the latent function values  $\mathbf{f}$ , which act as parameters.



- ▶  $p(\mathbf{y}|\mathbf{f})$  is the probability that we would see some random variables,  $\mathbf{y}$ , if we knew the latent function values  $\mathbf{f}$ , which act as parameters.
- ▶ Given the observed values for  $\mathbf{y}$  are fixed, it can also be seen as the likelihood that some latent function values,  $\mathbf{f}$ , would give rise to the observed values of  $\mathbf{y}$ . Note this is a *function* of  $\mathbf{f}$ , and doesn't integrate to 1 in  $\mathbf{f}$ .



- ▶  $p(\mathbf{y}|\mathbf{f})$  is the probability that we would see some random variables,  $\mathbf{y}$ , if we knew the latent function values  $\mathbf{f}$ , which act as parameters.
- ▶ Given the observed values for  $\mathbf{y}$  are fixed, it can also be seen as the likelihood that some latent function values,  $\mathbf{f}$ , would give rise to the observed values of  $\mathbf{y}$ . Note this is a *function* of  $\mathbf{f}$ , and doesn't integrate to 1 in  $\mathbf{f}$ .
- ▶ Often observations aren't observed by simple Gaussian corruptions of the underlying latent function,  $\mathbf{f}$ .

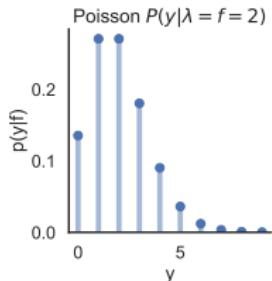
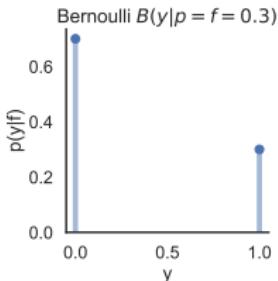
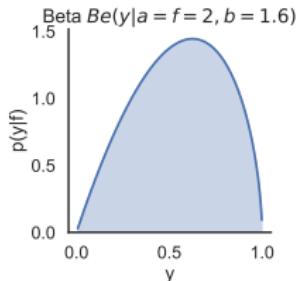
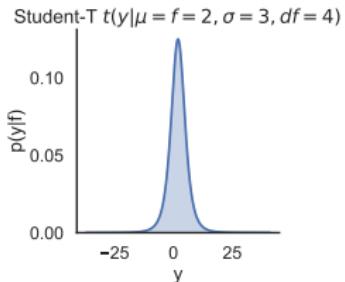
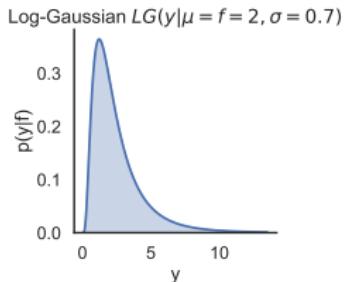
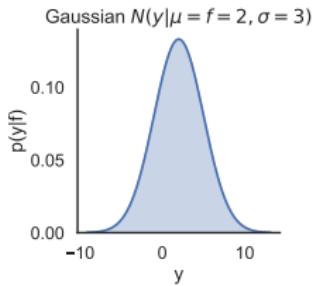


- ▶  $p(\mathbf{y}|\mathbf{f})$  is the probability that we would see some random variables,  $\mathbf{y}$ , if we knew the latent function values  $\mathbf{f}$ , which act as parameters.
- ▶ Given the observed values for  $\mathbf{y}$  are fixed, it can also be seen as the likelihood that some latent function values,  $\mathbf{f}$ , would give rise to the observed values of  $\mathbf{y}$ . Note this is a *function* of  $\mathbf{f}$ , and doesn't integrate to 1 in  $\mathbf{f}$ .
- ▶ Often observations aren't observed by simple Gaussian corruptions of the underlying latent function,  $\mathbf{f}$ .
- ▶ In the case of count data, binary data, etc, we need to choose a different likelihood function.

# Likelihood



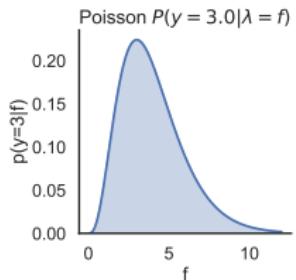
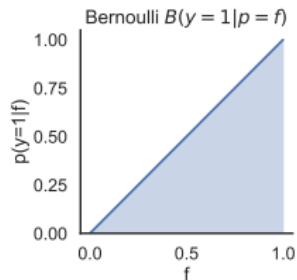
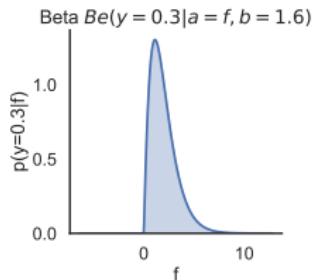
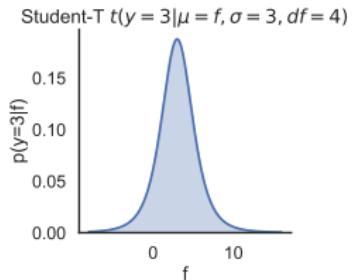
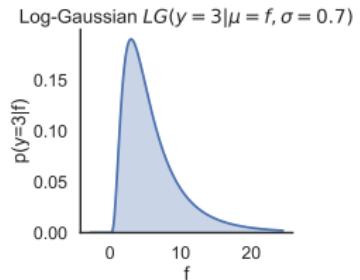
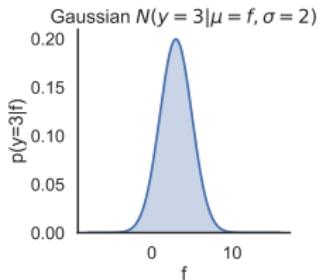
$p(y|f)$  as a function of  $y$ , with fixed  $f$



# Likelihood



$p(y|f)$  as a function of  $f$ , with fixed  $y$

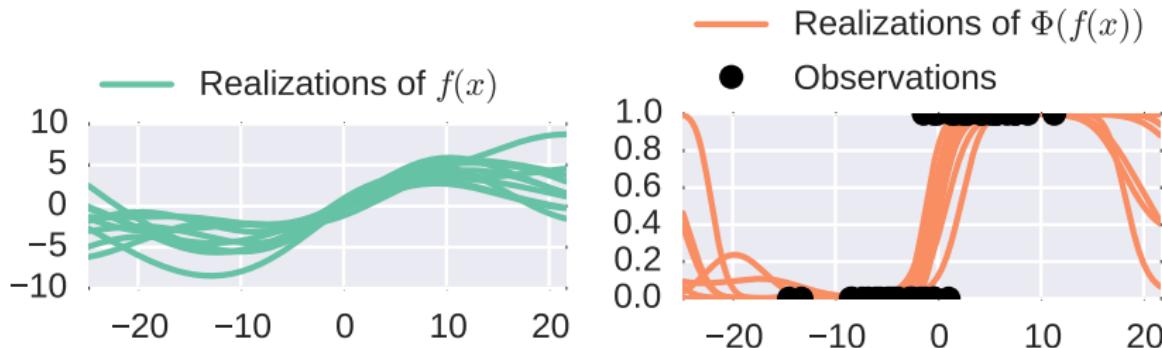




# Binary example

- ▶ Binary outcomes for  $\mathbf{y}_i, \mathbf{y}_i \in [0, 1]$ .
- ▶ Model the probability of  $\mathbf{y}_i = 1$  with transformation of GP, with Bernoulli likelihood.
- ▶ Probability of 1 must be between 0 and 1, thus use squashing transformation,  $\lambda(\mathbf{f}_i) = \Phi(\mathbf{f}_i)$ .

$$p(\mathbf{y}_i | \lambda(\mathbf{f}_i)) = \begin{cases} \lambda(\mathbf{f}_i), & \text{if } \mathbf{y}_i = 1 \\ 1 - \lambda(\mathbf{f}_i), & \text{if } \mathbf{y}_i = 0 \end{cases}$$

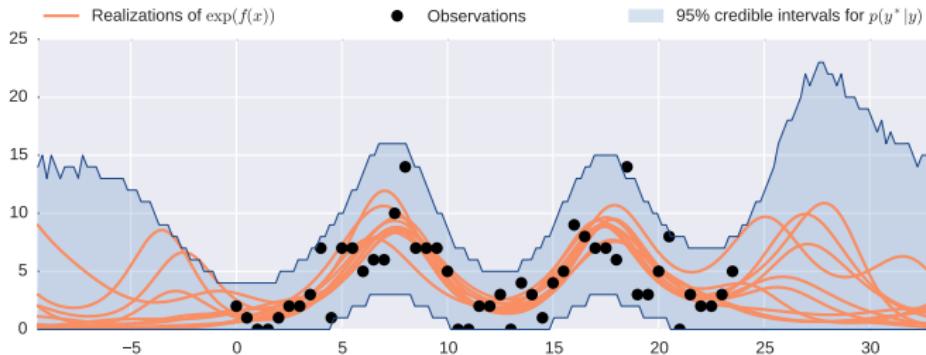




# Count data example

- ▶ Non-negative and discrete values only for  $y_i, y_i \in \mathbb{N}$ .
- ▶ Model the *rate* or *intensity*,  $\lambda$ , of events with a transformation of a Gaussian process.
- ▶ Rate parameter must remain positive, use transformation to maintain positiveness  $\lambda(f_i) = \exp(f_i)$  or  $\lambda(f_i) = f_i^2$

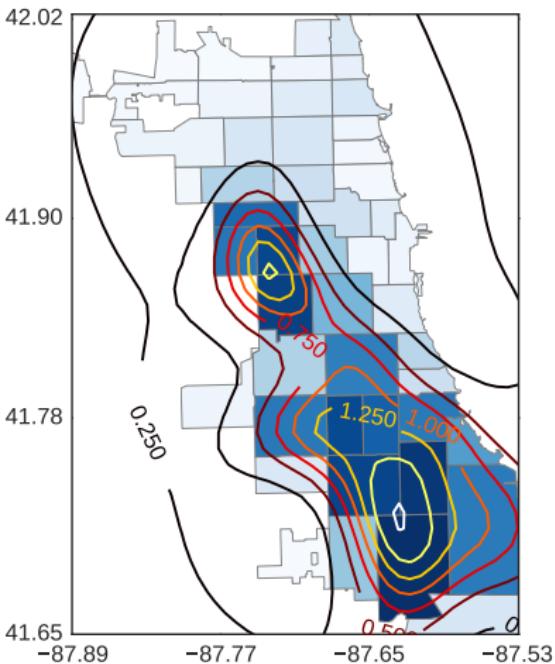
$$y_i \sim \text{Poisson}(y_i | \lambda_i = \lambda(f_i)) \quad \text{Poisson}(y_i | \lambda_i) = \frac{\lambda_i^{y_i}}{y_i!} e^{-\lambda_i}$$





# Application example

- ▶ Chicago crime counts.
- ▶ Same Poisson likelihood.
- ▶ 2D-input to kernel.





# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

Expectation propagation

Comparisons



# Non-Gaussian posteriors

- ▶ Exact computation of posterior is no longer analytically tractable due to non-conjugate Gaussian process prior to non-Gaussian likelihood,  $p(\mathbf{y}|\mathbf{f})$ .

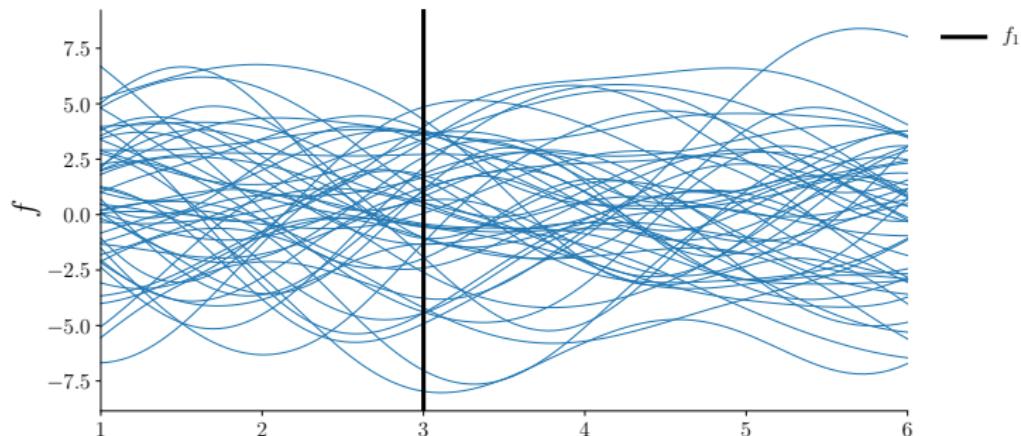
$$p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{f}) \prod_{i=1}^n p(\mathbf{y}_i|\mathbf{f}_i)}{\int p(\mathbf{f}) \prod_{i=1}^n p(\mathbf{y}_i|\mathbf{f}_i) d\mathbf{f}}$$

Why is it so difficult?

# Non-Gaussian posteriors illustrated



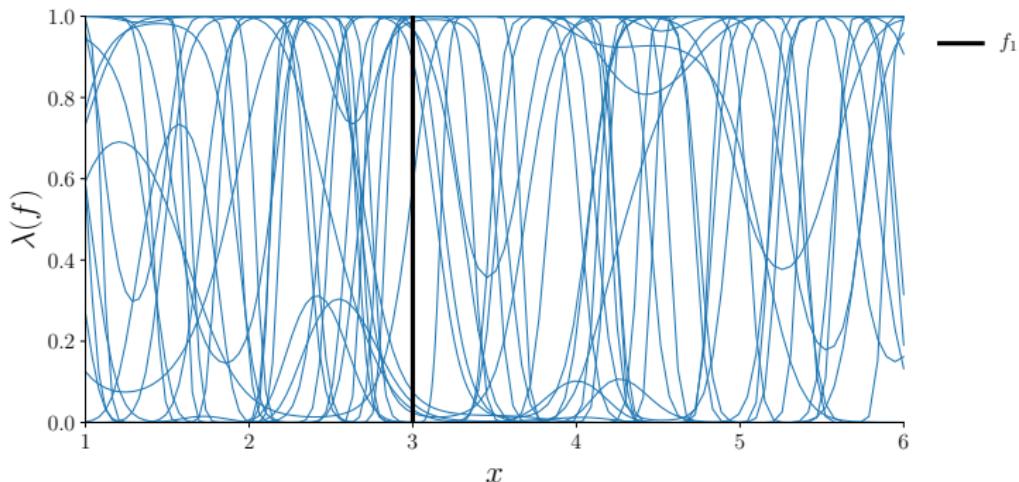
- ▶ Consider one observation,  $y_1 = 1$ , at input  $x_1$ .
- ▶ Can normalise easily with numerical integration,  
 $\int p(y_1 = 1|\lambda(f_1))p(f_1)df_1.$



