

CS-E4895 Gaussian Processes

Lecture 6: Classification

Arno Solin

Aalto University

Tuesday 14.3.2023

based on slides by ST John and Michael Riis Andersen

Roadmap for today

- 1 Beyond Gaussian noise
 - Classification vs. regression
 - Other likelihoods
 - Looking at the likelihood in more detail
- 2 Inference for arbitrary likelihoods
 - Posterior predictive distribution
 - Why is the posterior intractable?
- 3 Approximating the intractable
 - Gaussian approximations
 - Laplace approximation
 - Minimising divergences
 - Variational inference
- 4 Conclusion

Section 1

Beyond Gaussian noise

Regression vs. classification

- Response variable y is continuous in regression problems

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

- Response variable y is discrete in classification problems

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

- Classification problems

X = images,

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

X = X-ray scan,

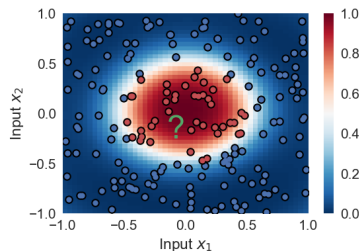
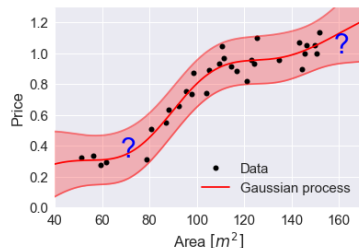
$y_n \in \{\text{tumor}, \text{no tumor}\}$

X = images of digits,

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

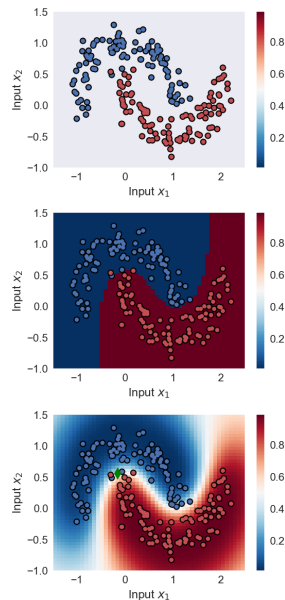
X = emails,

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



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



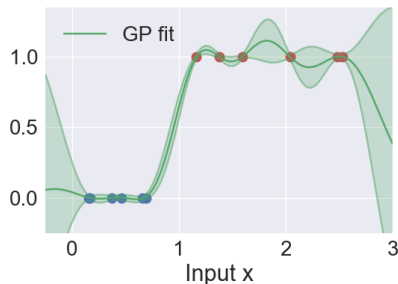
Why don't we use regression models for classification?

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

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

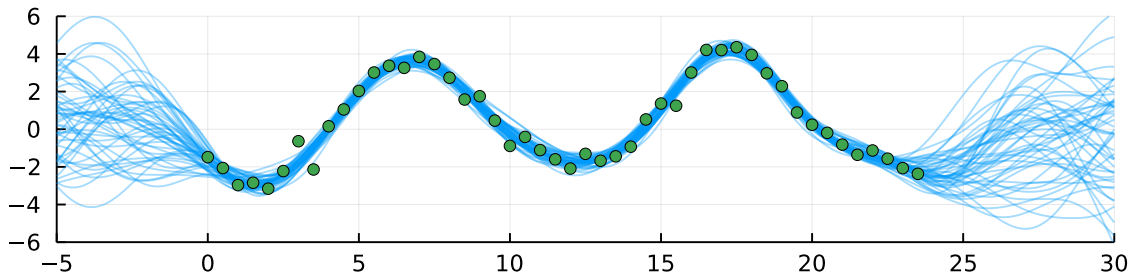
- What's wrong with simply using the GP regression model with labels: $y_n \in \{0, 1\}$:

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



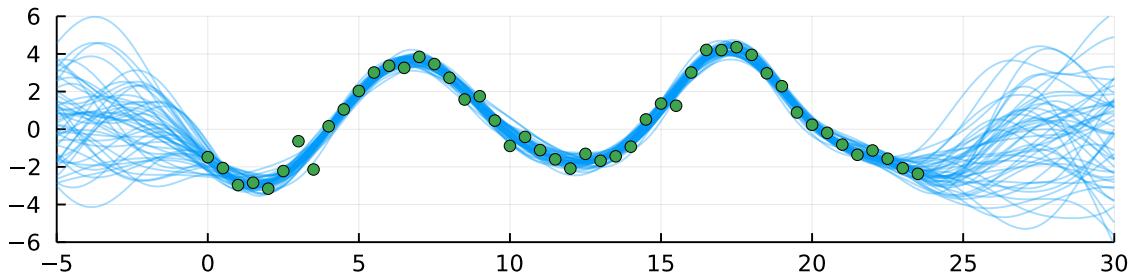
Recap: Gaussian noise model

$$y(x) = f(x) + \epsilon, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\text{noise}}^2)$$



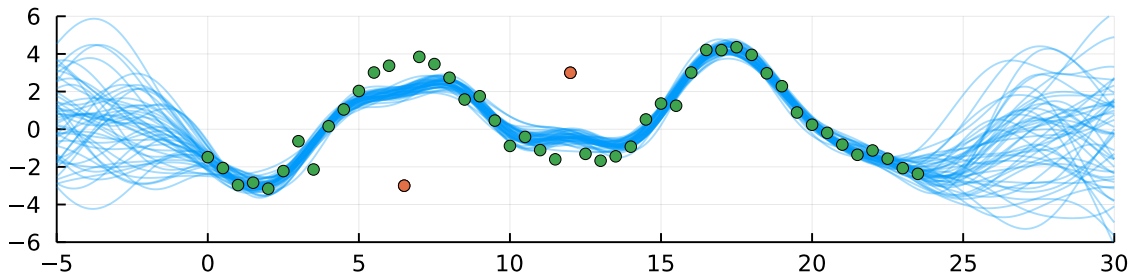
Recap: Gaussian noise model

$$y(x) = f(x) + \epsilon, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\text{noise}}^2)$$
$$p(y | f) = \mathcal{N}(y | f, \sigma_{\text{noise}}^2)$$

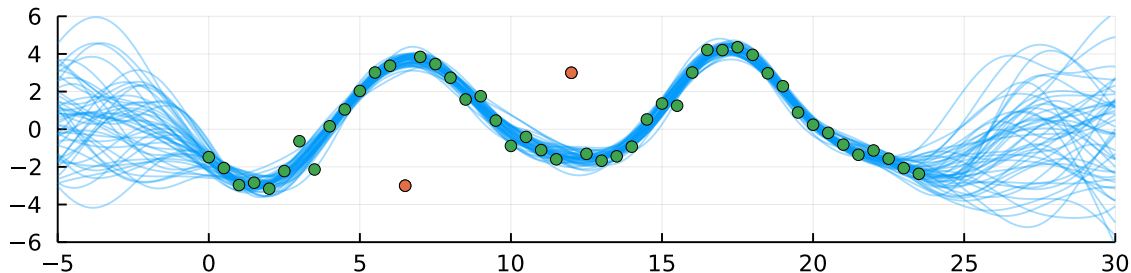


Misspecified Gaussian noise model

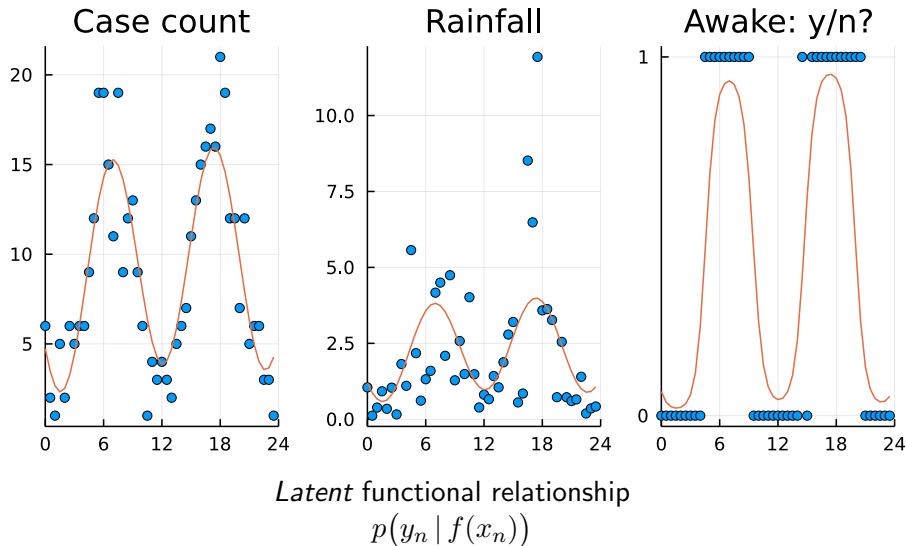
$$y(x) = f(x) + \epsilon, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_{\text{noise}}^2)$$
$$p(y | f) = \mathcal{N}(y | f, \sigma_{\text{noise}}^2)$$



