



UNIVERSITY OF  
CAMBRIDGE

# Machine Learning and the Physical World

## Lecture 4 : Inference and Latent Variables

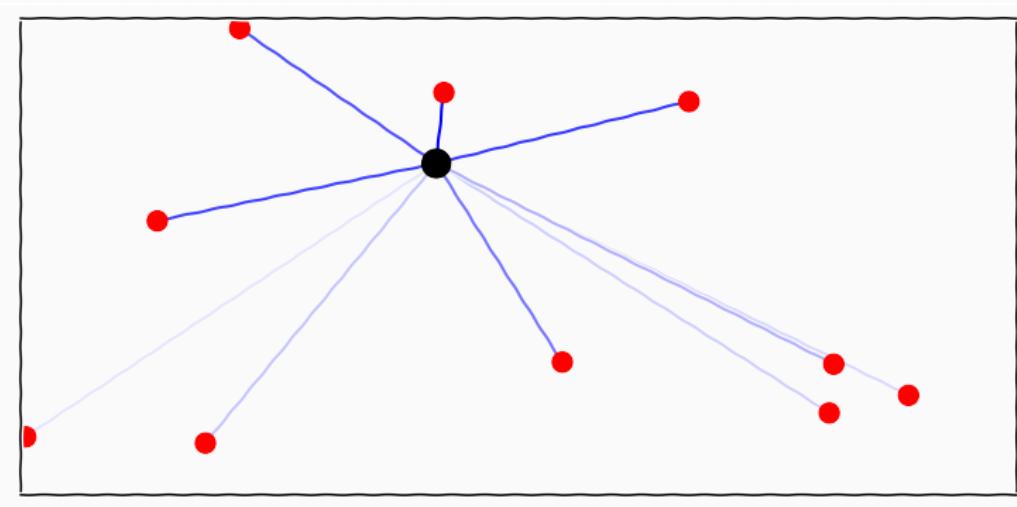
---

Carl Henrik Ek - [che29@cam.ac.uk](mailto:che29@cam.ac.uk)