# Non-Gaussian posteriors illustrated



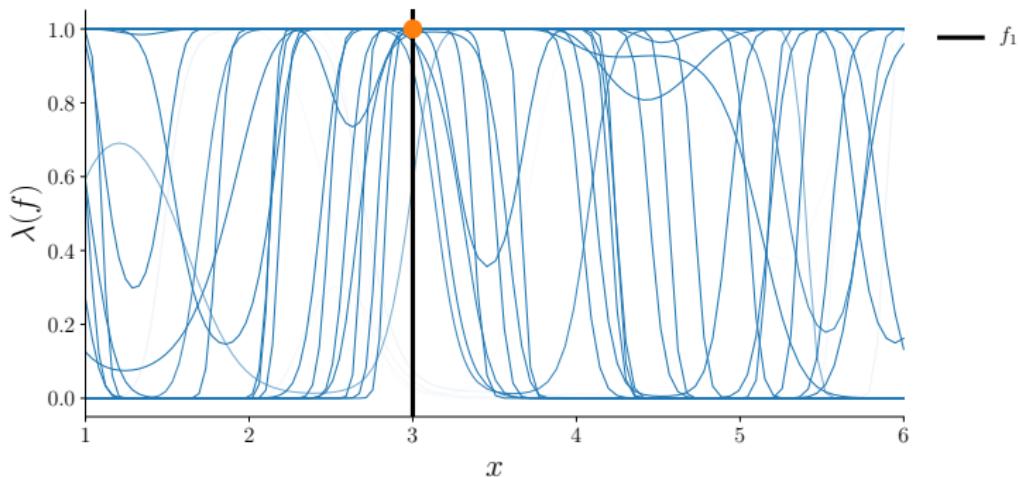
- ▶ Consider one observation,  $y_1 = 1$ , at input  $x_1$ .
- ▶ Can normalise easily with numerical integration,  
 $\int p(y_1 = 1|\lambda(f_1))p(f_1)df_1.$



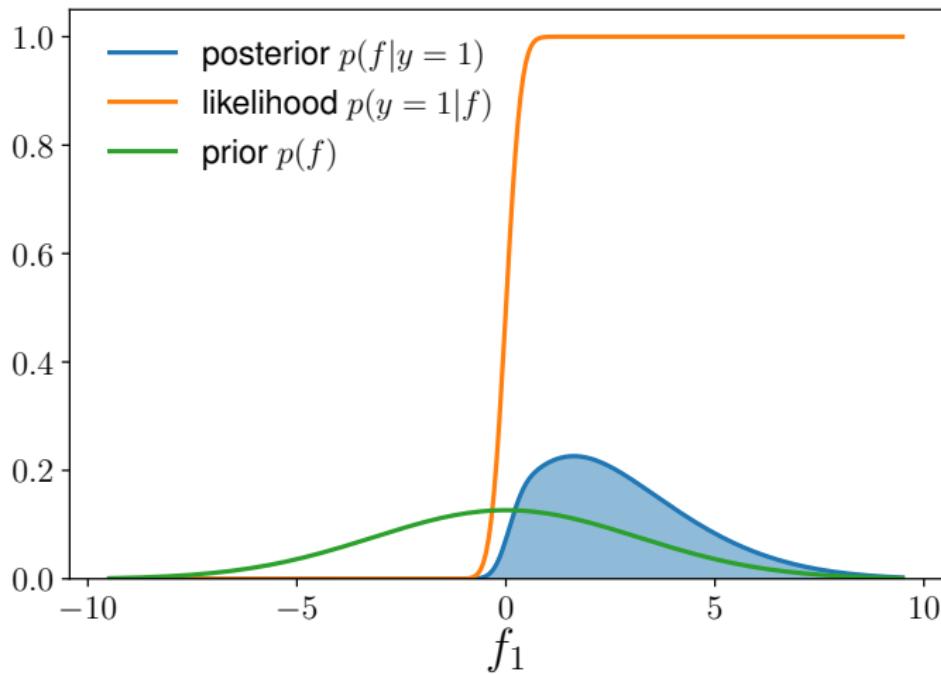
# Non-Gaussian posteriors illustrated



- ▶ Consider one observation,  $y_1 = 1$ , at input  $x_1$ .
- ▶ Can normalise easily with numerical integration,  
 $\int p(y_1 = 1|\lambda(f_1))p(f_1)df_1.$



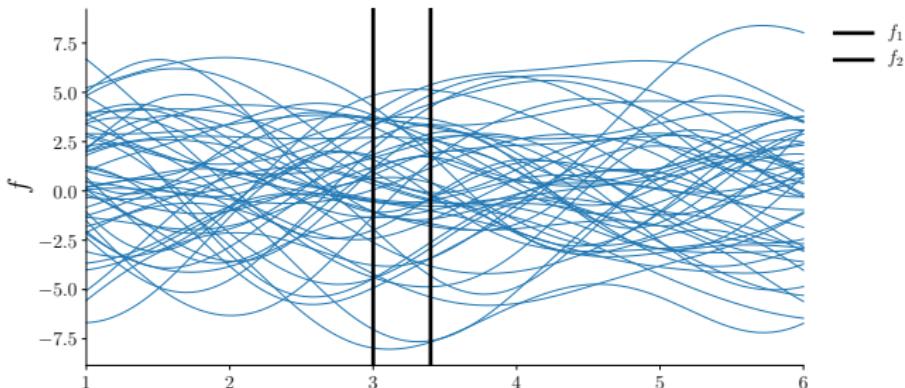
# Non-Gaussian posteriors illustrated





# Non-Gaussian posteriors illustrated

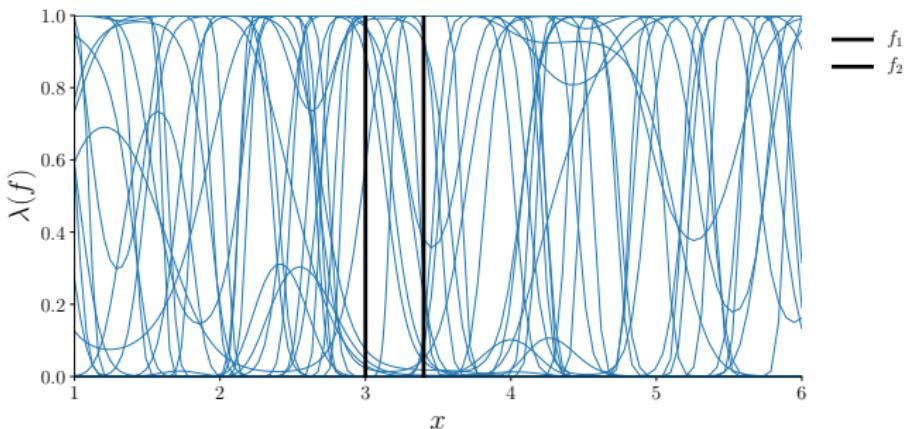
- Now two observations,  $y_1 = 1$  and  $y_2 = 1$  at  $x_1$  and  $x_2$
- Need to calculate the joint posterior,  
 $p(\mathbf{f}|\mathbf{y}) = p(f_1, f_2|y_1 = 1, y_2 = 1)$ .
- Requires 2D integral  
 $\int \int p(y_1 = 1, y_2 = 1 | \lambda(f_1), \lambda(f_2)) p(f_1, f_2) df_1 df_2.$





# Non-Gaussian posteriors illustrated

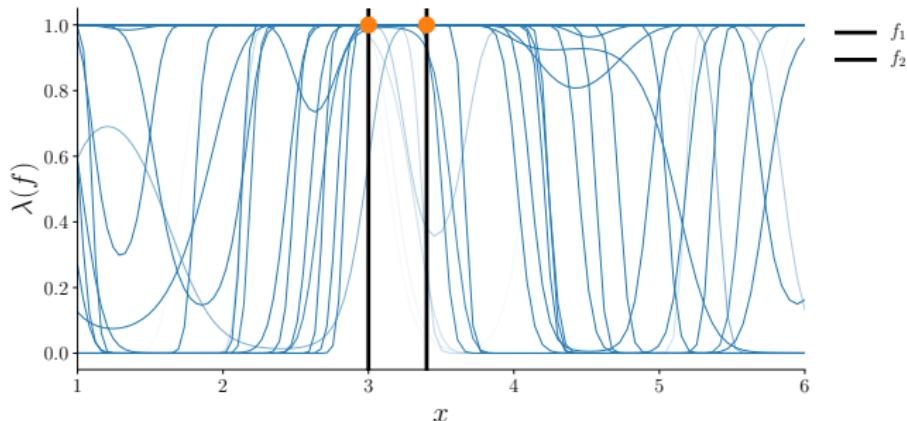
- Now two observations,  $y_1 = 1$  and  $y_2 = 1$  at  $x_1$  and  $x_2$
- Need to calculate the joint posterior,  
 $p(\mathbf{f}|\mathbf{y}) = p(f_1, f_2|y_1 = 1, y_2 = 1)$ .
- Requires 2D integral  
 $\int \int p(y_1 = 1, y_2 = 1 | \lambda(f_1), \lambda(f_2)) p(f_1, f_2) df_1 df_2.$





# Non-Gaussian posteriors illustrated

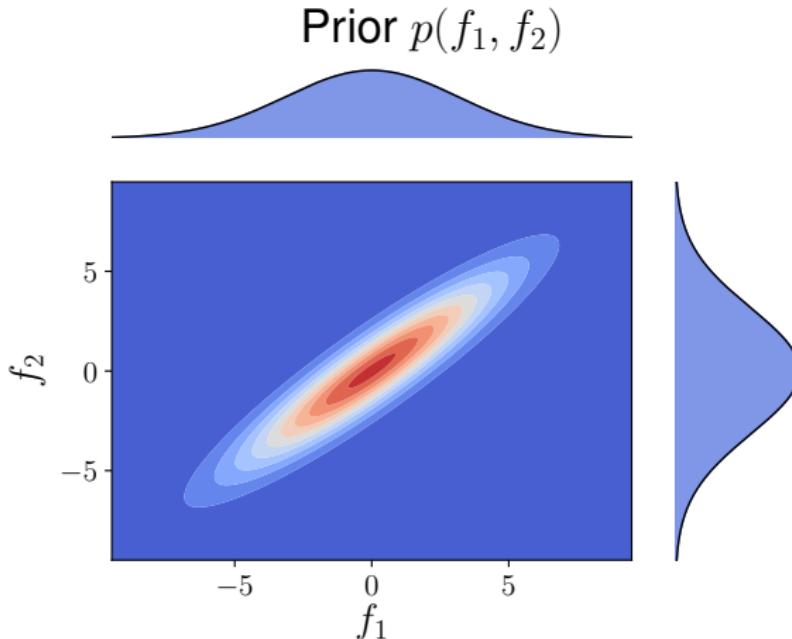
- Now two observations,  $y_1 = 1$  and  $y_2 = 1$  at  $x_1$  and  $x_2$
- Need to calculate the joint posterior,  
 $p(\mathbf{f}|\mathbf{y}) = p(f_1, f_2|y_1 = 1, y_2 = 1)$ .
- Requires 2D integral  
 $\int \int p(y_1 = 1, y_2 = 1 | \lambda(f_1), \lambda(f_2)) p(f_1, f_2) df_1 df_2.$



# Non-Gaussian posteriors illustrated



- ▶ To find the true posterior values, we need to perform a two dimensional integral.
- ▶ Still possible, but things are getting more difficult quickly.

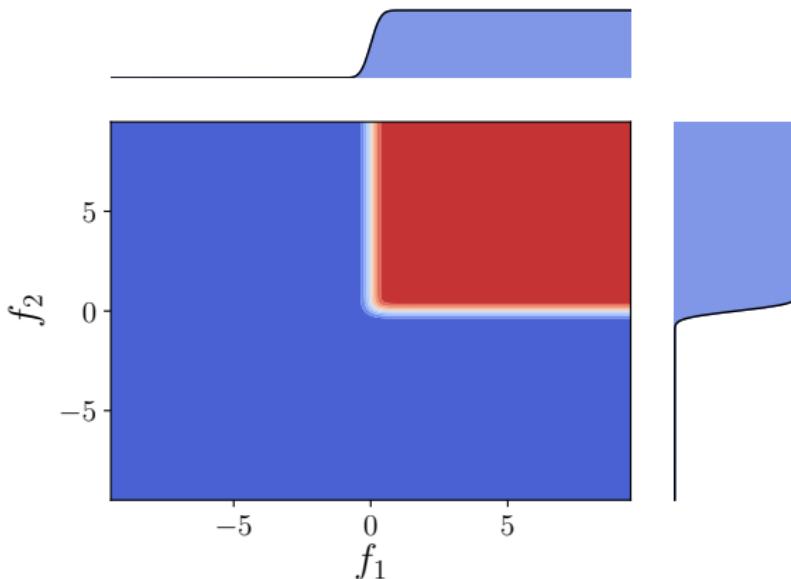


# Non-Gaussian posteriors illustrated



- ▶ To find the true posterior values, we need to perform a two dimensional integral.
- ▶ Still possible, but things are getting more difficult quickly.

Likelihood  $p(y_1 = 1, y_2 = 1 | f_1, f_2)$

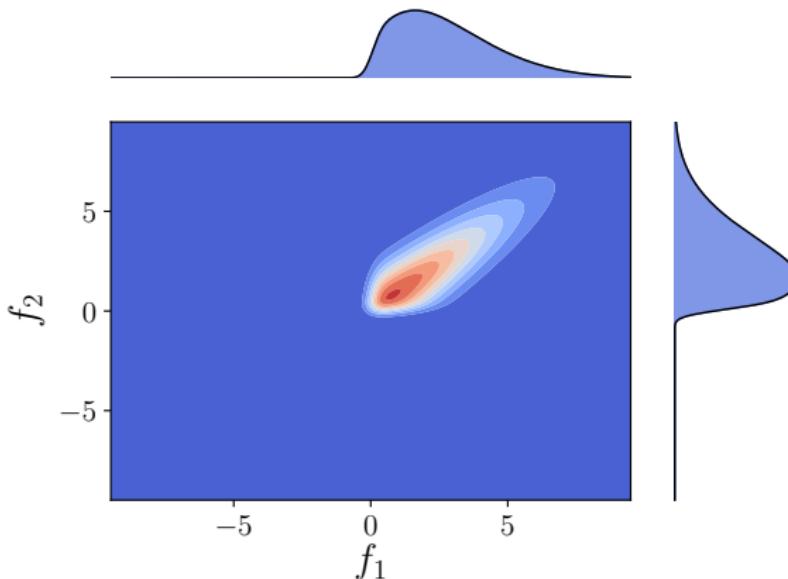


# Non-Gaussian posteriors illustrated



- ▶ To find the true posterior values, we need to perform a two dimensional integral.
- ▶ Still possible, but things are getting more difficult quickly.

True posterior  $p(f_1, f_2 | y_1 = 1, y_2 = 1)$



# Approaches to handling non-Gaussian posteriors



Generally fall into two areas:

- ▶ Sampling methods that obtain samples of the posterior.
- ▶ Approximation of the posterior with something of known form.

Today we will focus on the latter.



# Non-Gaussian posterior approximation



- ▶ Various methods to make a Gaussian approximation,  
 $p(\mathbf{f}|\mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu} = ?, \mathbf{C} = ?)$ .
- ▶ Only need to obtain an approximate posterior at the training locations.
- ▶ At test locations, the data only effects their probability via the posterior at these locations.

$$p(\mathbf{f}, \mathbf{f}^* | \mathbf{x}^*, \mathbf{x}, \mathbf{y}) = p(\mathbf{f}^* | \mathbf{f}, \mathbf{x}^*) p(\mathbf{f} | \mathbf{x}, \mathbf{y})$$

# Why do we want an the posterior anyway?



True posterior, posterior approximation, or samples are needed to make predictions at new locations,  $\mathbf{x}^*$ .

$$p(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) = \int p(\mathbf{f}^*|\mathbf{f}, \mathbf{x}^*) p(\mathbf{f}|\mathbf{y}, \mathbf{x}) d\mathbf{f}$$
$$q(\mathbf{f}^*|\mathbf{x}^*, \mathbf{x}, \mathbf{y}) = \int p(\mathbf{f}^*|\mathbf{f}, \mathbf{x}^*) q(\mathbf{f}|\mathbf{x}) d\mathbf{f}$$



# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

Expectation propagation

Comparisons

# Methods overview



Given choice of Gaussian approximation of posterior. How do we choose the parameter values  $\mu$  and  $C$ ?