Heavy-tailed noise model



Non-Gaussian observations



Likelihood

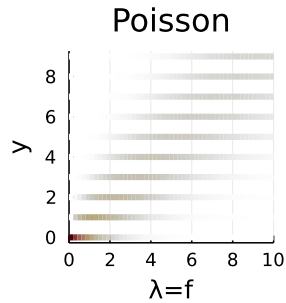
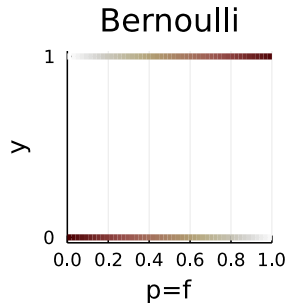
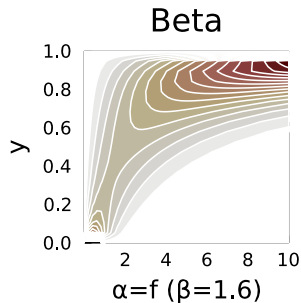
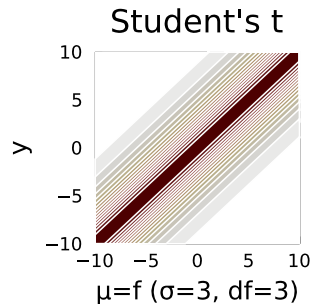
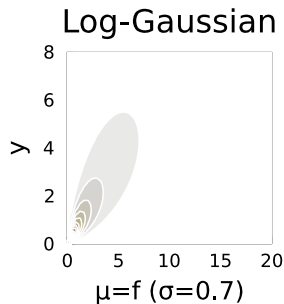
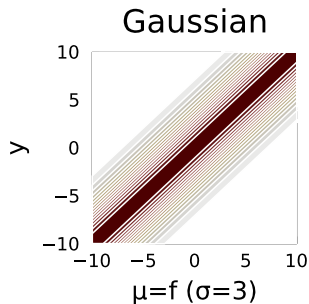
$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N p(y_n | f_n); \quad f_n = f(x_n)$$

factorizing

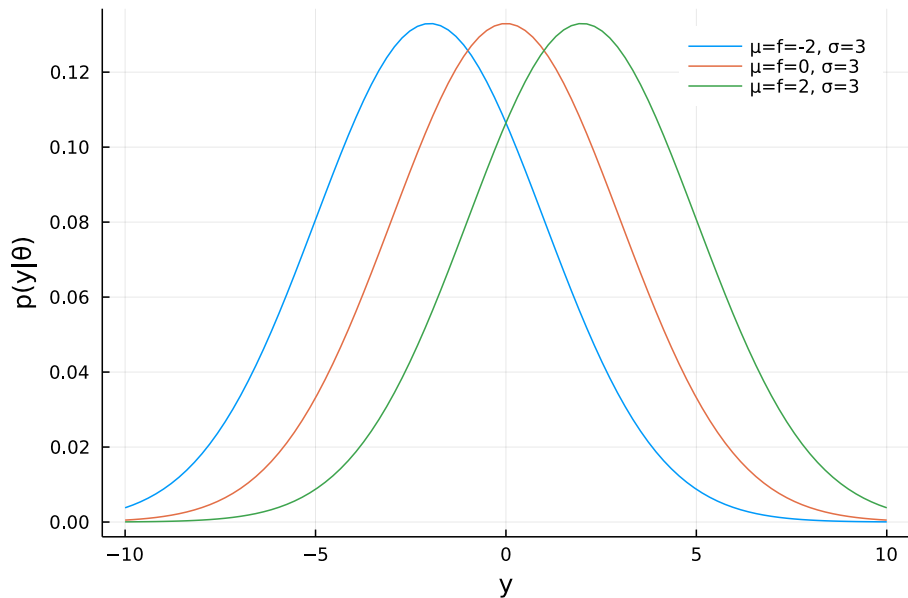
$$p(y | f)$$

Function of two arguments:

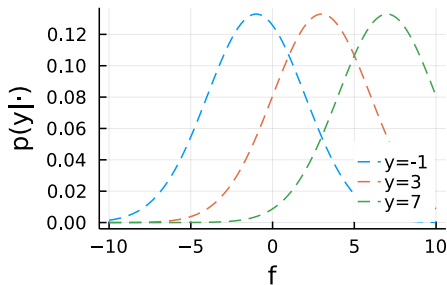
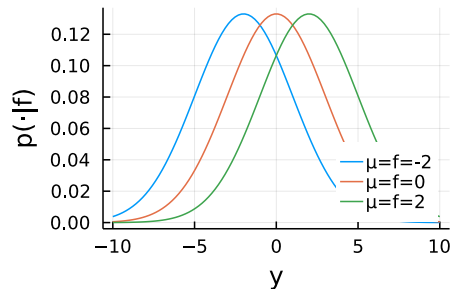
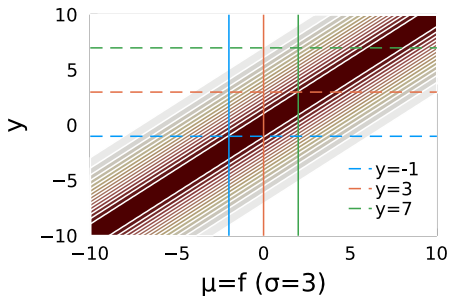
$$y \mapsto p(y | f), \quad f \mapsto p(y | f)$$



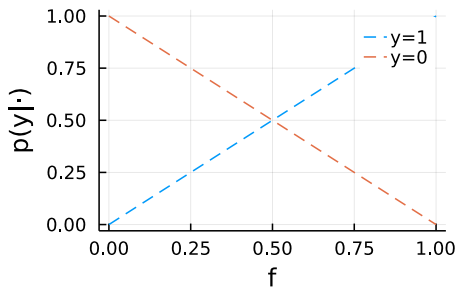
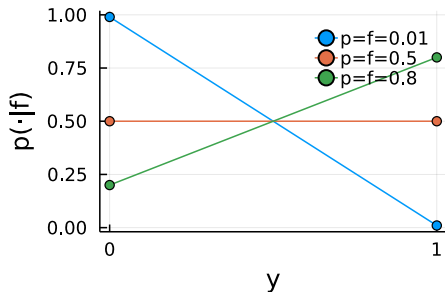
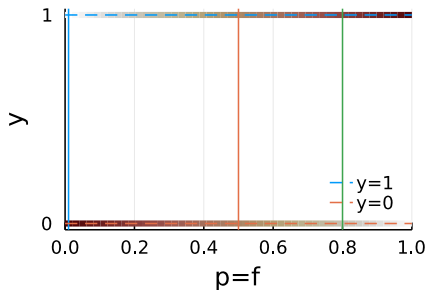
$p(y | f)$: Gaussian



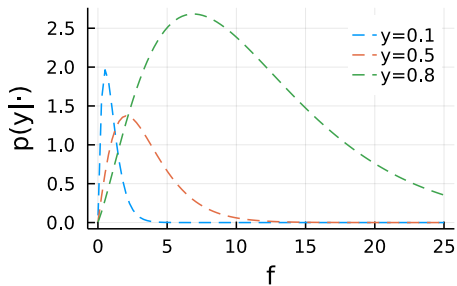
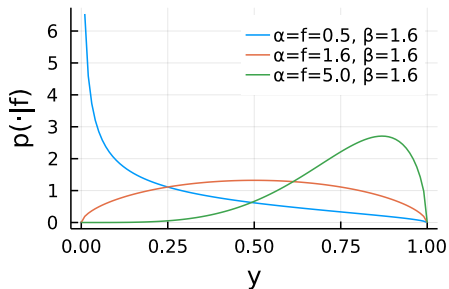
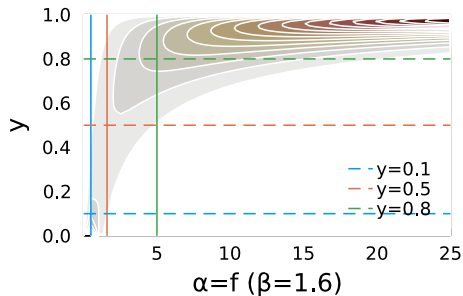
$p(y | f)$: Gaussian

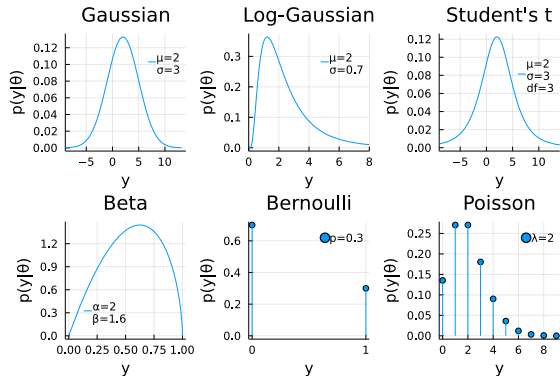
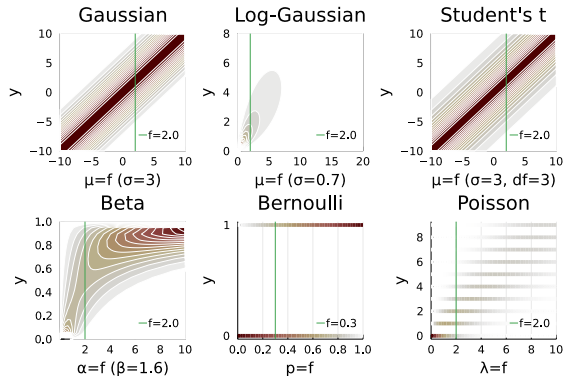


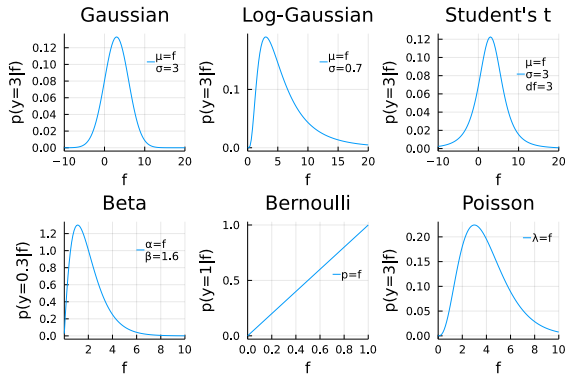
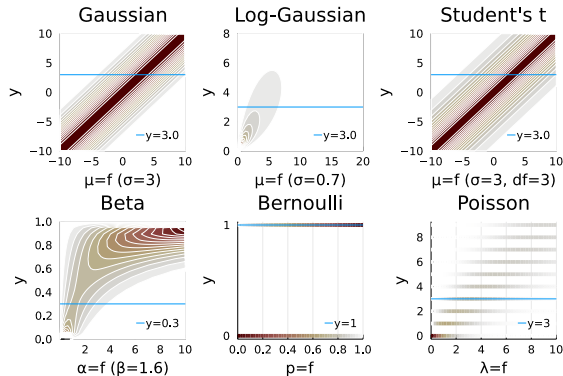
$p(y | f)$: Bernoulli



$p(y | f)$: Beta







Two important aspects of likelihoods:

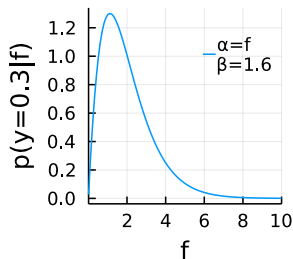
- 1 link functions
- 2 log-concavity

Link functions

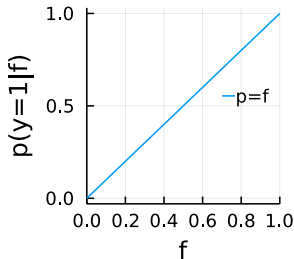
$$\mathbb{E}[y] = \theta \in (0 \dots \infty)$$

$$f \sim \mathcal{N} \quad \in (-\infty \dots \infty)$$

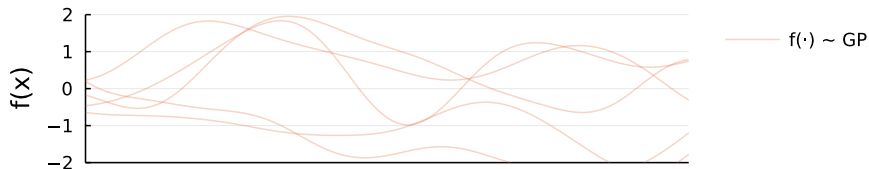
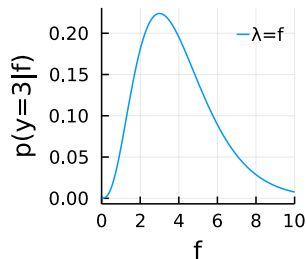
Beta



Bernoulli



Poisson



Link functions

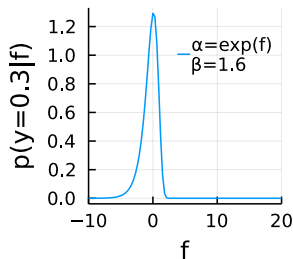
$$\mathbb{E}[y] = \theta \in (0 \dots \infty)$$

$$f \sim \mathcal{N} \quad \in (-\infty \dots \infty)$$

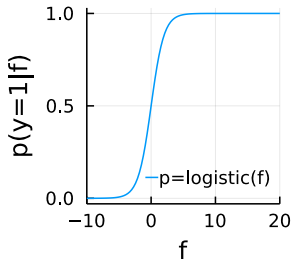
$$\text{link}(\theta) = f$$

$$\theta = \text{invlink}(f)$$

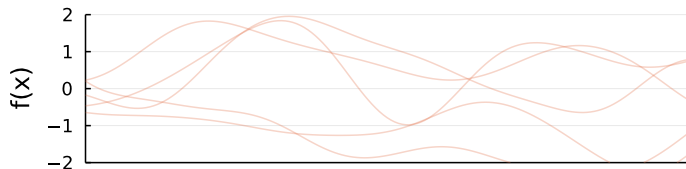
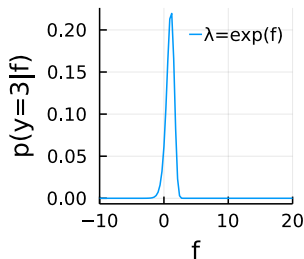
Beta



Bernoulli



Poisson



Link functions

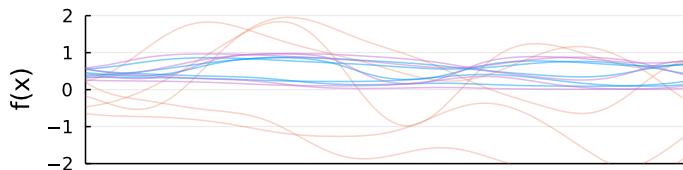
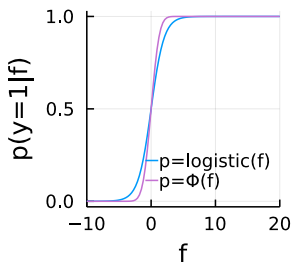
$$\mathbb{E}[y] = \theta \in (0 \dots \infty)$$

$$f \sim \mathcal{N} \quad \in (-\infty \dots \infty)$$

$$\text{link}(\theta) = f$$

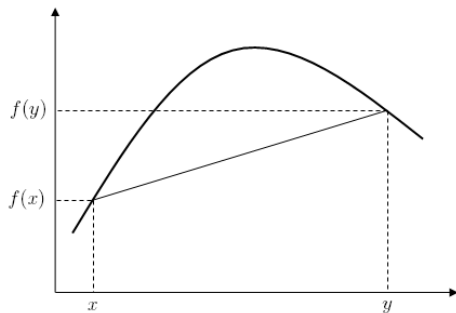
$$\theta = \text{invlink}(f)$$

Bernoulli



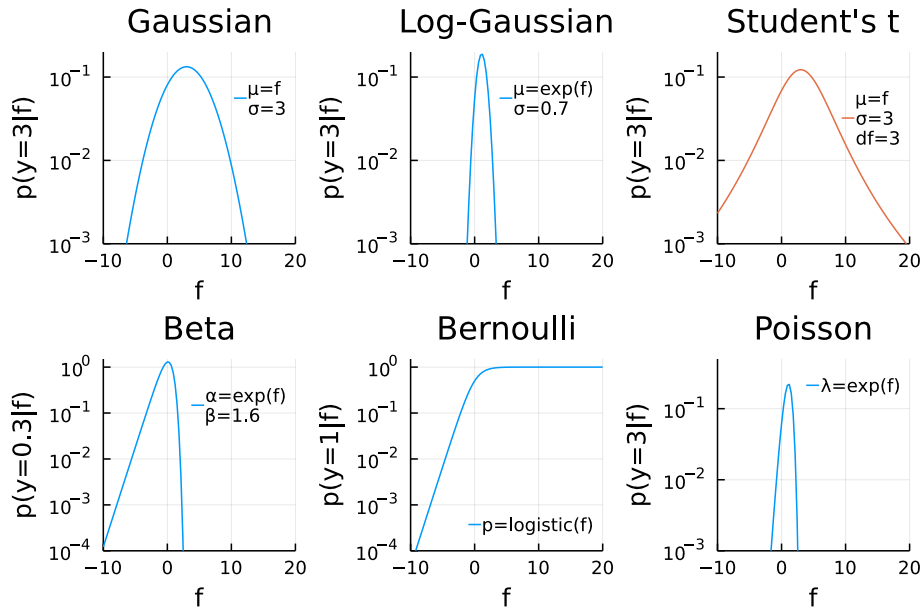
$f(\cdot) \sim \text{GP}$
 $\text{logistic}(f(\cdot))$
 $\Phi(f(\cdot))$

(Log-)concavity



$$f(\alpha x + (1 - \alpha)y) \geq \alpha f(x) + (1 - \alpha)f(y)$$

Log-concavity of likelihoods



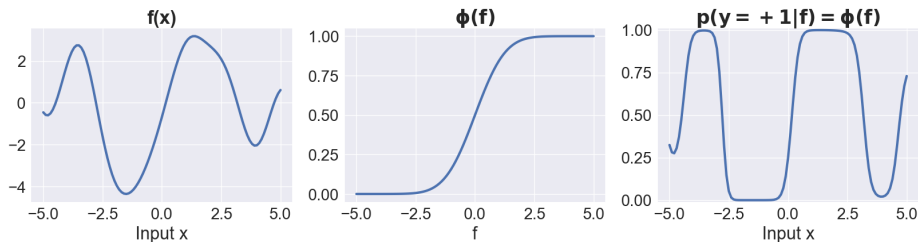
Section 2

Inference for arbitrary likelihoods

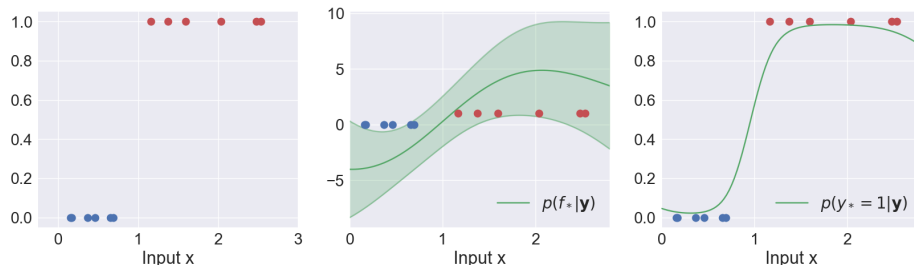
- 1 Beyond Gaussian noise
- 2 Inference for arbitrary likelihoods
 - Posterior predictive distribution
 - Why is the posterior intractable?
- 3 Approximating the intractable
- 4 Conclusion

GP classification: Connecting the dots

- We map the unknown function $f(x)$ through the squashing function



- Example re-visited



Predictive distribution at new test point \mathbf{x}_*

- ① Joint model:

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

- ② Posterior distribution at training points:

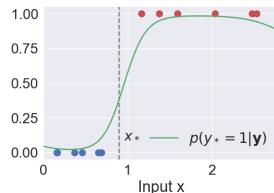
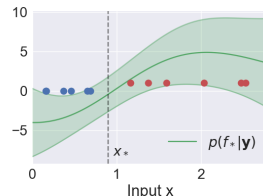
$$p(\mathbf{f} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f})}{p(\mathbf{y})} \approx q(\mathbf{f})$$

- ③ Posterior of f_* for new test point \mathbf{x}_* :

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

- ④ Predictive distribution

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



Analytically intractable distributions!