There a number of different methods in which to choose how to set the parameters of our Gaussian approximation.



# Parameters effect - mean



# Parameters effect - variance



# How to choose the parameters?

Two approaches that we might take:

- ▶ Is to match the mean and variance at some point, for example the mode.
  - ▶ Attempt to minimise some divergence measure between the approximate distribution and the true distribution.
- 
- ▶ Laplace takes the former
  - ▶ Variational bayes takes the latter
  - ▶ EP kind of takes the latter



# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

Expectation propagation

Comparisons



# Laplace approximation

Task: for some generic random variable,  $z$ , and data,  $y$ , find a good approximation to difficult to compute posterior distribution,  $p(z|y)$ .

Laplace approach: fit a Gaussian by matching the curvature at the modal point of the posterior.

- ▶ Use a second-order taylor expansion around the mode of the log-posterior.
- ▶ Use the expansion to find an equivalent Gaussian in the probability space.



# Laplace approximation

- ▶ Log of a Gaussian distribution,  $q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \mathbf{C})$ , is a quadratic function of  $\mathbf{f}$ .
- ▶ A second-order taylor expansion is an approximation of a function using only quadratic terms.
- ▶ Laplace approximation expands the un-normalised posterior, and then uses it to set the linear and quadratic terms of the log  $q(\mathbf{f})$ .
- ▶ The first and second derivatives of the form of the log-posterior, at the mode, will match the derivatives of the approximate Gaussian at this same point.



## Second-order taylor expansion

$$p(\mathbf{f}|\mathbf{y}) = \frac{1}{Z} h(\mathbf{f})$$



# Second-order taylor expansion

$$p(\mathbf{f}|\mathbf{y}) = \frac{1}{Z} h(\mathbf{f})$$

$$\log p(\mathbf{f}|\mathbf{y}) = \log \frac{1}{Z} + \log h(\mathbf{f})$$



# Second-order taylor expansion

$$p(\mathbf{f}|\mathbf{y}) = \frac{1}{Z} h(\mathbf{f})$$

$$\begin{aligned}\log p(\mathbf{f}|\mathbf{y}) &= \log \frac{1}{Z} + \log h(\mathbf{f}) \\ &\approx \log \frac{1}{Z} + \log h(\mathbf{a}) + \frac{d \log h(\mathbf{a})}{d\mathbf{a}} (\mathbf{f} - \mathbf{a}) \\ &\quad + \frac{1}{2} (\mathbf{f} - \mathbf{a})^\top \frac{d^2 \log h(\mathbf{a})}{d\mathbf{a}^2} (\mathbf{f} - \mathbf{a})\end{aligned}$$



## Second-order taylor expansion

$$p(\mathbf{f}|\mathbf{y}) = \frac{1}{Z} h(\mathbf{f})$$

$$\begin{aligned}\log p(\mathbf{f}|\mathbf{y}) &= \log \frac{1}{Z} + \log h(\mathbf{f}) \\ &\approx \log \frac{1}{Z} + \log h(\mathbf{a}) + \frac{d \log h(\mathbf{a})}{d\mathbf{a}} (\mathbf{f} - \mathbf{a}) \\ &\quad + \frac{1}{2} (\mathbf{f} - \mathbf{a})^\top \frac{d^2 \log h(\mathbf{a})}{d\mathbf{a}^2} (\mathbf{f} - \mathbf{a})\end{aligned}$$

In our case we want to make our expansion around the mode,  $\hat{\mathbf{f}}$ :

$$\left. \frac{d \log h(\mathbf{a})}{d\mathbf{a}} \right|_{\mathbf{a}=\hat{\mathbf{f}}} = 0$$



# Second-order taylor expansion

$$\begin{aligned}\log p(\mathbf{f}|\mathbf{y}) &\approx \log \frac{1}{Z} + \log h(\hat{\mathbf{f}}) + \frac{d \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}} (\mathbf{f} - \hat{\mathbf{f}}) \\ &+ \frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} (\mathbf{f} - \hat{\mathbf{f}})\end{aligned}$$

# Second-order taylor expansion



$$\begin{aligned}\log p(\mathbf{f}|\mathbf{y}) &\approx \log \frac{1}{Z} + \log h(\hat{\mathbf{f}}) + \frac{d \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}} (\mathbf{f} - \hat{\mathbf{f}}) \\ &\quad + \frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} (\mathbf{f} - \hat{\mathbf{f}}) \\ p(\mathbf{f}|\mathbf{y}) &= \frac{1}{Z} h(\hat{\mathbf{f}}) \exp \left\{ -\frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \left( -\frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} \right) (\mathbf{f} - \hat{\mathbf{f}}) \right\}\end{aligned}$$



## Second-order taylor expansion

$$\begin{aligned}\log p(\mathbf{f}|\mathbf{y}) &\approx \log \frac{1}{Z} + \log h(\hat{\mathbf{f}}) + \frac{d \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}} (\mathbf{f} - \hat{\mathbf{f}}) \\ &\quad + \frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} (\mathbf{f} - \hat{\mathbf{f}}) \\ p(\mathbf{f}|\mathbf{y}) &= \frac{1}{Z} h(\hat{\mathbf{f}}) \exp \left\{ -\frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \left( -\frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} \right) (\mathbf{f} - \hat{\mathbf{f}}) \right\} \\ &= \mathcal{N} \left( \mathbf{f} | \hat{\mathbf{f}}, \left( -\frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} \right)^{-1} \right)\end{aligned}$$

# Laplace appoximation for Gaussian processes



In our case,  $h(\mathbf{f}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f})$ , so we need to evaluate

$$\begin{aligned}-\frac{d^2 \log h(\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} &= -\frac{d^2(\log p(\mathbf{y}|\hat{\mathbf{f}}) + \log p(\hat{\mathbf{f}}))}{d\hat{\mathbf{f}}^2} \\&= -\frac{d^2 \log p(\mathbf{y}|\hat{\mathbf{f}})}{d\hat{\mathbf{f}}^2} + \mathbf{K}^{-1} \\&\triangleq \mathbf{W} + \mathbf{K}^{-1}\end{aligned}$$

giving a posterior approximation:

$$p(\mathbf{f}|\mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}\left(\mathbf{f}|\hat{\mathbf{f}}, (\mathbf{W} + \mathbf{K}^{-1})^{-1}\right)$$

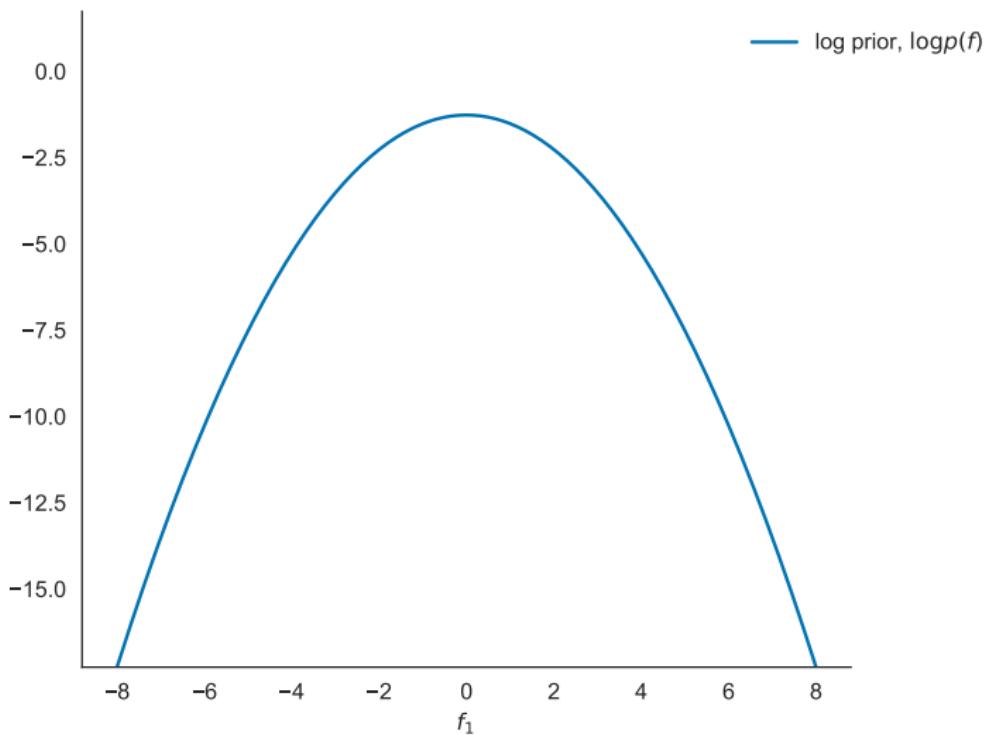
# Laplace approximation - algorithm overview



- ▶ Find the mode,  $\hat{\mathbf{f}}$  of the true log posterior, via Newton's method.
- ▶ Use second-order Taylor expansion around this modal value.
- ▶ Form Gaussian approximation setting the mean equal to the posterior mode,  $\hat{\mathbf{f}}$ , and matching the curvature.
- ▶  $p(\mathbf{f}|\mathbf{y}) \approx q(\mathbf{f}|\boldsymbol{\mu}, \mathbf{C}) = \mathcal{N}\left(\mathbf{f}|\hat{\mathbf{f}}, (\mathbf{K}^{-1} + \mathbf{W})^{-1}\right)$
- ▶  $\mathbf{W} \triangleq -\frac{d^2 \log p(\mathbf{y}|\mathbf{f})}{d\hat{\mathbf{f}}^2}$ .
- ▶ For factorizing likelihoods (most),  $\mathbf{W}$  is diagonal.

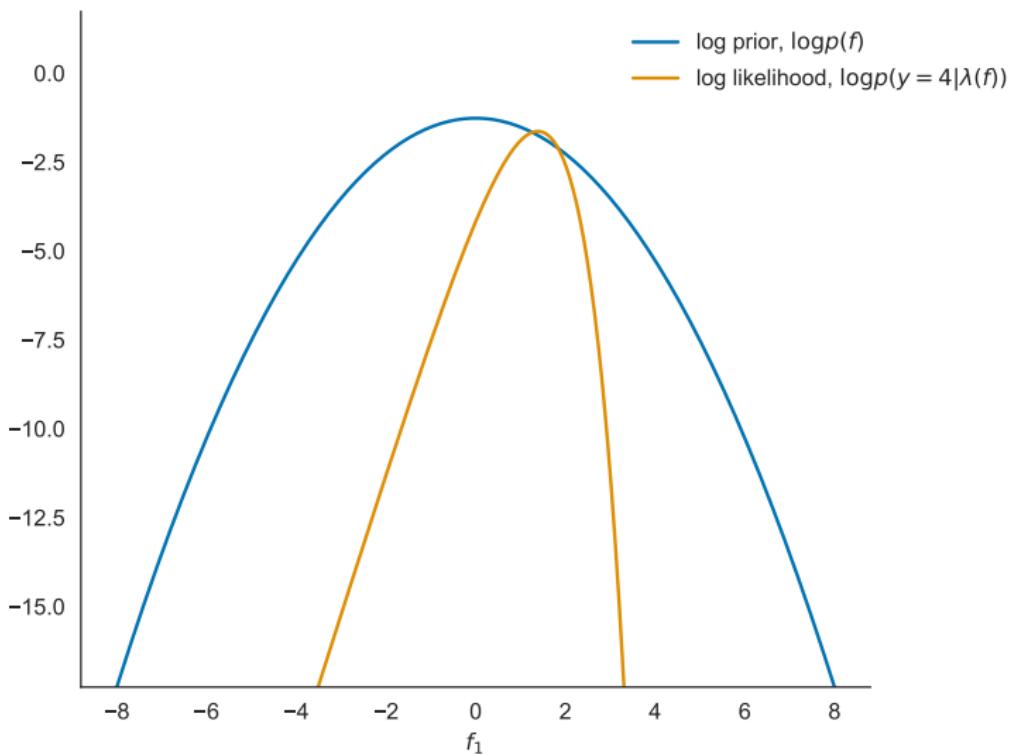


# Visualization of Laplace



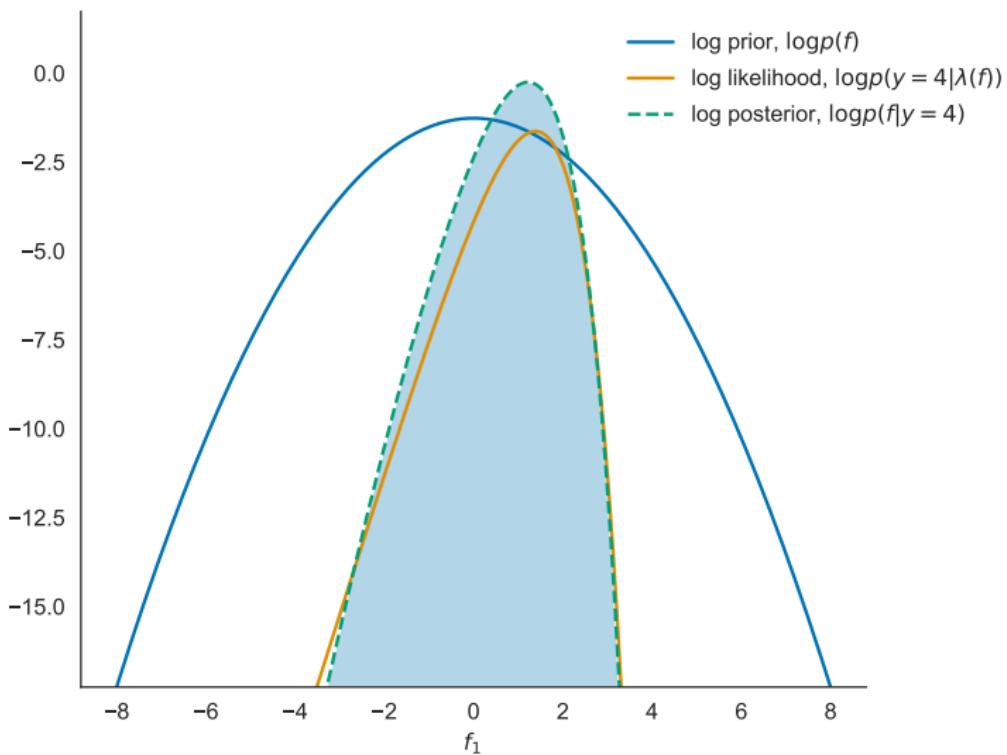


# Visualization of Laplace



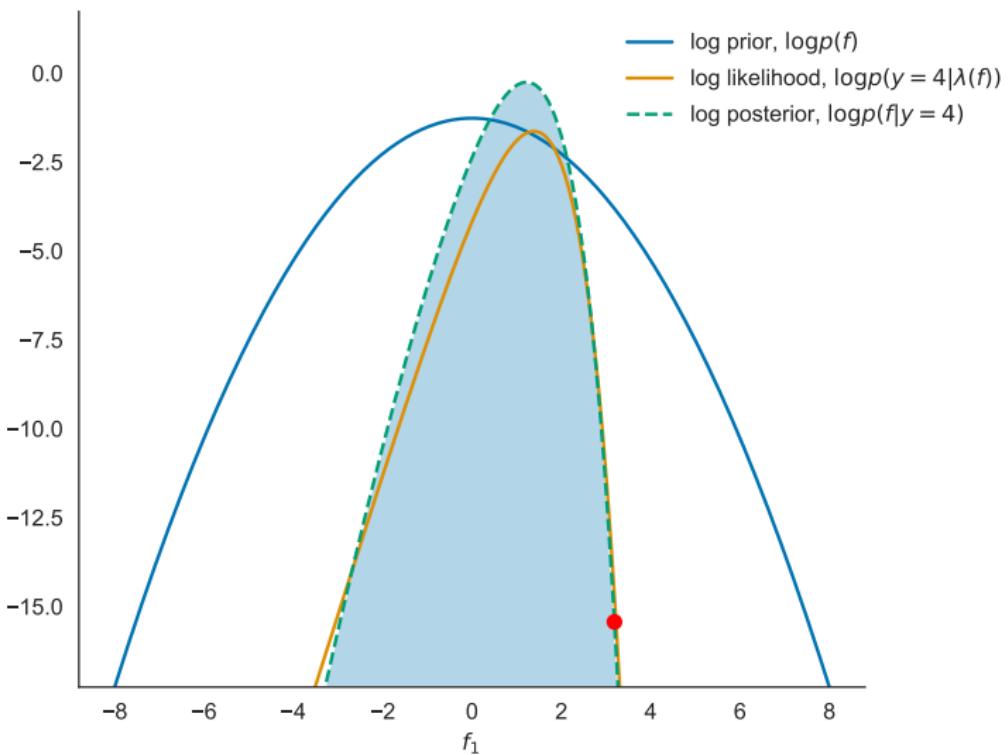


# Visualization of Laplace



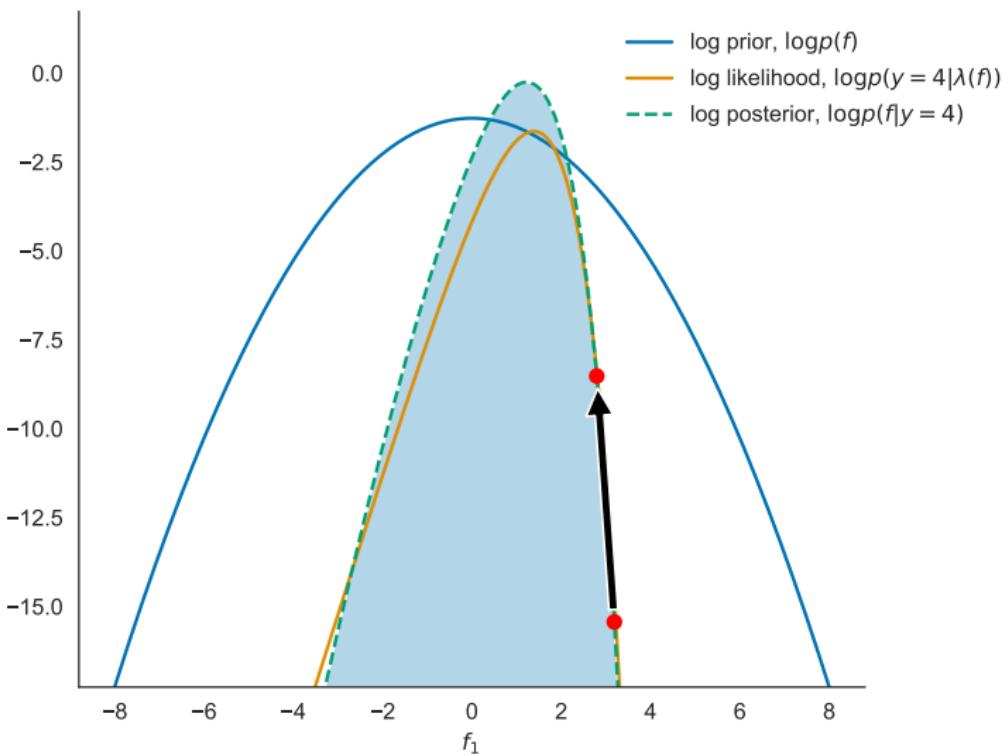


# Visualization of Laplace



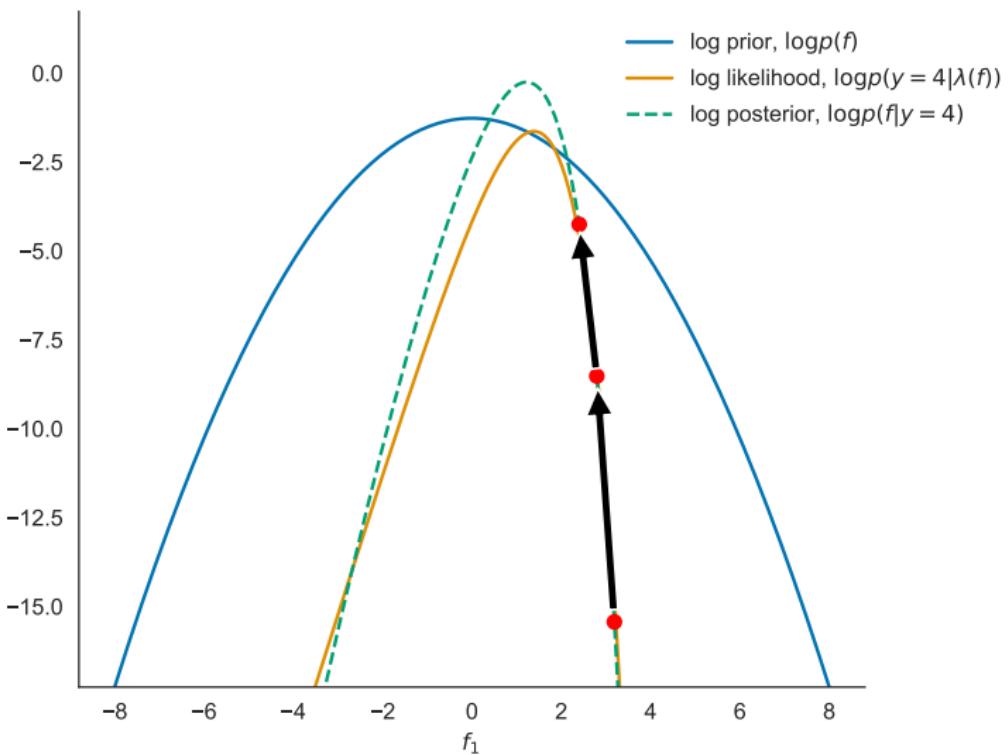


# Visualization of Laplace



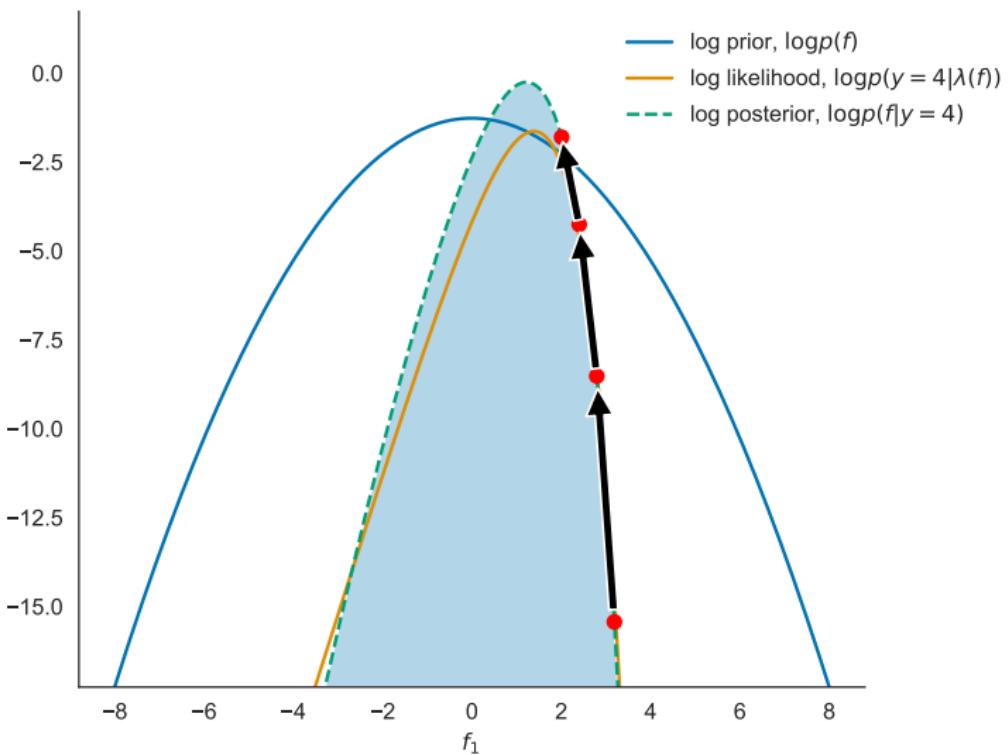


# Visualization of Laplace



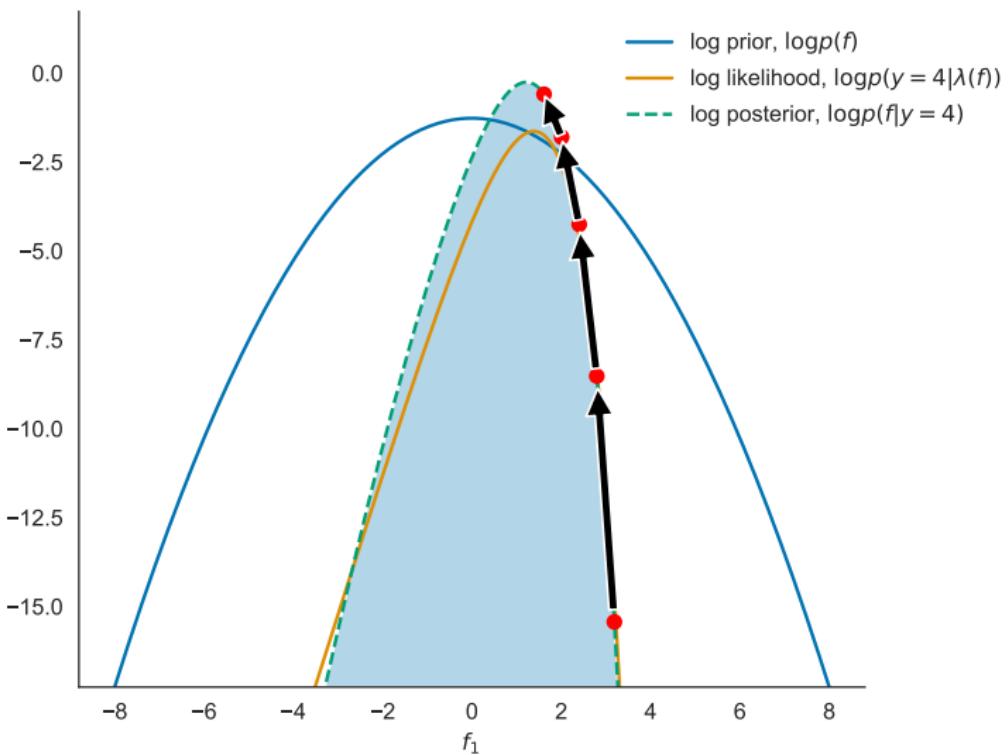


# Visualization of Laplace



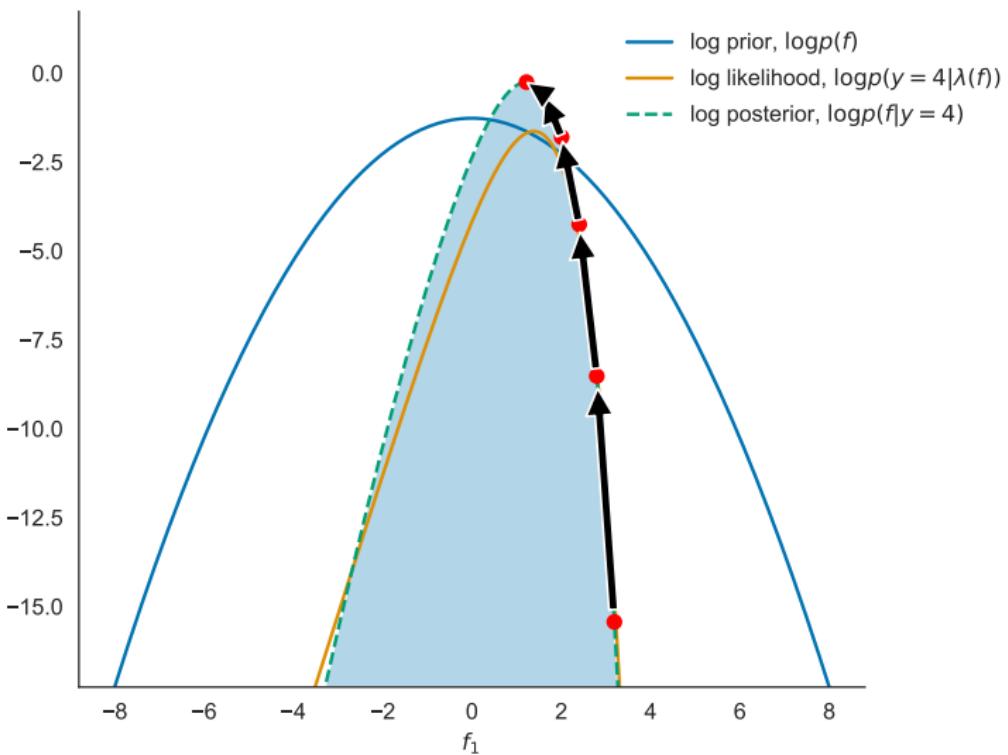


# Visualization of Laplace



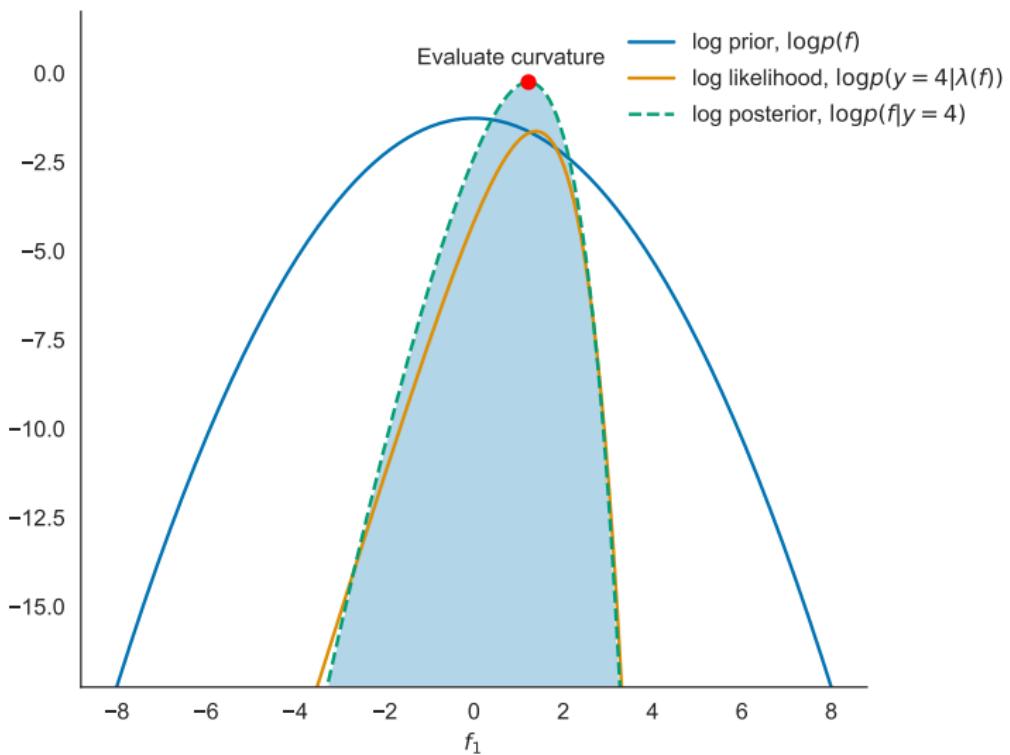


# Visualization of Laplace



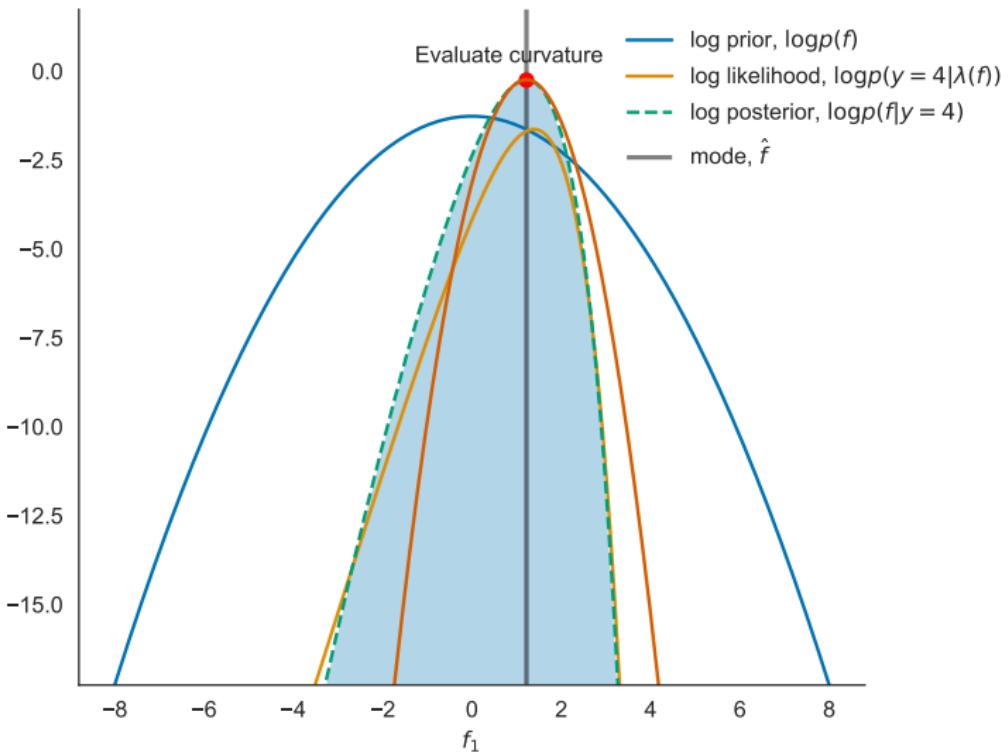


# Visualization of Laplace



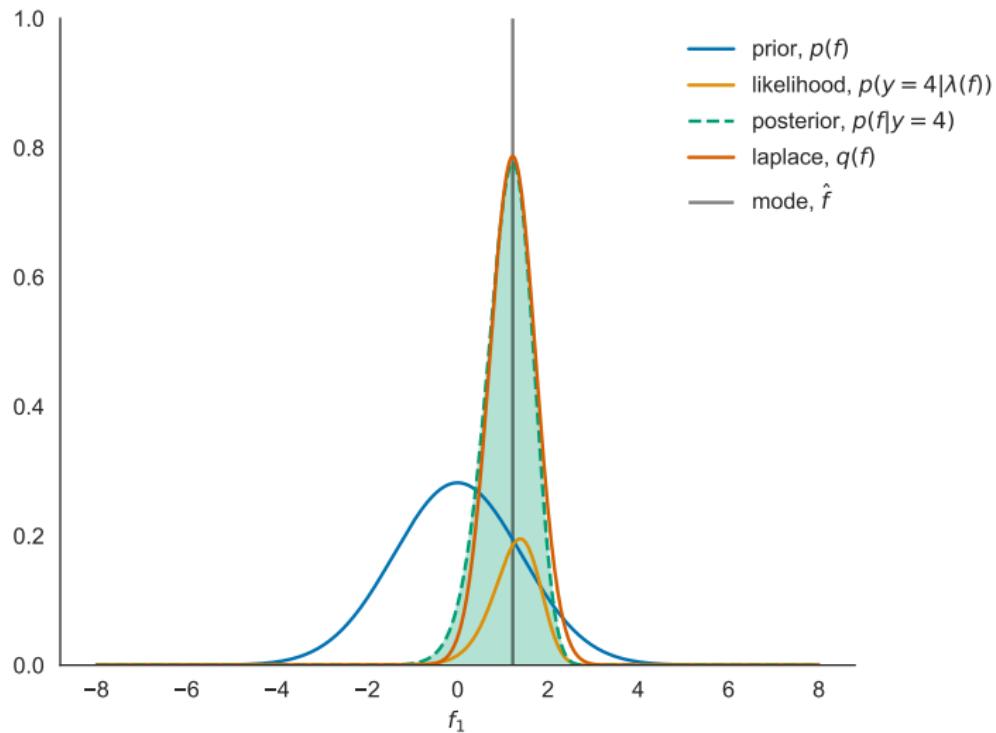


# Visualization of Laplace



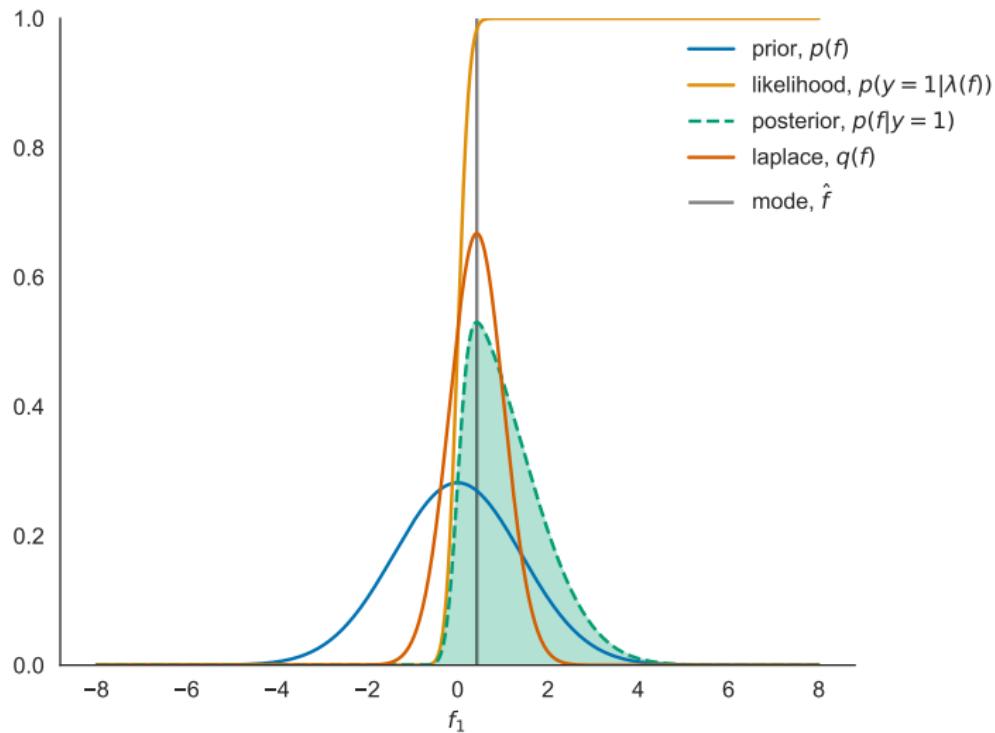


# Visualization of Laplace





# Visualise of Laplace - Bernoulli





# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

Expectation propagation

Comparisons