Posterior predictions

At new point x^* :

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

At training data:

$$p(\mathbf{f} | \mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{f} | \mathbf{x}) \prod_{n=1}^N p(y_n | f(x_n))}{\int p(\mathbf{f}' | \mathbf{x}) \prod_{n=1}^N p(y_n | f'(x_n)) d\mathbf{f}'}$$

$$p(\mathbf{f} | \mathbf{y}) = \frac{1}{Z} p(\mathbf{f}) \prod_{n=1}^N p(y_n | f_n)$$

$$Z = p(\mathbf{y} | \mathcal{M}) = \int p(\mathbf{f} | \mathcal{M}) \prod_{n=1}^N p(y_n | f_n, \mathcal{M}) d\mathbf{f}$$

“marginal likelihood” or “evidence” given model \mathcal{M}

Posterior at training points

$$p(\mathbf{f} | \mathbf{y}) = \frac{1}{Z} p(\mathbf{f}) \prod_{n=1}^N p(y_n | f_n)$$

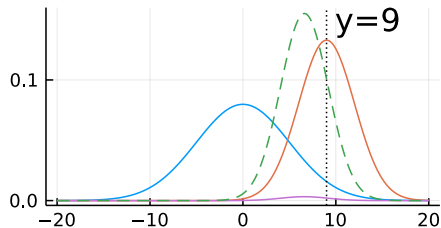
Gaussian (process) prior $p(f(\cdot)) \dots$

& Gaussian likelihood: conjugate case \rightarrow posterior Gaussian

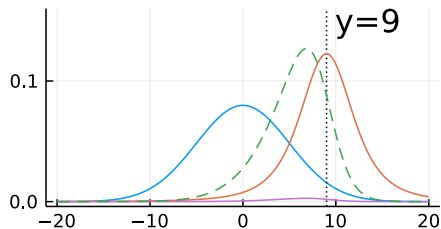
& non-Gaussian $p(y | f) \rightarrow p(\mathbf{f} | \mathbf{y})$ also non-Gaussian, intractable

1D examples

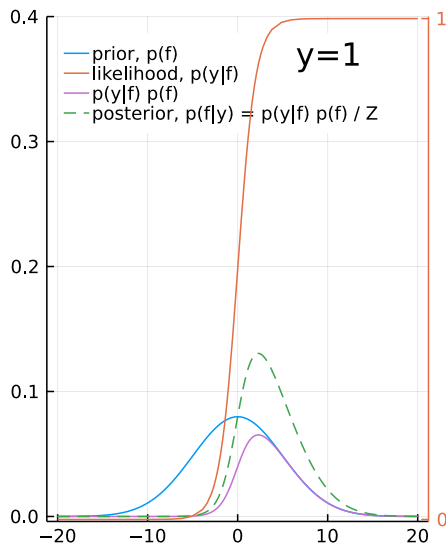
Gaussian



Student's t

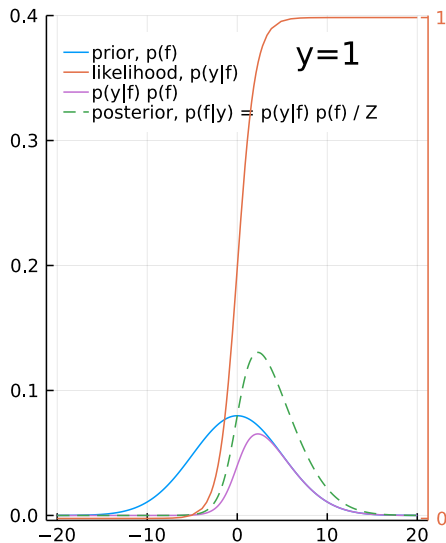


Bernoulli

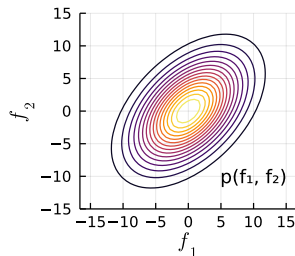


Bernoulli example in 2D

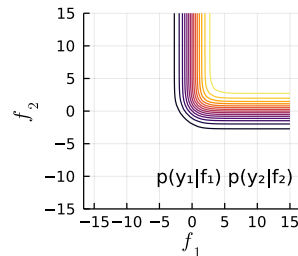
Bernoulli



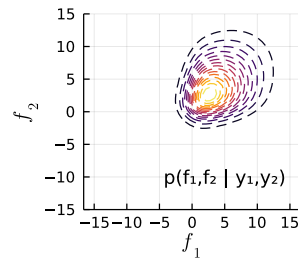
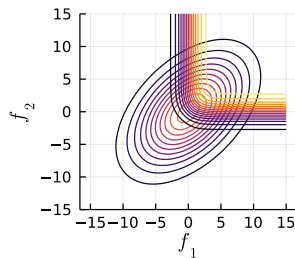
Prior



Likelihood



Posterior



Posterior for N observations

$$p(\mathbf{f} | \mathbf{y}) = \frac{p(\mathbf{f}) \prod_{n=1}^N p(y_n | f_n)}{\int p(\mathbf{f}') \prod_{n=1}^N p(y_n | f'_n) d\mathbf{f}'}$$

$$f_1 = f(x_1)$$

$$f_2 = f(x_2)$$

$$\vdots$$

$$f_N = f(x_N)$$

Summary so far

- What is the likelihood $p(y | f)$?
- When is it non-Gaussian?
- Why does the posterior $p(f | y)$ become intractable?

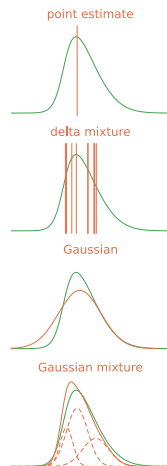
Section 3

Approximating the intractable

- 1 Beyond Gaussian noise
- 2 Inference for arbitrary likelihoods
- 3 Approximating the intractable
 - Gaussian approximations
 - Laplace approximation
 - Minimising divergences
 - Variational inference
- 4 Conclusion

Approximating distributions

- Delta distribution
 - Point estimate
- Mixture of delta distributions
 - Markov Chain Monte Carlo (MCMC)
 - Neural network ensembles. . .
- **Gaussian distribution**
 - Laplace
 - Variational Bayes/Variational Inference (VB / VI)
 - Expectation Propagation (EP), PowerEP, . . .
- Mixture of Gaussians
- . . .



Approximating the exact posterior with Gaussian

Approximating the posterior at observations:

$$p(\mathbf{f} | \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mu = ?, \Sigma = ?)$$

Predictions at new points:

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

What does this mean for Gaussian processes?

Choosing μ and Σ for $q(\mathbf{f})$

$$p(\mathbf{f} | \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mu = ?, \Sigma = ?)$$

locally: match mean &
variance at point

globally: minimise divergence

**Laplace
approximation**

Variational
inference (VI)

Expectation
Propagation (EP)

Laplace approximation

Idea: log of Gaussian pdf = quadratic polynomial

$$p_{\mathcal{N}}(\mathbf{f}) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{f} - \mu)^\top \Sigma^{-1}(\mathbf{f} - \mu)\right)$$

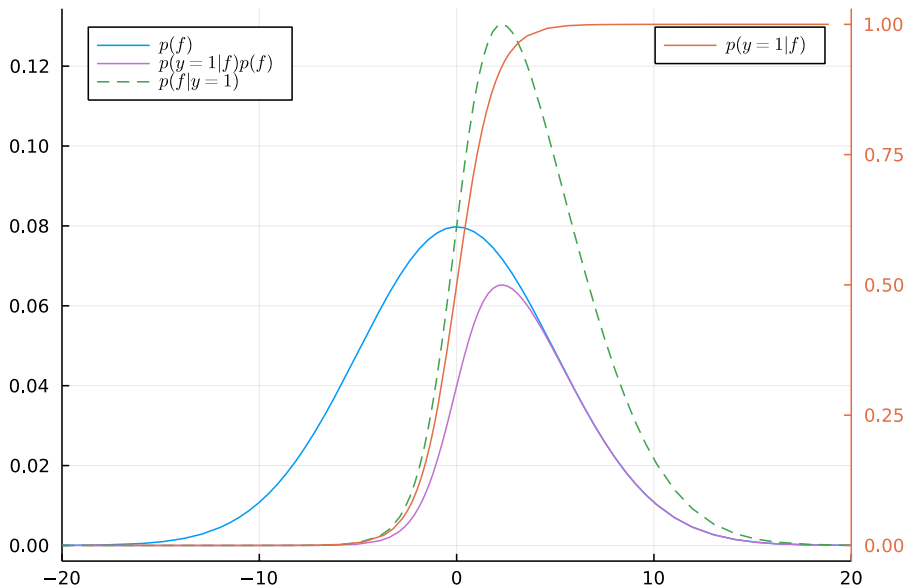
Quadratic polynomial through approximation:

2nd-order Taylor expansion of log of $h(f) = p(y|f)p(f)$ at \hat{f}

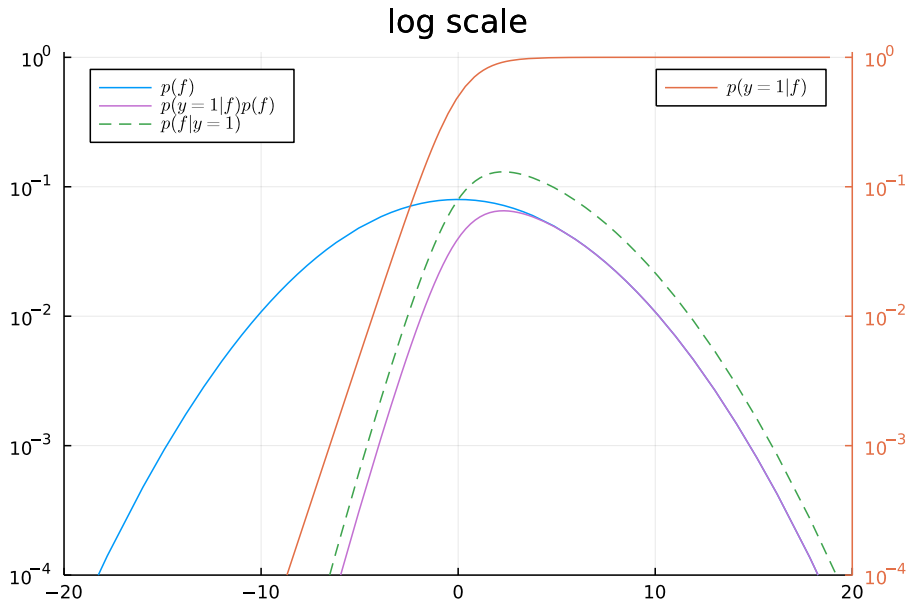
$$g(x + \delta) \approx g(x) + \left(\frac{dg}{dx}(x)\right)\delta + \frac{1}{2!}\left(\frac{d^2g}{dx^2}(x)\right)\delta^2$$