# Variational Bayes (VB)



Task: for some generic random variable,  $z$ , and data,  $y$ , find a good approximation to difficult to compute posterior distribution,  $p(z|y)$ .

VB approach: minimise a divergence measure between an approximate posterior,  $q(z)$  and true posterior,  $p(z|y)$ .

- ▶ KL divergence,  $\text{KL}(q(z) \parallel p(z|y))$ .
- ▶ Minimize this with respect to parameters of  $q(z)$ .

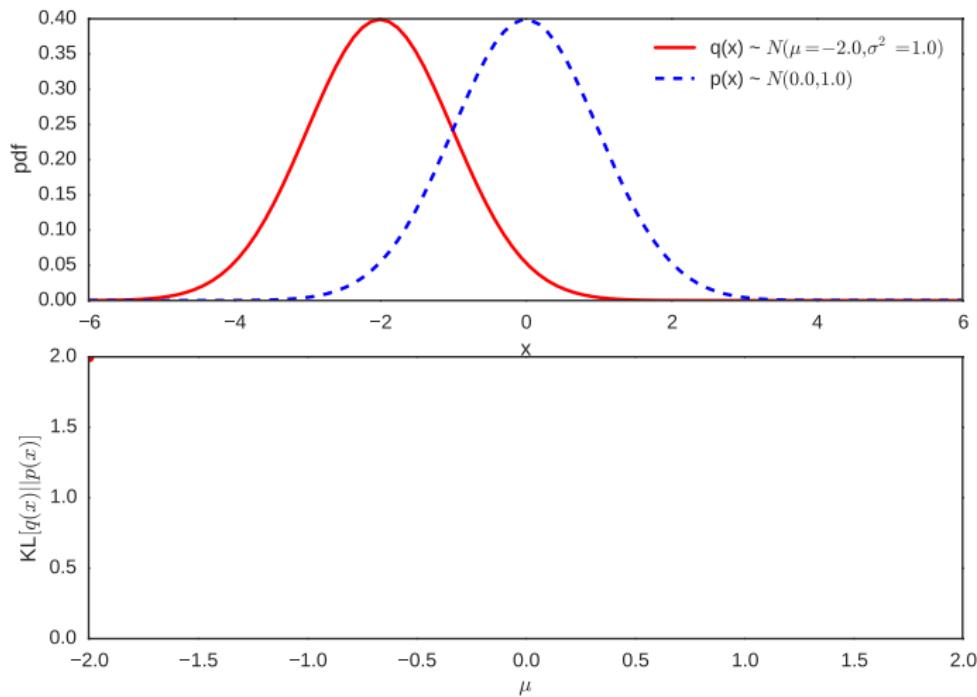
# KL divergence



- ▶ General for any two distributions  $q(\mathbf{x})$  and  $p(\mathbf{x})$ .
- ▶  $\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x}))$  is the average additional amount of information lost when  $p(\mathbf{x})$  is used to approximate  $q(\mathbf{x})$ . It's a measure of divergence of one distribution to another.
- ▶  $\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \left\langle \log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\rangle_{q(\mathbf{x})}$
- ▶ Always 0 or positive, not symmetric.
- ▶ Lets look at how it changes with response to changes in the approximating distribution.