- 1 Find **mode** of posterior
2nd-order gradient optimisation (e.g. Newton's method)
- 2 Match **curvature** (Hessian) at mode

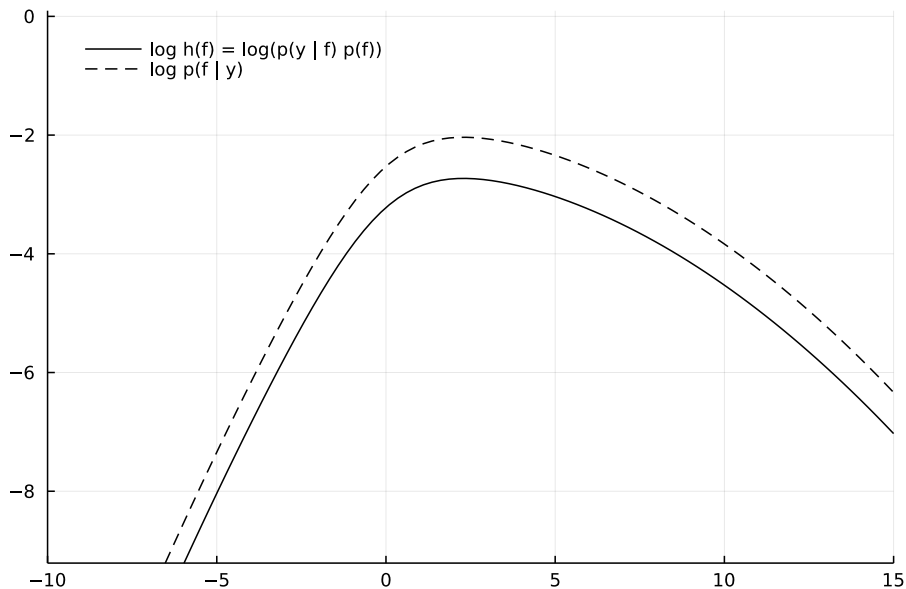
$$p(f | y) = \frac{1}{Z} p(y | f) p(f)$$



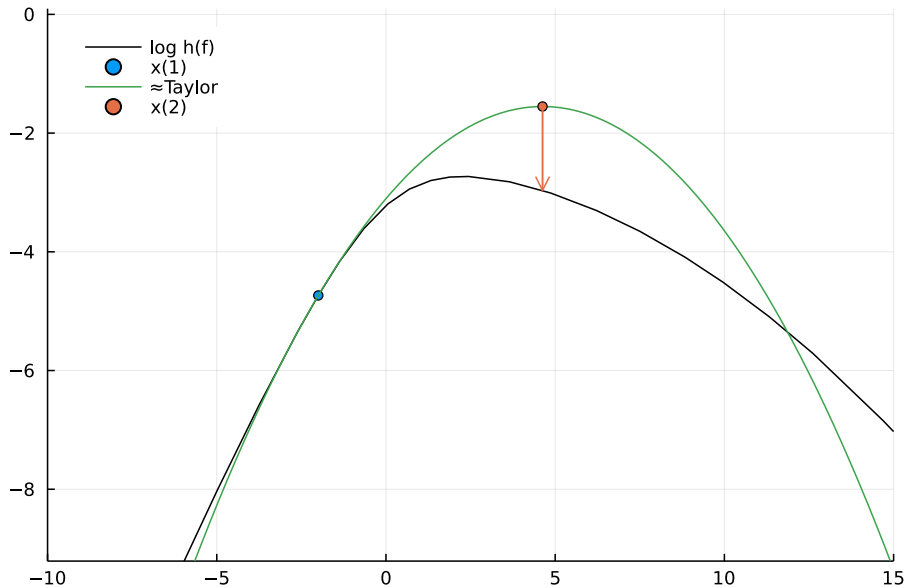
$$\log p(f | y) = -\log Z + \log p(y | f) + \log p(f)$$