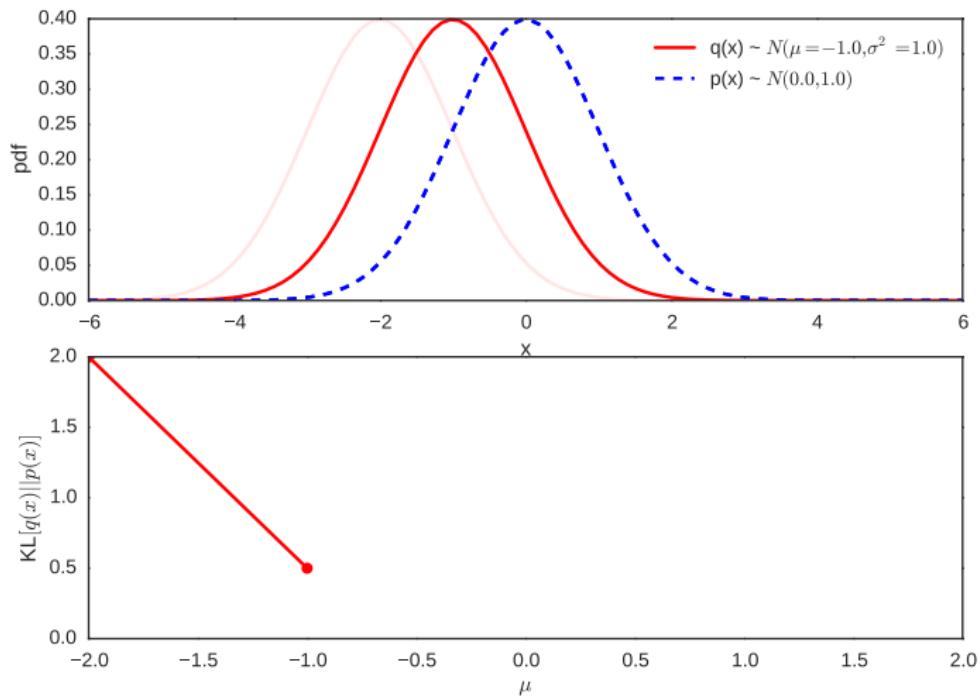


# KL varying mean



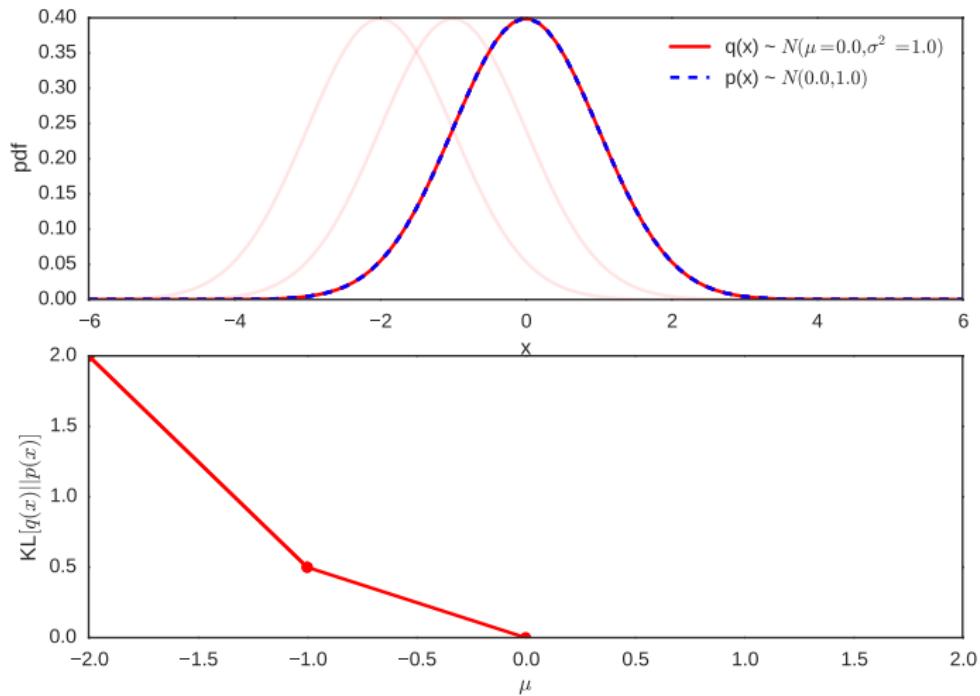


# KL varying mean



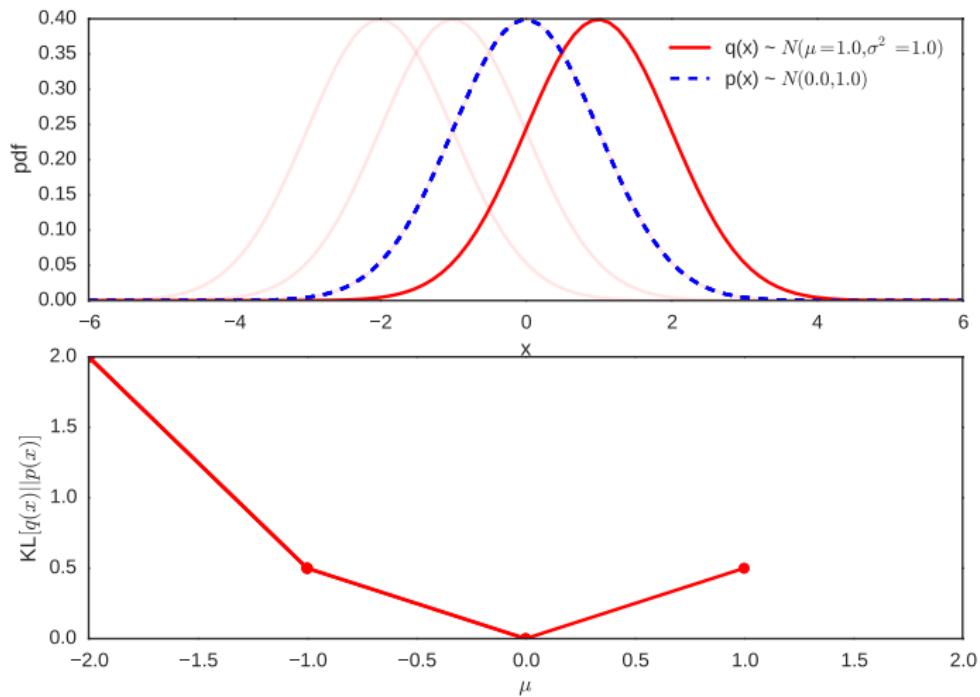


# KL varying mean



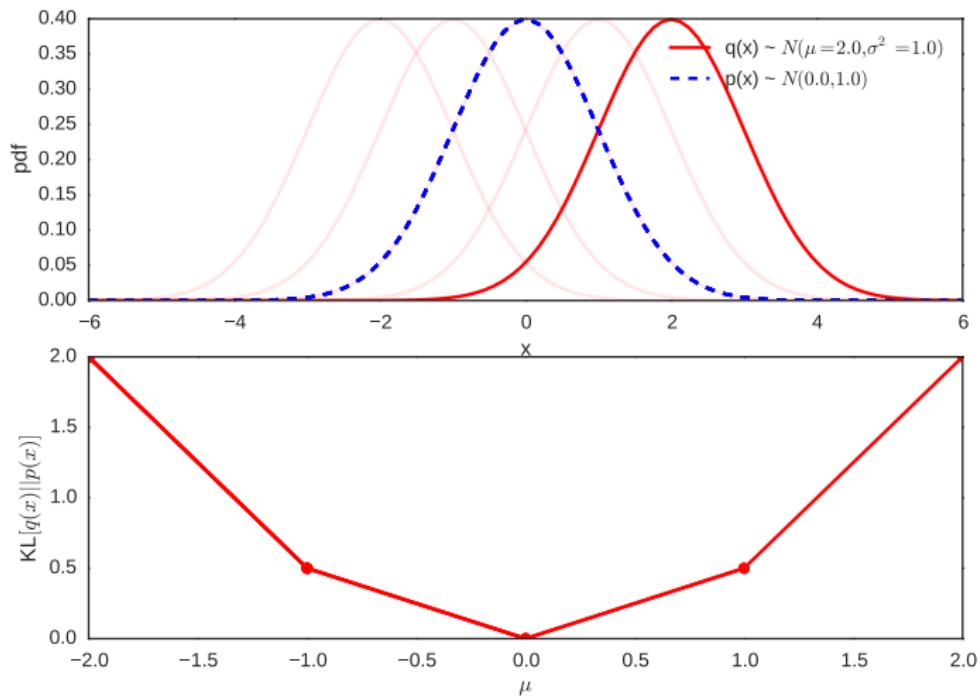


# KL varying mean



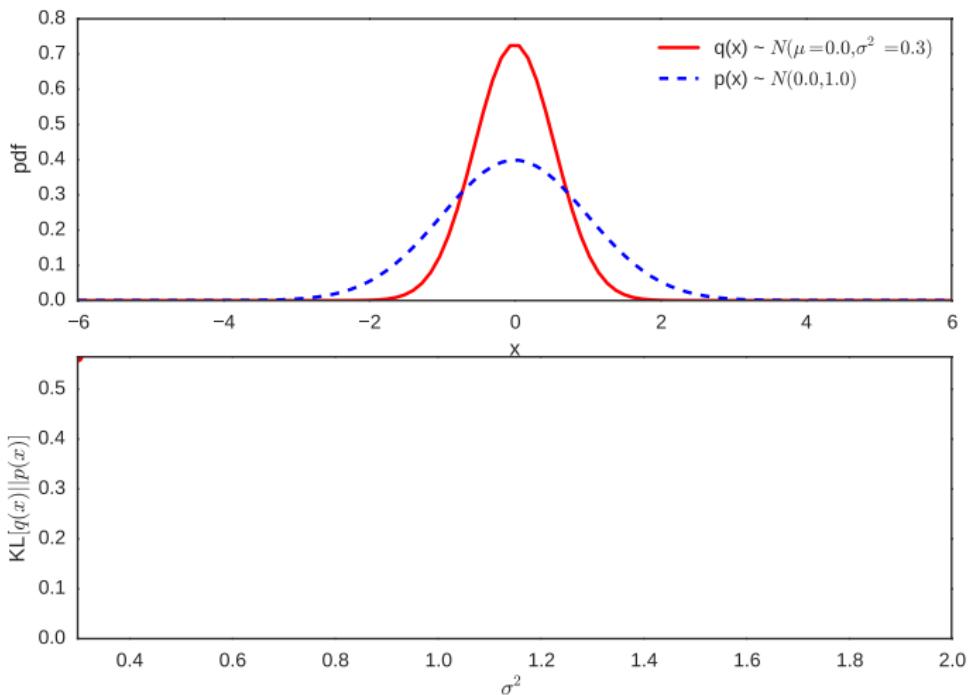


# KL varying mean



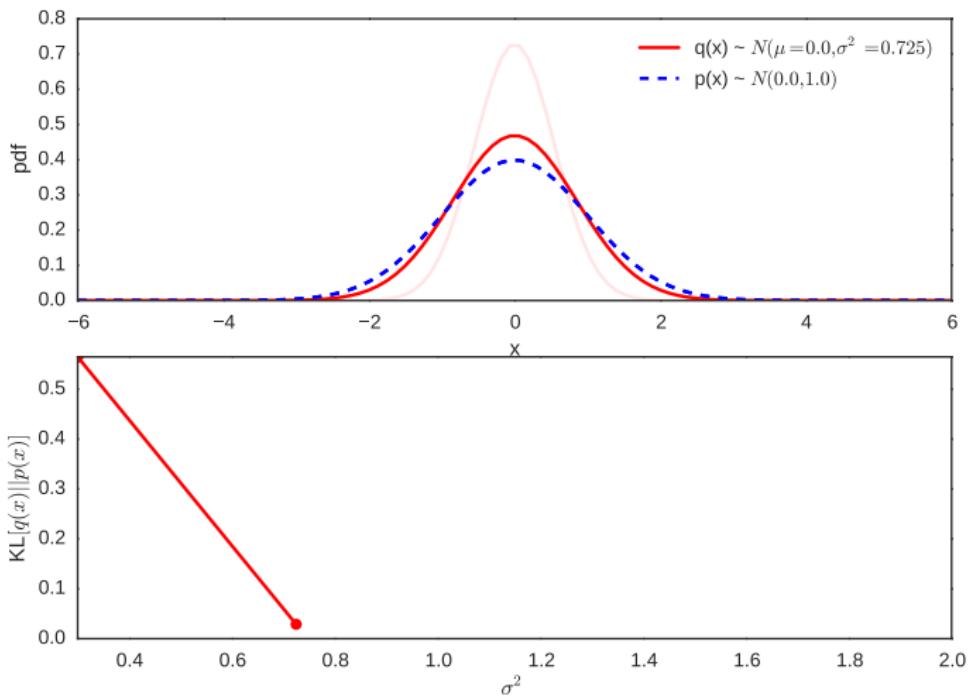


# KL varying variance



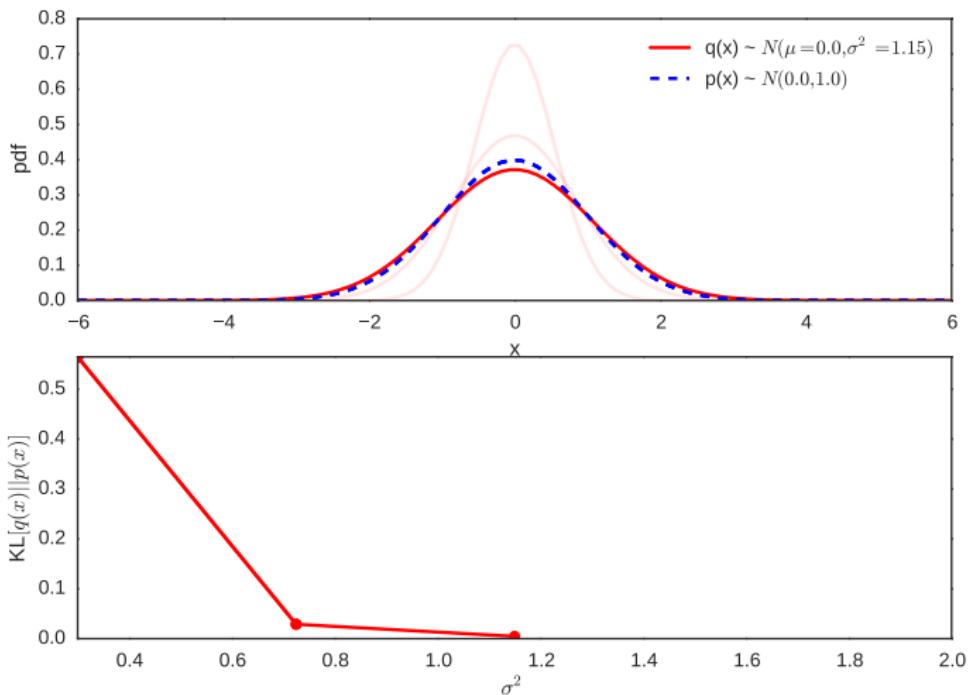


# KL varying variance



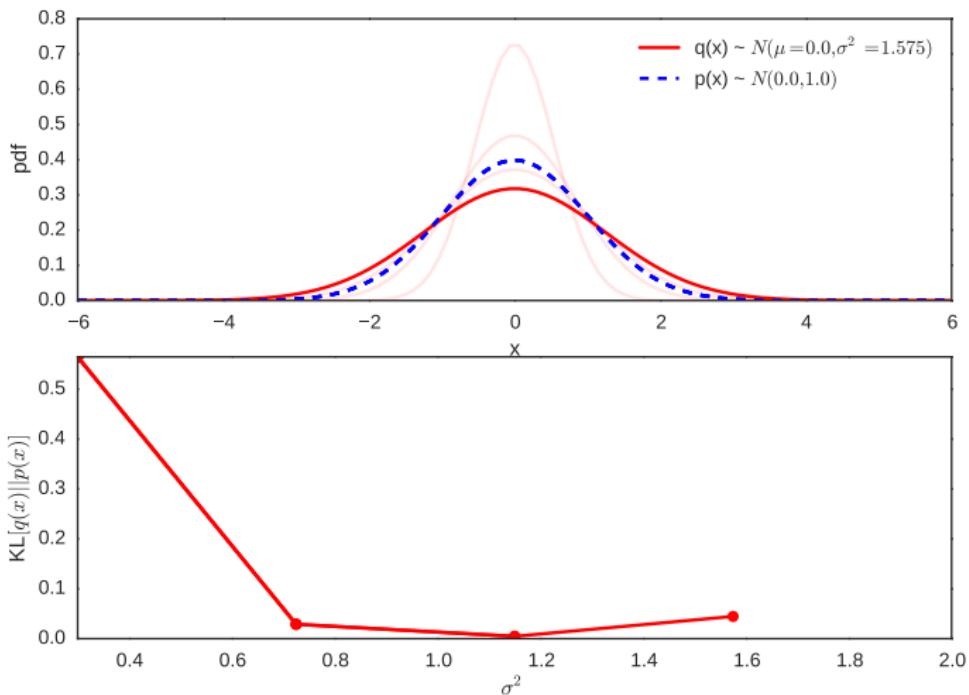


# KL varying variance



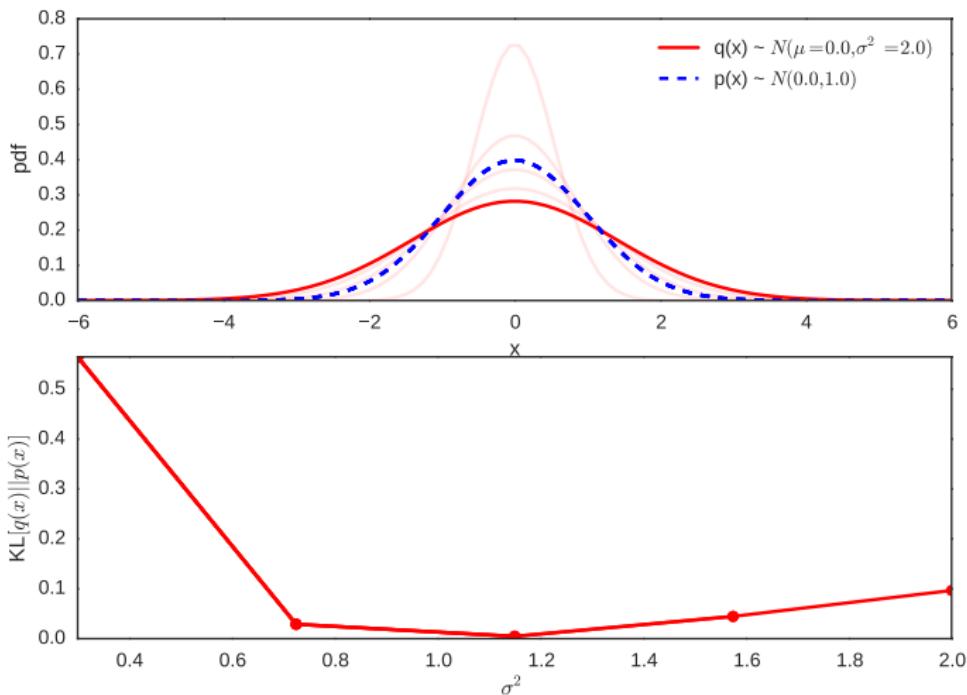


# KL varying variance





# KL varying variance



# Variational Bayes



Don't have access to or can't compute for computational reasons:  $p(z|y)$  or  $p(y)$ , and hence  $\text{KL}(q(z) \parallel p(z|y))$

How can we minimize something we can't compute?