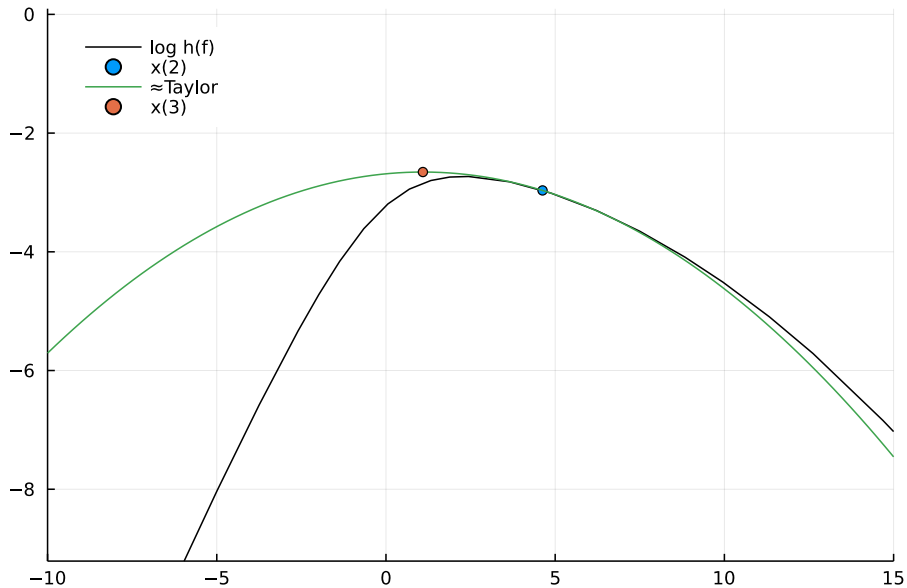
$$\log p(f | y) = -\log Z + \log h(f)$$



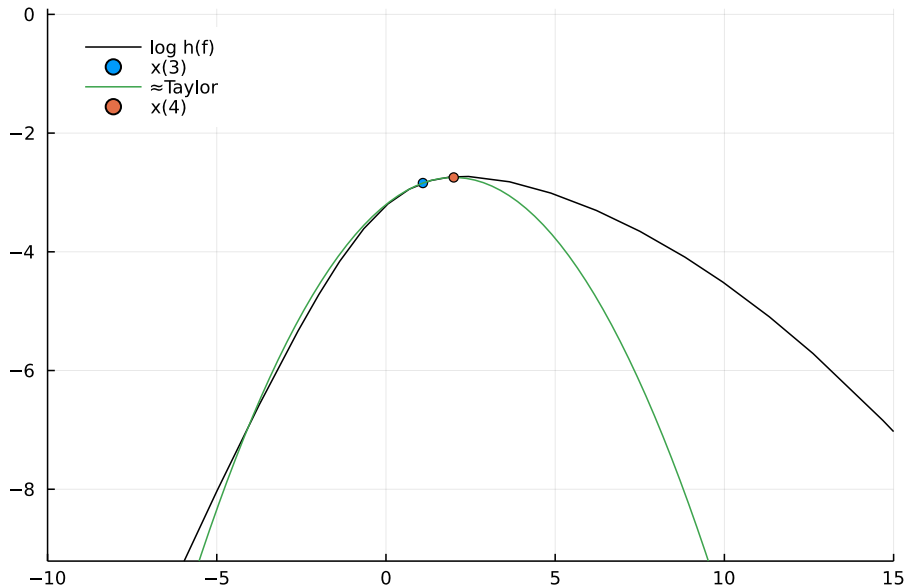
Newton's method



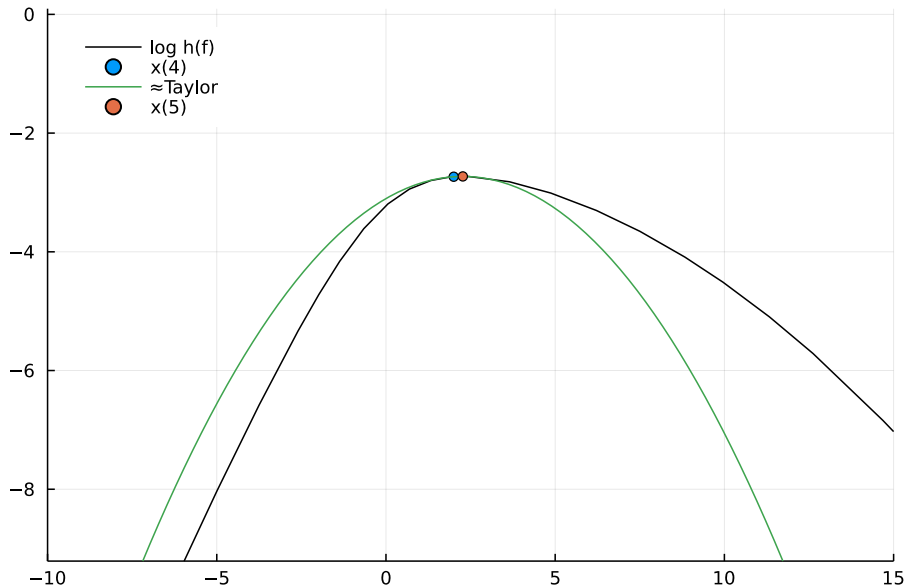
Newton's method



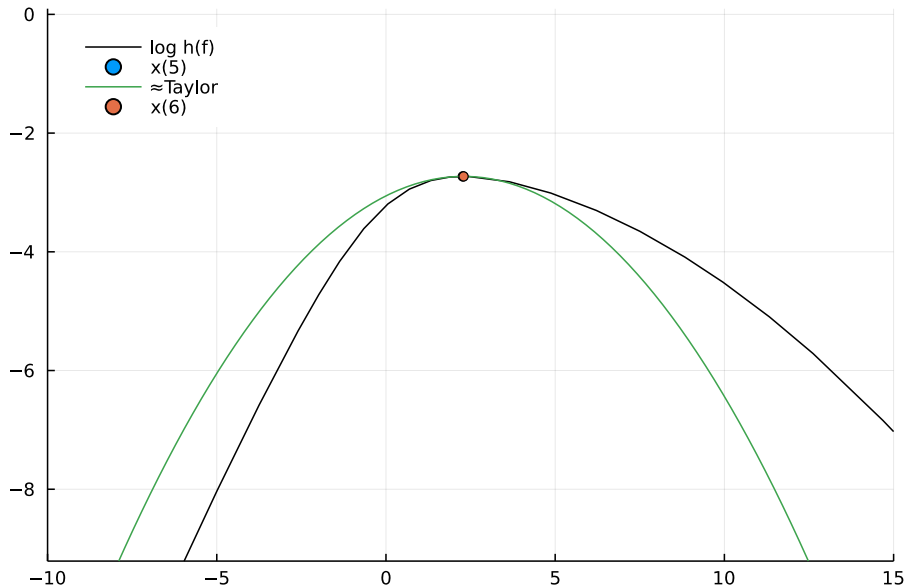
Newton's method



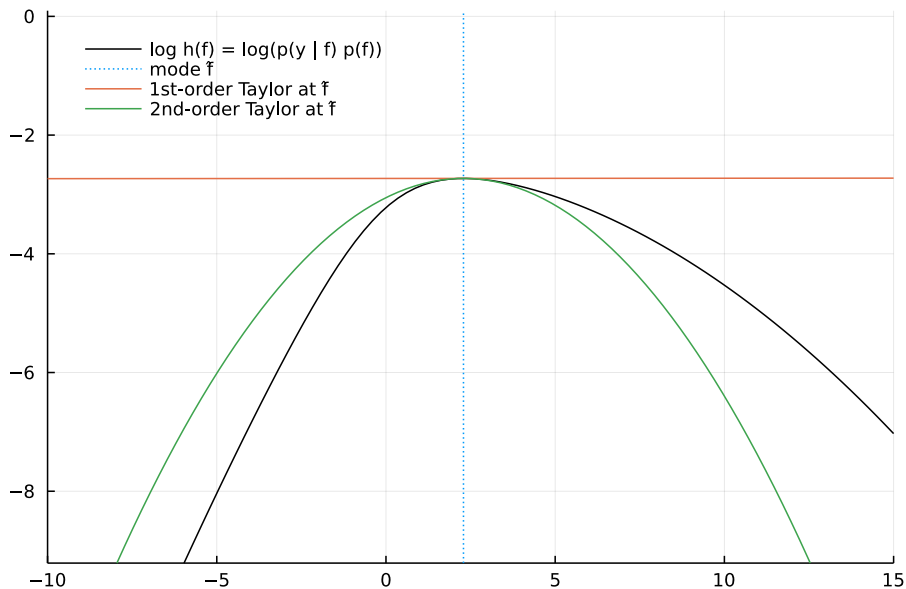
Newton's method



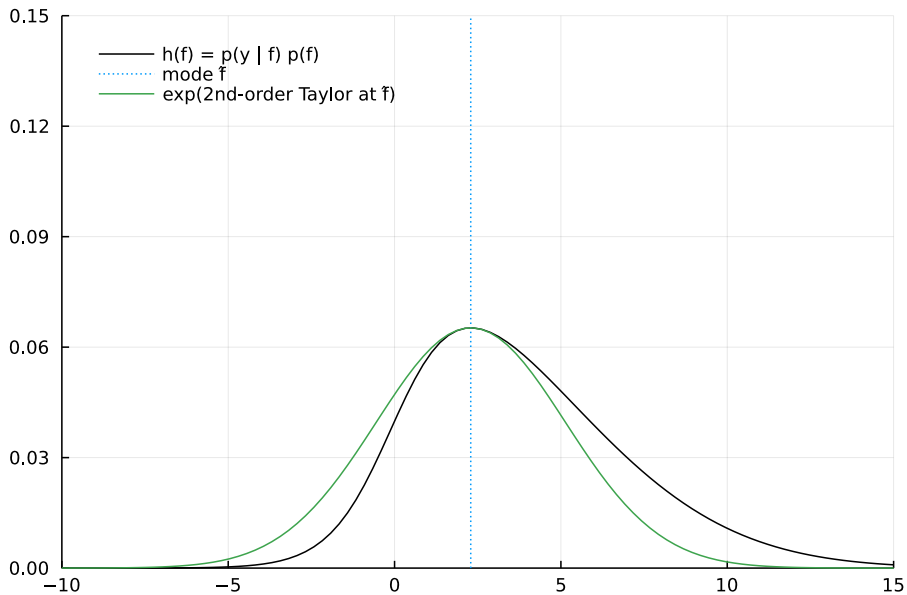
Newton's method



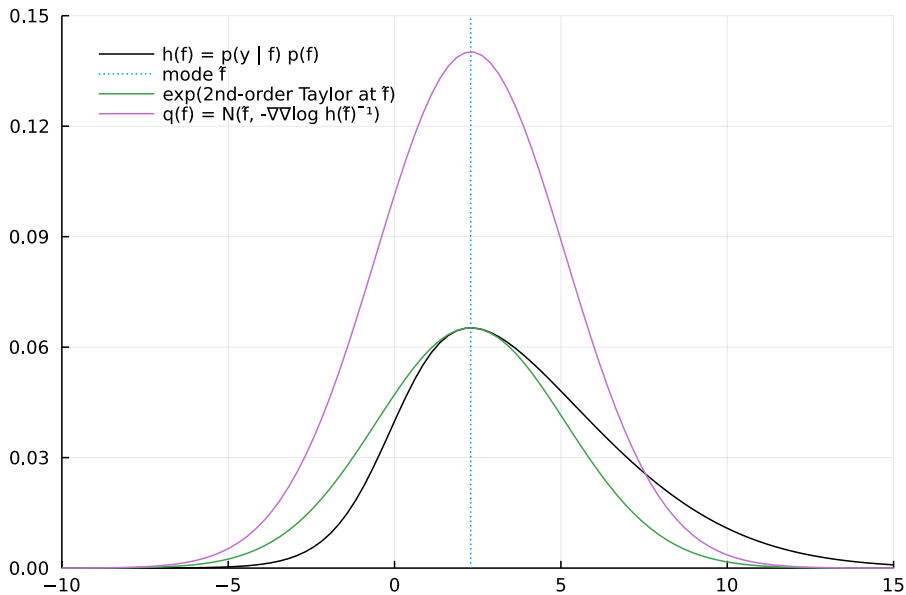
$$\log p(f | y) + \log Z = \log h(f) \approx \mathcal{O}(f^2)$$



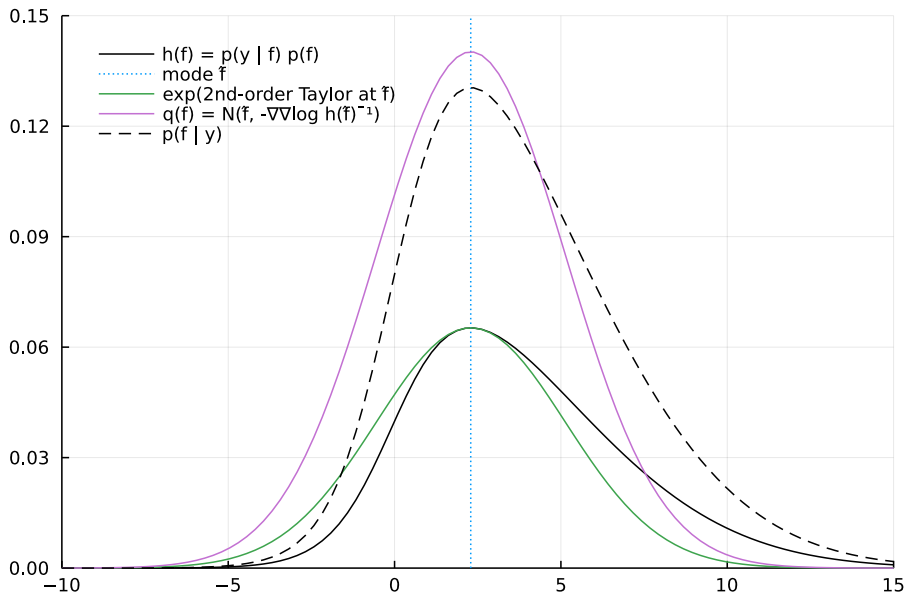
$$p(f | y) Z \approx \exp(\mathcal{O}(f^2))$$



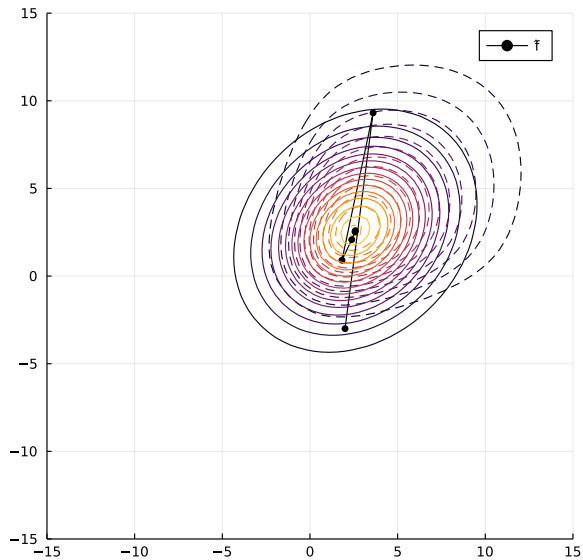
$$p(f | y) \approx \mathcal{N}(f | \hat{f}, -(\mathrm{d}^2 \log h / \mathrm{d} f^2)^{-1})$$



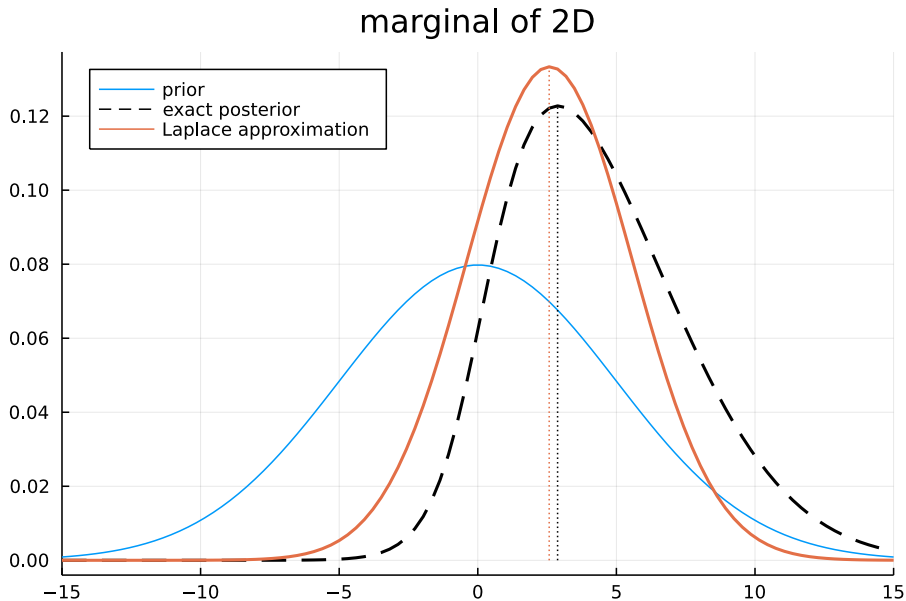
$$p(f | y) \approx \mathcal{N}(f | \hat{f}, -(\mathrm{d}^2 \log h / \mathrm{d} f^2)^{-1}) = q(f)$$



Laplace in 2D example



Laplace in 2D: marginals



Marginal likelihood approximation (I)

- Finally, we need the marginal likelihood in order to do model selection

$$\begin{aligned} p(\mathbf{y}) &= \int p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f}) \, d\mathbf{f} \\ &= \int \exp [\log p(\mathbf{y} \mid \mathbf{f}) + \log p(\mathbf{f})] \, d\mathbf{f} \end{aligned}$$

- Let's define $\psi(\mathbf{f}) = \log h(\mathbf{f}) = \log(p(\mathbf{y} \mid \mathbf{f})p(\mathbf{f}))$

$$\psi(\mathbf{f}) = \log p(\mathbf{y} \mid \mathbf{f}) + \log p(\mathbf{f}) = \log p(\mathbf{y} \mid \mathbf{f}) - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f}$$

- Second order Taylor approximation around the mode $\hat{\mathbf{f}}$

$$\psi(\mathbf{f}) \approx \psi(\hat{\mathbf{f}}) - \frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \mathbf{A} (\mathbf{f} - \hat{\mathbf{f}})$$

- Substituting back

$$p(\mathbf{y}) \approx q(\mathbf{y}) = \int \exp \left[\psi(\hat{\mathbf{f}}) - \frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \mathbf{A} (\mathbf{f} - \hat{\mathbf{f}}) \right] d\mathbf{f}$$

Marginal likelihood approximation (II)

- We have

$$\begin{aligned} p(\mathbf{y}) \approx q(\mathbf{y}) &= \int \exp \left[\psi(\hat{\mathbf{f}}) - \frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \mathbf{A} (\mathbf{f} - \hat{\mathbf{f}}) \right] d\mathbf{f} \\ &= \int \exp \left[\psi(\hat{\mathbf{f}}) \right] \exp \left[-\frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \mathbf{A} (\mathbf{f} - \hat{\mathbf{f}}) \right] d\mathbf{f} \end{aligned}$$

- $\exp \left[\psi(\hat{\mathbf{f}}) \right]$ does not depend on \mathbf{f} :

$$p(\mathbf{y}) \approx q(\mathbf{y}) = \exp \left[\psi(\hat{\mathbf{f}}) \right] \int \exp \left[-\frac{1}{2} (\mathbf{f} - \hat{\mathbf{f}})^\top \mathbf{A} (\mathbf{f} - \hat{\mathbf{f}}) \right] d\mathbf{f}$$

- The integral evaluates to the normalization constant of a Gaussian

$$p(\mathbf{y}) \approx q(\mathbf{y}) = \exp \left[\psi(\hat{\mathbf{f}}) \right] (2\pi)^{\frac{N}{2}} |\mathbf{A}^{-1}|^{\frac{1}{2}}$$

- We substitute in the expression for $\exp \left[\psi(\hat{\mathbf{f}}) \right]$:

$$q(\mathbf{y}) = \exp \left[\log p(\mathbf{y} | \hat{\mathbf{f}}) - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \hat{\mathbf{f}}^\top \mathbf{K}^{-1} \hat{\mathbf{f}} \right] (2\pi)^{\frac{N}{2}} |\mathbf{A}^{-1}|^{\frac{1}{2}}$$

Marginal likelihood approximation (III)

- Taking the log of $q(\mathbf{y})$

$$\begin{aligned}\log q(\mathbf{y}) &= \log p(\mathbf{y}|\hat{\mathbf{f}}) - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \hat{\mathbf{f}}^\top \mathbf{K}^{-1} \hat{\mathbf{f}} + \frac{N}{2} \log(2\pi) + \frac{1}{2} \log |\mathbf{A}^{-1}| \\ &= \log p(\mathbf{y} | \hat{\mathbf{f}}) - \frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \hat{\mathbf{f}}^\top \mathbf{K}^{-1} \hat{\mathbf{f}} + \frac{1}{2} \log |\mathbf{A}^{-1}|\end{aligned}$$

- We can now use the fact that $|\mathbf{A}^{-1}| = |\mathbf{A}|^{-1}$ to get

$$\log q(\mathbf{y}) = \log p(\mathbf{y} | \hat{\mathbf{f}}) - \frac{1}{2} \hat{\mathbf{f}}^\top \mathbf{K}^{-1} \hat{\mathbf{f}} - \frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \log |\mathbf{A}|$$

- Finally, recall that $\mathbf{A} = \mathbf{K}^{-1} + \mathbf{W}$

$$\log q(\mathbf{y}) = \log p(\mathbf{y} | \hat{\mathbf{f}}) - \frac{1}{2} \hat{\mathbf{f}}^\top \mathbf{K}^{-1} \hat{\mathbf{f}} - \frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \log |\mathbf{K}^{-1} + \mathbf{W}|$$

- We optimize $\log q(\mathbf{y})$ using gradient based methods to choose hyperparameters

Laplace approximation: important properties

- Find mode: Newton's method
 - Match curvature (Hessian) at mode
 - “Point estimate++”
- + Simple, fast
- Poor approximation if mode is not representative (e.g., Bernoulli)
 - May not converge for non-log-concave likelihoods

Choosing μ and Σ for $q(\mathbf{f})$

$$p(\mathbf{f} | \mathbf{y}) \approx q(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mu = ?, \Sigma = ?)$$

locally: match mean &
variance at point

globally: minimise divergence

Laplace
approximation

**Variational
inference (VI)**

Expectation
Propagation (EP)

Why variational inference

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

Variational inference: The big picture

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

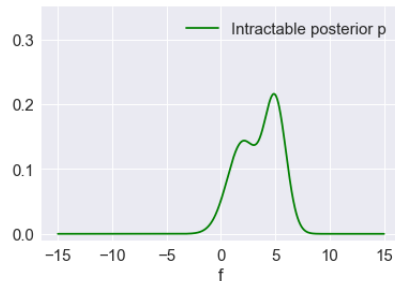
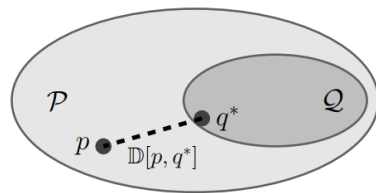
- 1 Define some “simple” family of distributions \mathcal{Q} .
- 2 Define some way to compute a “distance” $\mathbb{D}[p, q]$ between intractable distribution p and each distribution $q \in \mathcal{Q}$

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

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

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

- 4 Use q^* as an approximation of p



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



How to “measure distances” between distributions?

Here: *Kullback-Leibler divergence*

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

Important properties:

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

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

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

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

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

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

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

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

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

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

Likelihood term

Integral separates for a factorizing likelihood:

$$\begin{aligned} & \int q(\mathbf{f}) [\log p(\mathbf{y} | \mathbf{f})] d\mathbf{f} \\ &= \sum_{n=1}^N \int q(f_n) [\log p(y_n | f_n)] df_n \end{aligned}$$

Sum over 1D integrals

Each integral is a Gaussian expectation of the log likelihood

- Analytic for some (e.g., Exponential, Gamma, Poisson)
- Fast and accurate to approximate numerically (e.g., Gauss–Hermite quadrature)
- Monte Carlo (e.g., multi-class classification)

Take away #2: We can tractably optimize the bound for non-Gaussian likelihoods

Gauss–Hermite Quadrature

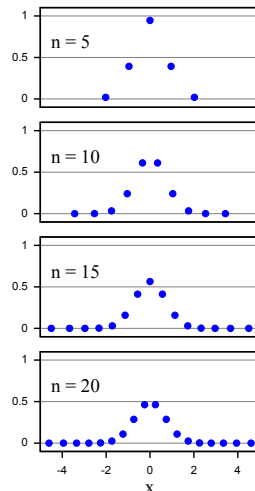
$$\int q(f_n) [\log p(y_n | f_n)] df_n, \quad q(f_n) = \mathcal{N}(m_n, S_n)$$

Gauss–Hermite quadrature can be applied:

$$\mathbb{E}_{q(f_n)} [\log p(y_n | f_n)] \approx \sum_{j=1}^C w_j \log p(y_n | f_j),$$

$$w_j = \frac{2^{C-1} C! \sqrt{\pi}}{C^2 [H_{C-1}(f_j)]^2}$$

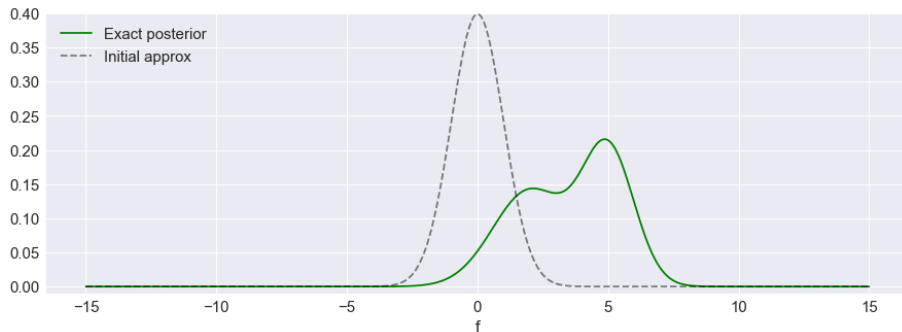
Gauss–Hermite is exact if $\log p(y_n | f_n)$ is polynomial of order less than C