- ▶ Can compute  $q(z)$  and  $p(y|z)$  for any  $z$ .
- ▶  $q(z)$  is parameterised by 'variational parameters'.
- ▶ True posterior using Bayes rule,  $p(z|y) = \frac{p(y|z)p(z)}{p(y)}$ .
- ▶  $p(y)$  doesn't change when variational parameters are changed.



# Variational Bayes - Derivation

$$\text{KL}(q(z) \parallel p(z|y))$$

# Variational Bayes - Derivation



$$\text{KL}(q(z) \parallel p(z|y)) \\ = \int q(z) \left[ \log \frac{q(z)}{p(z|y)} \right] dz$$



# Variational Bayes - Derivation

$$\text{KL}(q(z) \parallel p(z|y))$$

$$= \int q(z) \left[ \log \frac{q(z)}{p(z|y)} \right] dz$$

$$= \int q(z) \left[ \log \frac{q(z)}{p(z)} - \log p(y|z) + \log p(y) \right] dz$$

# Variational Bayes - Derivation



$$\begin{aligned} & \text{KL}(q(z) \parallel p(z|y)) \\ &= \int q(z) \left[ \log \frac{q(z)}{p(z|y)} \right] dz \\ &= \int q(z) \left[ \log \frac{q(z)}{p(z)} - \log p(y|z) + \log p(y) \right] dz \\ &= \text{KL}(q(z) \parallel p(z)) - \int q(z) [\log p(y|z)] dz + \log p(y) \end{aligned}$$

# Variational Bayes - Derivation



$$\text{KL}(q(z) \parallel p(z|y))$$

$$= \int q(z) \left[ \log \frac{q(z)}{p(z|y)} \right] dz$$

$$= \int q(z) \left[ \log \frac{q(z)}{p(z)} - \log p(y|z) + \log p(y) \right] dz$$

$$= \text{KL}(q(z) \parallel p(z)) - \int q(z) [\log p(y|z)] dz + \log p(y)$$

$$\log p(y) = \int q(z) [\log p(y|z)] dz - \text{KL}(q(z) \parallel p(z)) + \text{KL}(q(z) \parallel p(z|y))$$



# Variational Bayes - Derivation

$$\begin{aligned}\log p(y) &= \int q(z) [\log p(y|z)] dz - \text{KL}(q(z) \| p(z)) + \text{KL}(q(z) \| p(z|y)) \\ &\geq \int q(z) [\log p(y|z)] dz - \text{KL}(q(z) \| p(z))\end{aligned}$$

- ▶ Tractable terms give lower bound on  $\log p(y)$  as  $\text{KL}(q(z) \| p(z|y))$  always positive.
- ▶ Adjust variational parameters of  $q(z)$  to make tractable terms as large as possible, thus  $\text{KL}(q(z) \| p(z|y))$  as small as possible.

# VB optimisation illustration



# Variational Bayes for Gaussian processes



- ▶ Make a Gaussian approximation,  $q(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \mathbf{C})$ , as similar possible to true posterior,  $p(\mathbf{f}|\mathbf{y})$ .
- ▶ Treat  $\boldsymbol{\mu}$  and  $\mathbf{C}$  as ‘variational parameters’, effecting quality of approximation.

$$\begin{aligned}\text{KL}(q(\mathbf{f}) \parallel p(\mathbf{f}|\mathbf{y})) &= \left\langle \log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right\rangle_{q(\mathbf{f})} \\ &= \left\langle \log \frac{q(\mathbf{f})}{p(\mathbf{f})} - \log p(\mathbf{y}|\mathbf{f}) + \log p(\mathbf{y}) \right\rangle_{q(\mathbf{f})} \\ &= \text{KL}(q(\mathbf{f}) \parallel p(\mathbf{f})) - \langle \log p(\mathbf{y}|\mathbf{f}) \rangle_{q(\mathbf{f})} + \log p(\mathbf{y}) \\ \log p(\mathbf{y}) &= \langle \log p(\mathbf{y}|\mathbf{f}) \rangle_{q(\mathbf{f})} - \text{KL}(q(\mathbf{f}) \parallel p(\mathbf{f})) + \text{KL}(q(\mathbf{f}) \parallel p(\mathbf{f}|\mathbf{y}))\end{aligned}$$



$$\begin{aligned}\log p(\mathbf{y}) &= \langle \log p(\mathbf{y}|\mathbf{f}) \rangle_{q(\mathbf{f})} - \text{KL}(q(\mathbf{f}) \| p(\mathbf{f})) + \text{KL}(q(\mathbf{f}) \| p(\mathbf{f}|\mathbf{y})) \\ &\geq \langle \log p(\mathbf{y}|\mathbf{f}) \rangle_{q(\mathbf{f})} - \text{KL}(q(\mathbf{f}) \| p(\mathbf{f}))\end{aligned}$$

- ▶ Adjust variational parameters  $\mu$  and  $C$  to make tractable terms as large as possible, thus  $\text{KL}(q(\mathbf{f}) \| p(\mathbf{f}|\mathbf{y}))$  as small as possible.
- ▶  $\langle \log p(\mathbf{y}|\mathbf{f}) \rangle_{q(\mathbf{f})}$  with factorizing likelihood can be done with a series of  $n$  1 dimensional integrals.
- ▶ In practice, can reduce the number of variational parameters by reparameterizing  $C = (\mathbf{K}_{\mathbf{ff}} - 2\Lambda)^{-1}$  by noting that the bound is constant in off diagonal terms of  $C$ .

# VB optimisation illustration for Gaussian processes





# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

Variational bayes

Expectation propagation

Comparisons



# Expectation propagation

$$p(\mathbf{f}|\mathbf{y}) \propto p(\mathbf{f}) \prod_{i=1}^n p(\mathbf{y}_i|\mathbf{f}_i)$$

$$q(\mathbf{f}|\mathbf{y}) \triangleq \frac{1}{Z_{ep}} p(\mathbf{f}) \prod_{i=1}^n t_i(\mathbf{f}_i | \tilde{Z}_i, \tilde{\mu}_i, \tilde{\sigma}_i^2) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$t_i \triangleq \tilde{Z}_i \mathcal{N}(\mathbf{f}_i | \tilde{\mu}_i, \tilde{\sigma}_i^2)$$

- ▶ Individual likelihood terms,  $p(\mathbf{y}_i|\mathbf{f}_i)$ , replaced by independent un-normalised 1D Gaussians,  $t_i$ .
- ▶ Uses an iterative algorithm to update  $t_i$ 's, such that the approximate marginal moments of  $q(\mathbf{f}_i)$ , match the marginal moments if  $t_i$  was replaced with the true likelihood  $p(\mathbf{y}_i|\mathbf{f}_i)$ .

# Expectation propagation



1. From the approximate current posterior,  $q(\mathbf{f}|\mathbf{y})$ , leave out one of the local likelihoods,  $t_i$ , and marginalise  $\mathbf{f}_j$  where  $j \neq i$ , giving rise to the approximate marginal with the contribution of one data removed; known as the *cavity distribution*,  $q_{-i}(\mathbf{f}_i)$ .



# Expectation propagation

1. From the approximate current posterior,  $q(\mathbf{f}|\mathbf{y})$ , leave out one of the local likelihoods,  $t_i$ , and marginalise  $\mathbf{f}_j$  where  $j \neq i$ , giving rise to the approximate marginal with the contribution of one data removed; known as the *cavity distribution*,  $q_{-i}(\mathbf{f}_i)$ .
2. Combine resulting cavity distribution,  $q_{-i}(\mathbf{f}_i)$ , with exact likelihood contribution,  $p(\mathbf{y}_i|\mathbf{f}_i)$ , giving a non-Gaussian un-normalized distribution,  $\hat{q}(\mathbf{f}_i) \triangleq p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i)$ .



# Expectation propagation

1. From the approximate current posterior,  $q(\mathbf{f}|\mathbf{y})$ , leave out one of the local likelihoods,  $t_i$ , and marginalise  $\mathbf{f}_j$  where  $j \neq i$ , giving rise to the approximate marginal with the contribution of one data removed; known as the *cavity distribution*,  $q_{-i}(\mathbf{f}_i)$ .
2. Combine resulting cavity distribution,  $q_{-i}(\mathbf{f}_i)$ , with exact likelihood contribution,  $p(\mathbf{y}_i|\mathbf{f}_i)$ , giving a non-Gaussian un-normalized distribution,  $\hat{q}(\mathbf{f}_i) \triangleq p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i)$ .
3. Choose a un-normalized Gaussian approximation to this distribution,  $\mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i$ , by finding moments of  $\hat{q}(\mathbf{f}_i)$ .



# Expectation propagation

1. From the approximate current posterior,  $q(\mathbf{f}|\mathbf{y})$ , leave out one of the local likelihoods,  $t_i$ , and marginalise  $\mathbf{f}_j$  where  $j \neq i$ , giving rise to the approximate marginal with the contribution of one data removed; known as the *cavity distribution*,  $q_{-i}(\mathbf{f}_i)$ .
2. Combine resulting cavity distribution,  $q_{-i}(\mathbf{f}_i)$ , with exact likelihood contribution,  $p(\mathbf{y}_i|\mathbf{f}_i)$ , giving a non-Gaussian un-normalized distribution,  $\hat{q}(\mathbf{f}_i) \triangleq p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i)$ .
3. Choose a un-normalized Gaussian approximation to this distribution,  $\mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i$ , by finding moments of  $\hat{q}(\mathbf{f}_i)$ .
4. Replace parameters of  $t_i$  with those that produce the same moments as this approximation.



# Expectation propagation

1. From the approximate current posterior,  $q(\mathbf{f}|\mathbf{y})$ , leave out one of the local likelihoods,  $t_i$ , and marginalise  $\mathbf{f}_j$  where  $j \neq i$ , giving rise to the approximate marginal with the contribution of one data removed; known as the *cavity distribution*,  $q_{-i}(\mathbf{f}_i)$ .
2. Combine resulting cavity distribution,  $q_{-i}(\mathbf{f}_i)$ , with exact likelihood contribution,  $p(\mathbf{y}_i|\mathbf{f}_i)$ , giving a non-Gaussian un-normalized distribution,  $\hat{q}(\mathbf{f}_i) \triangleq p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i)$ .
3. Choose a un-normalized Gaussian approximation to this distribution,  $\mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i$ , by finding moments of  $\hat{q}(\mathbf{f}_i)$ .
4. Replace parameters of  $t_i$  with those that produce the same moments as this approximation.
5. This minimizes  $\text{KL}\left(p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i) \parallel \mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i\right)$



# Expectation propagation

1. From the approximate current posterior,  $q(\mathbf{f}|\mathbf{y})$ , leave out one of the local likelihoods,  $t_i$ , and marginalise  $\mathbf{f}_j$  where  $j \neq i$ , giving rise to the approximate marginal with the contribution of one data removed; known as the *cavity distribution*,  $q_{-i}(\mathbf{f}_i)$ .
2. Combine resulting cavity distribution,  $q_{-i}(\mathbf{f}_i)$ , with exact likelihood contribution,  $p(\mathbf{y}_i|\mathbf{f}_i)$ , giving a non-Gaussian un-normalized distribution,  $\hat{q}(\mathbf{f}_i) \triangleq p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i)$ .
3. Choose a un-normalized Gaussian approximation to this distribution,  $\mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i$ , by finding moments of  $\hat{q}(\mathbf{f}_i)$ .
4. Replace parameters of  $t_i$  with those that produce the same moments as this approximation.
5. This minimizes  $\text{KL}\left(p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i) \parallel \mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i\right)$
6. Choose another  $i$  and start again. Repeat to convergence.



# Expectation propagation - in math

Step 1. First choose a local likelihood contribution,  $i$ , to leave out, and find the marginal cavity distribution,

$$\begin{aligned} q(\mathbf{f}|\mathbf{y}) &\propto p(\mathbf{f}) \prod_{j=1}^n t_j(\mathbf{f}_j) \rightarrow \frac{p(\mathbf{f}) \prod_{j=1}^n t_j(\mathbf{f}_j)}{t_i(\mathbf{f}_i)} \rightarrow p(\mathbf{f}) \prod_{j \neq i}^n t_j(\mathbf{f}_j) \\ &\rightarrow \int p(\mathbf{f}) \prod_{j \neq i}^n t_j(\mathbf{f}_j) d\mathbf{f}_{j \neq i} \triangleq q_{-i}(\mathbf{f}_i) \end{aligned}$$

Step 2.  $p(\mathbf{y}_i|\mathbf{f}_i)q_{-i}(\mathbf{f}_i) \triangleq \hat{q}(\mathbf{f}_i)$

Step 3.  $\hat{q}(\mathbf{f}_i) \approx \mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)\hat{Z}_i$

Step 4: Compute parameters of  $t_i(\mathbf{f}_i|\tilde{Z}_i, \tilde{\mu}_i, \tilde{\sigma}_i^2)$  making moments of  $q(\mathbf{f}_i)$  match those of  $\hat{Z}_i \mathcal{N}\left(\mathbf{f}_i|\hat{\mu}_i, \hat{\sigma}_i^2\right)$ .