1D Toy example I

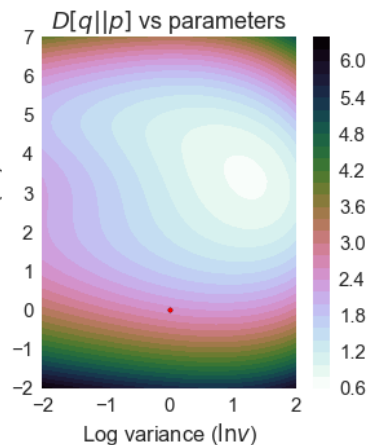
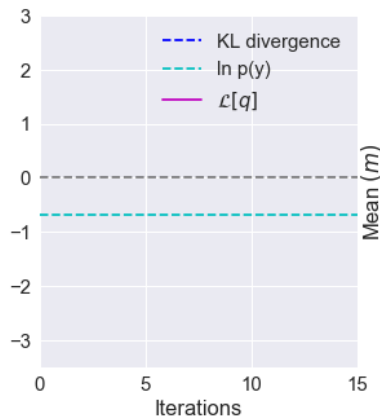
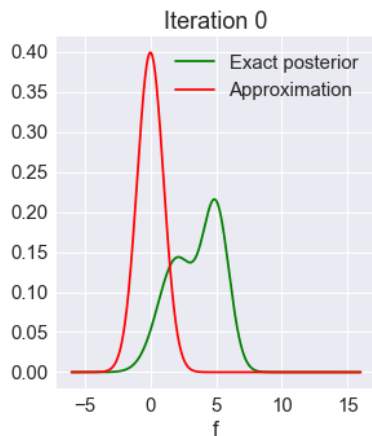
Assume model $p(y, f)$ with some intractable posterior $p(f | y)$

- Variational approximation for $p(f | y)$
- In 1D: \mathcal{Q} is the set of univariate Gaussians, i.e. $q_{\lambda}(x) = \mathcal{N}(x | m, v)$, and $\lambda = \{m, v\}$
- Initialization: $q(f) = \mathcal{N}(f | 0, 1)$

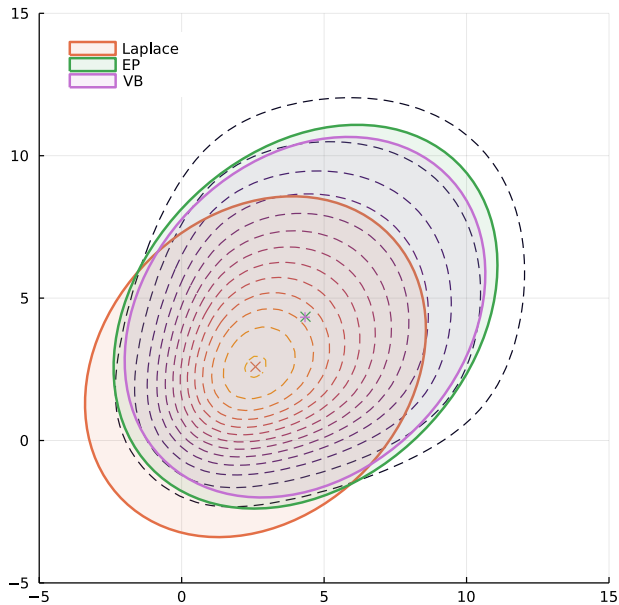


1D Toy example II

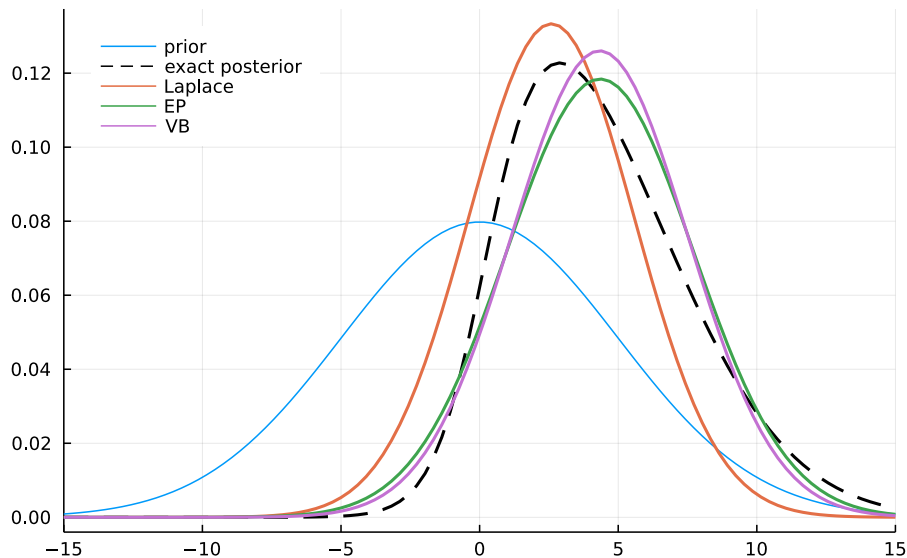
- Gradient ascent: $\lambda_{i+1} = \lambda_i + \eta \nabla_{\lambda} \mathcal{L}[\lambda]$
- $\log p(\mathbf{y}) = \mathcal{L}[\lambda] + \mathbb{D}[q_{\lambda}(\mathbf{f}) \parallel p(\mathbf{f} \mid \mathbf{y})] \geq \mathcal{L}[\lambda]$



Comparison 2D



marginal of 2D



Variational Bayes: Important properties

- Principled: directly minimising divergence from true posterior
- Mode-seeking (e.g., multi-modal posterior: fits just one, if q is unimodal)
- + Minimises a true lower bound \rightarrow convergence
- Underestimates variance

Section 4

Conclusion

- 1 Beyond Gaussian noise
- 2 Inference for arbitrary likelihoods
- 3 Approximating the intractable
- 4 Conclusion

Posterior distribution for f_*

- Now we know how to compute the approximate posterior $q(\mathbf{f} \mid \mathbf{y})$
- Let's now consider the posterior distribution for f_*

$$\begin{aligned} p(f_* \mid y) &= \int p(f_* \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{y}) \, d\mathbf{f} \\ &= \int \mathcal{N}\left(f_* \mid \mathbf{k}_{f_* \mathbf{f}} \mathbf{K}^{-1} \mathbf{f}, k_{f_* f_*} - \mathbf{k}_{f_* \mathbf{f}} \mathbf{K}^{-1} \mathbf{k}_{f_* \mathbf{f}}^\top\right) p(\mathbf{f} \mid \mathbf{y}) \, d\mathbf{f} \\ &\approx \int \mathcal{N}\left(f_* \mid \mathbf{k}_{f_* \mathbf{f}} \mathbf{K}^{-1} \mathbf{f}, k_{f_* f_*} - \mathbf{k}_{f_* \mathbf{f}} \mathbf{K}^{-1} \mathbf{k}_{f_* \mathbf{f}}^\top\right) \mathcal{N}\left(\mathbf{f} \mid \hat{\mathbf{f}}, \mathbf{A}^{-1}\right) \, d\mathbf{f} \\ &= \mathcal{N}\left(f_* \mid \underbrace{\mathbf{k}_{f_* \mathbf{f}} \mathbf{K}^{-1} \hat{\mathbf{f}}}_{\mu_*}, \underbrace{k_{f_* f_*} - \mathbf{k}_{f_* \mathbf{f}} (\mathbf{K} + \mathbf{W}^{-1})^{-1} \mathbf{k}_{f_* \mathbf{f}}^\top}_{\sigma_*^2}\right) \\ &= \mathcal{N}(f_* \mid \mu_*, \sigma_*^2) \end{aligned}$$

Predictive distribution

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

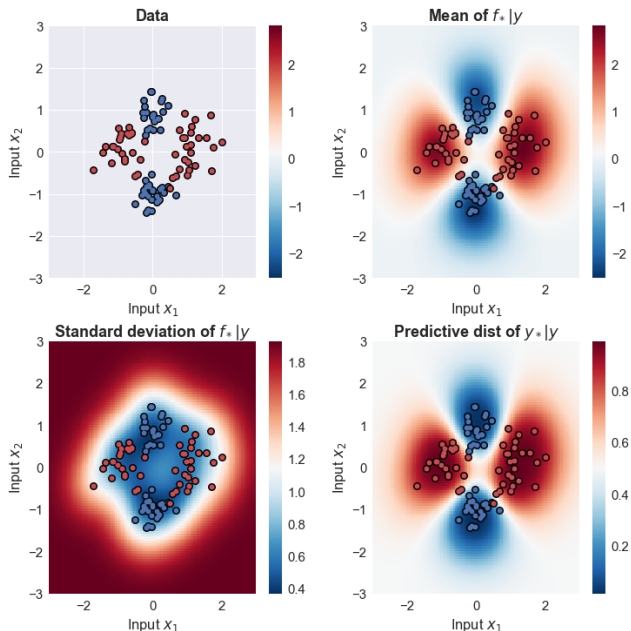
$$\begin{aligned} p(y_* = 1 | \mathbf{y}) &= \int p(y_* | f_*) p(f_* | \mathbf{y}) \mathrm{d}f_* \\ &= \int \phi(y_* \cdot f_*) p(f_* | \mathbf{y}) \mathrm{d}f_* \\ &\approx \int \phi(y_* \cdot f_*) q(f_*) \mathrm{d}f_* \\ &= \int \phi(y_* \cdot f_*) \mathcal{N}(f_* | \mu_*, \sigma_*^2) \mathrm{d}f_*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

Question

- What can we say about predictive distributions for y_* when μ_* is positive? (or negative?)
- How does uncertainty of posterior distribution of f_* influence the predictions for y_* ?
What happens as σ_*^2 approaches ∞ ?

Gaussian process classification example

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



End of today's lecture

- GPs can be used for all kinds of response variables
- Likelihood with parameters modulated by latent GP
- Non-Gaussian likelihood: non-Gaussian posterior, inference no longer exact
- Approximations: We covered Laplace and Variational Inference