# Outline

Motivation

Non-Gaussian posteriors

Approximate methods

Laplace approximation

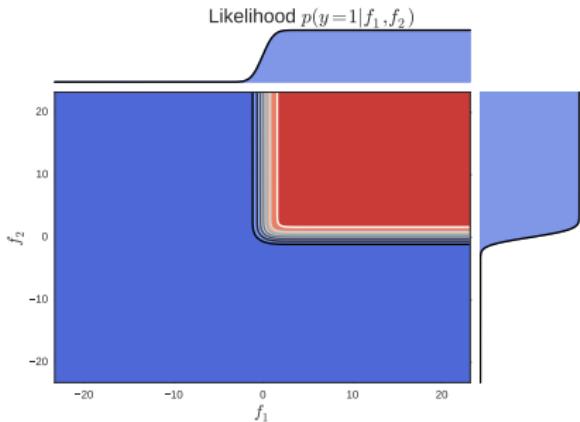
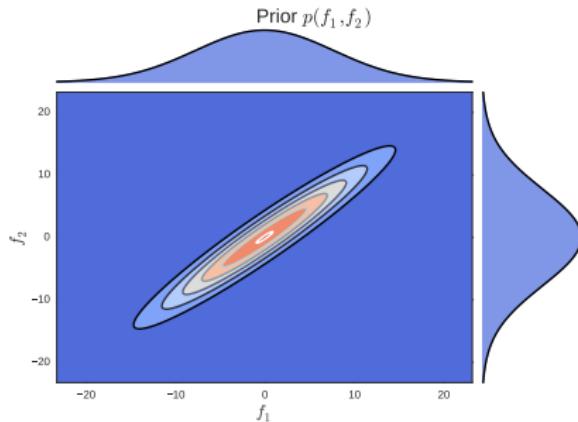
Variational bayes

Expectation propagation

Comparisons



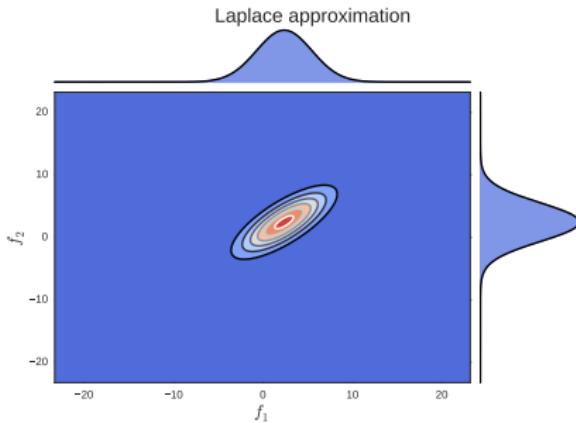
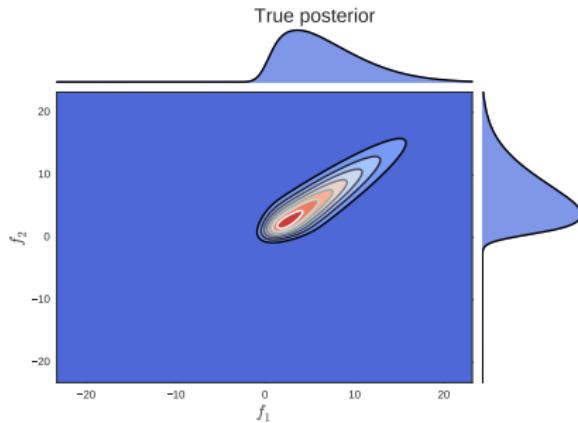
# Comparing posterior approximations



- ▶ Gaussian prior between two function values  $\{f_1, f_2\}$ , at  $\{x_1, x_2\}$  respectively.
- ▶ Bernoulli likelihood,  $y_1 = 1$  and  $y_2 = 1$ .



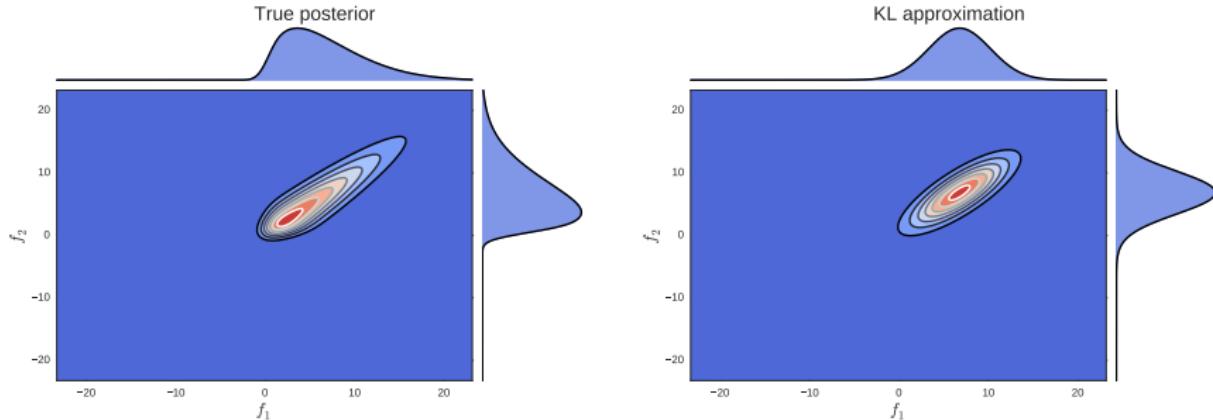
# Comparing posterior approximations



- ▶  $p(\mathbf{f}|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f})}{p(\mathbf{y})}$
- ▶ True posterior is non-Gaussian.
- ▶ Laplace approximates with a Gaussian at the mode of the posterior.



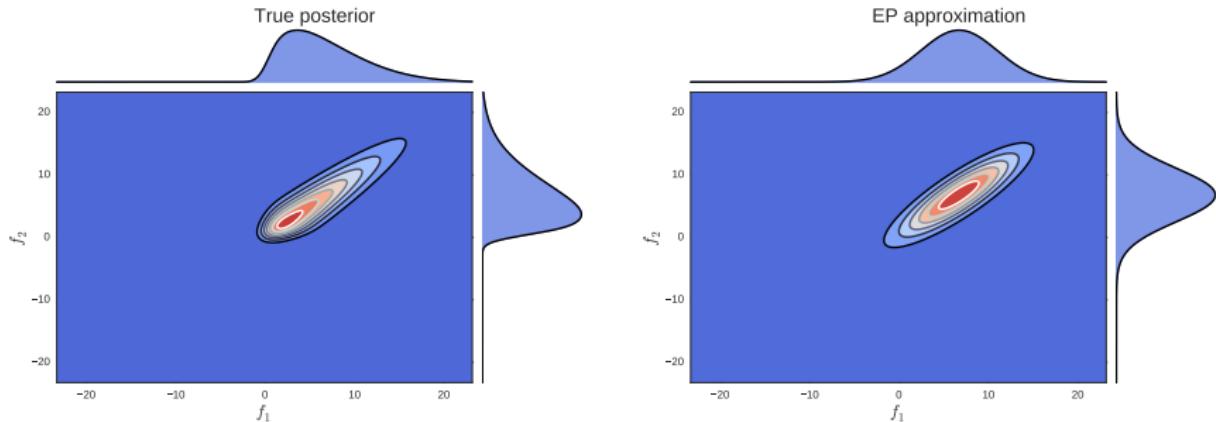
# Comparing posterior approximations



- ▶ True posterior is non-Gaussian.
- ▶ VB approximate with a Gaussian that has minimal KL divergence,  $\text{KL}(q(\mathbf{f}) \parallel p(\mathbf{f}|\mathbf{y}))$ .
- ▶ This leads to distributions that avoid regions in which  $p(\mathbf{f}|\mathbf{y})$  is small.
- ▶ It has a large penalty for assigning density where there is none.

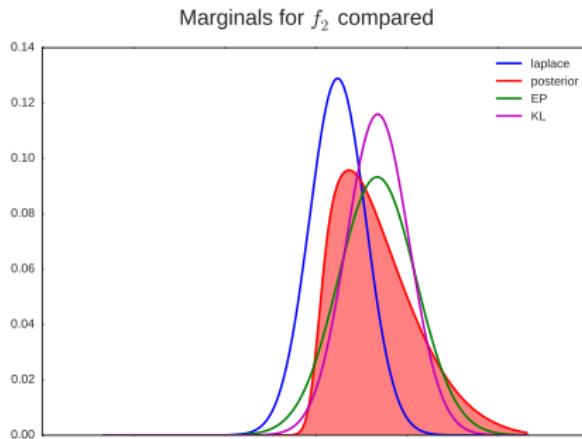


# Comparing posterior approximations



- ▶ True posterior is non-Gaussian.
- ▶ EP tends to try and put density where  $p(\mathbf{f}|\mathbf{y})$  is large
- ▶ Cares less about assigning density density where there is none. Contrasts to VB method.

# Comparing posterior marginal approximations



- ▶ Laplace: Poor approximation.
- ▶ VB: Avoids assigning density to areas where there is none, at the expense of areas where there is some (right tail).
- ▶ EP: Assigns density to areas with density, at the expense of areas where there is none (left tail).



## Laplace approximation

- ▶ Pros
  - ▶ Very fast.
- ▶ Cons
  - ▶ Poor approximation if the mode does not well describe the posterior, for example Bernoulli likelihood.
- ▶ When
  - ▶ When the posterior *is* well characterized by its mode, for example Poisson.



## Variational Bayes

- ▶ Pros
  - ▶ Principled in that we are directly optimizing a measure of divergence between an approximation and true distribution.
  - ▶ Lends itself to sparse extensions.
- ▶ Cons
  - ▶ Requires factorizing likelihoods to avoid  $n$  dimensional integral.
  - ▶ As seen, can result in underestimating the variance, i.e. becomes overconfident.
- ▶ When
  - ▶ Applicable to a range of likelihoods, but is known in some cases to underestimate variance, might need to be careful if you wish to be conservative with predictive uncertainty.
  - ▶ In conjunction with sparse methods.



## EP method

- ▶ Pros
  - ▶ Very effective for certain likelihoods (classification).
  - ▶ Also lends itself to sparse approximations.
- ▶ Cons
  - ▶ Standard algorithm is slow though possible to extend to sparse case.
  - ▶ Convergence issues for some likelihoods.
  - ▶ Must be able to match moments.
- ▶ When
  - ▶ Binary data (Nickisch and Rasmussen, 2008; Kuß, 2006), perhaps with truncated likelihood (censored data) (Vanhatalo et al., 2015).
  - ▶ In conjunction with sparse methods.



## MCMC methods

- ▶ Pros
  - ▶ Theoretical limit gives true distribution
- ▶ Cons
  - ▶ Can be very slow
- ▶ When
  - ▶ If time is not an issue, but exact accuracy is.
  - ▶ If you are unsure whether a different approximation is appropriate, can be used as a “ground truth”

# Conclusion



- ▶ Many real world tasks require non-Gaussian observation models.
- ▶ Non-Gaussian likelihoods cause complications in applying our framework.
- ▶ Several different ways to deal with the problem. Many are based on Gaussian approximations.
- ▶ Different methods have their own advantages and disadvantages.

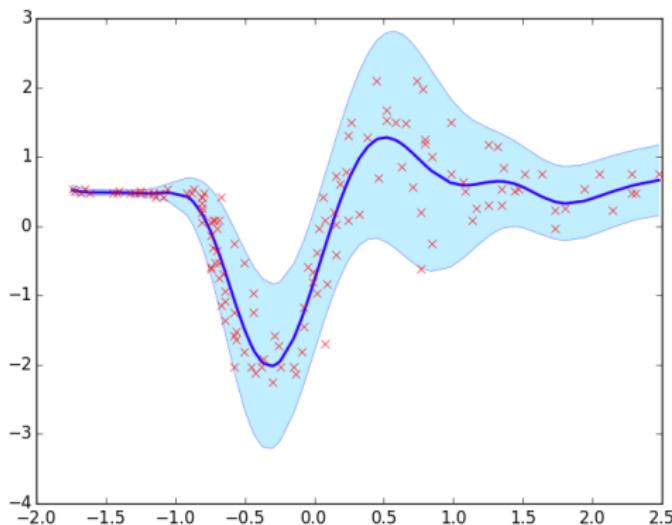
# Questions



Thanks for listening.

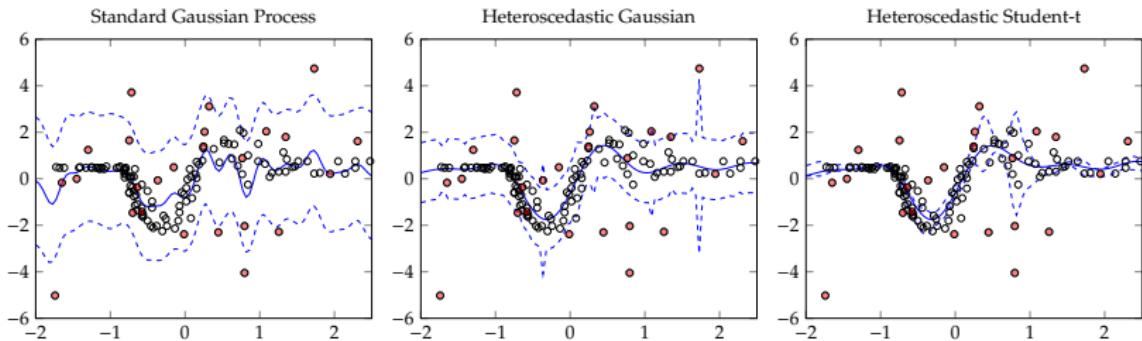
Any questions?

## Bonus - Hetroscedastic likelihoods



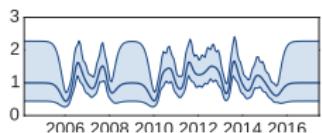
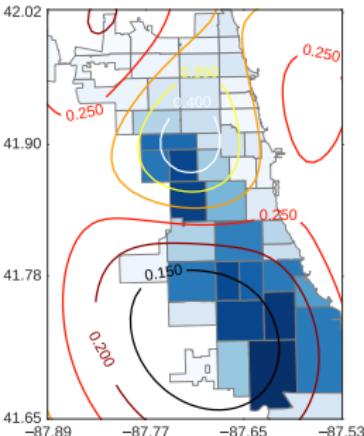
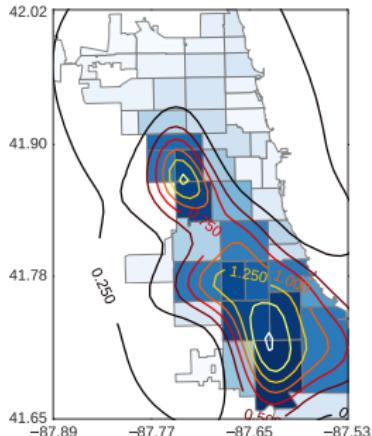
- ▶ Likelihood whose parameters are governed by two known functions,  $\mathbf{f}$  and  $\mathbf{g}$ .
- ▶  $p(\mathbf{y}|\mathbf{f}, \mathbf{g}) = \mathcal{N}(\mathbf{y}|\mu = \mathbf{f}, \sigma^2 = \exp(\mathbf{g}))$

# Bonus - non-Gaussian hetroscedastic likelihoods



- ▶ Likelihood whose parameters are governed by two known functions,  $\mathbf{f}$  and  $\mathbf{g}$ .
- ▶  $p(\mathbf{y}|\mathbf{f}, \mathbf{g}) = t(\mathbf{y}|\mu = \mathbf{f}, \sigma^2 = \exp(\mathbf{g}), \nu = 3.0)$

# Bonus - non-Gaussian heteroscedastic likelihoods



- ▶  $\Lambda(\mathbf{x}, \mathbf{t}) = \lambda_1(\mathbf{x})\mu_1(\mathbf{t}) + \lambda_2(\mathbf{x})\mu_2(\mathbf{t})$

# References I



- Hensman, J., Matthews, A. G. D. G., and Ghahramani, Z. (2015). Scalable variational gaussian process classification. In *In 18th International Conference on Artificial Intelligence and Statistics*, pages 1–9, San Diego, California, USA.
- Kuß, M. (2006). *Gaussian Process Models for Robust Regression, Classification, and Reinforcement Learning*. PhD thesis, TU Darmstadt.
- Nickisch, H. and Rasmussen, C. E. (2008). Approximations for Binary Gaussian Process Classification. *Journal of Machine Learning Research*, 9:2035–2078.
- Vanhatalo, J., Riihimaki, J., Hartikainen, J., Jylankki, P., Tolvanen, V., and Vehtari, A. (2015). Gpstuff.  
<http://mloss.org/software/view/451/>.