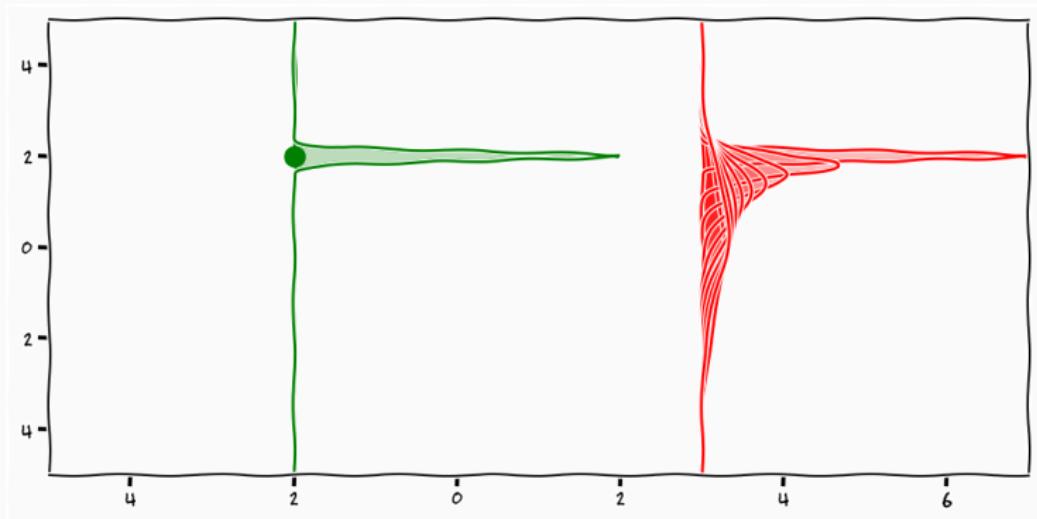
16th of October, 2020

<http://carlhenrik.com>

# Non-parametrics



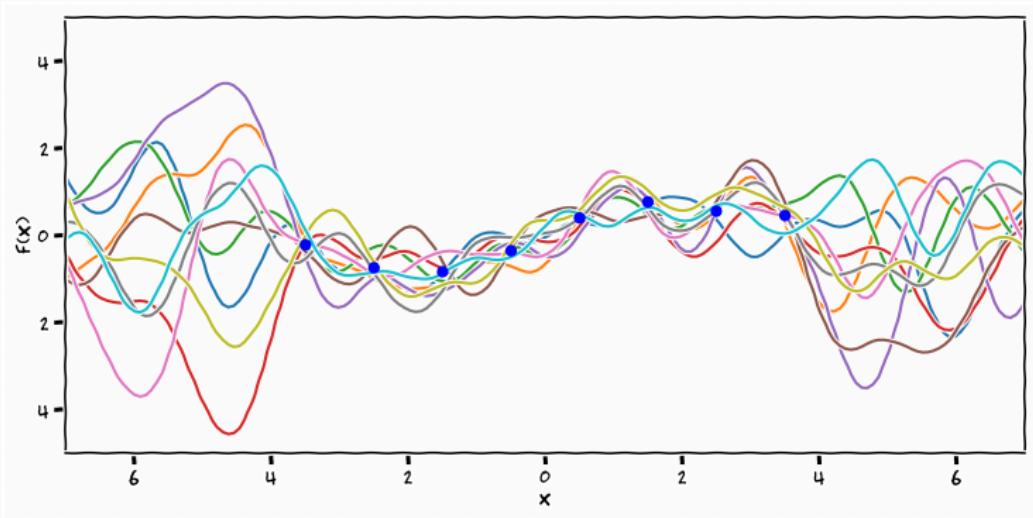
# Non-parametric functions



## Gaussian Processes: Formalism II

$$p(\mathbf{f}) = \mathcal{N} \left( \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_N \\ \vdots \end{bmatrix} \middle| \begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_N) \\ \vdots \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) & \dots \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) & \dots \\ \vdots & \vdots & \dots & \vdots & \ddots \end{bmatrix} \right)$$

# Posterior Samples



## Gaussian processes

$$\begin{array}{ccc} \mathcal{GP}(\cdot, \cdot) & M \in \mathbb{R}^{\infty \times N} & \mathcal{N}(\cdot, \cdot) \\ & \rightarrow & \\ \infty & & N \end{array}$$

The Gaussian distribution is the projection of the infinite Gaussian process

# Predictive Posteriors

- Gaussian Process - *Non-parametric formulation*

$$p(y_* \mid \mathbf{y}, \mathbf{x}_*, \mathbf{X}, \theta) = \mathcal{N}(y_* \mid \mu_*, K_*)$$

$$\mu_* = k(x_*, x)k(x, x)^{-1}\mathbf{y}$$

$$K_* = k(x_*, x_*) - k(x_*, x)k(x, x)^{-1}k(x, x_*)$$

# Predictive Posteriors

- Gaussian Process - *Non-parametric formulation*

$$p(y_* \mid \mathbf{y}, \mathbf{x}_*, \mathbf{X}, \theta) = \mathcal{N}(y_* \mid \mu_*, K_*)$$

$$\mu_* = k(x_*, x)k(x, x)^{-1}\mathbf{y}$$

$$K_* = k(x_*, x_*) - k(x_*, x)k(x, x)^{-1}k(x, x_*)$$

- Linear Regression - *Parametric formulation*

$$\begin{aligned} p(y_* | \mathbf{y}, \mathbf{x}_*, \mathbf{X}, \alpha, \beta) &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}((y_* \mid \mu_*(x_*), \Sigma_*(x_*))) \end{aligned}$$

$$\mu_* = (\beta(\alpha\mathbf{I} + \beta\mathbf{x}^T\mathbf{x})^{-1}\mathbf{x}\mathbf{y})^T \mathbf{x}_*$$

$$= \mathbf{x}_*^T \mathbf{x} (\beta(\alpha\mathbf{I} + \beta\mathbf{x}^T\mathbf{x}))^{-1} \mathbf{y}$$

$$\Sigma_* = \frac{1}{\beta} + \mathbf{x}_*^T (\alpha\mathbf{I} + \beta\mathbf{x}^T\mathbf{x})^{-1} \mathbf{x}$$

## Regression Models: Mean

---

$$k(x_*, x)k(x, x)^{-1}\mathbf{y}$$

$$\mathbf{x}_*^T \mathbf{x} \left( \beta(\alpha \mathbf{I} + \beta \mathbf{x}^T \mathbf{x}) \right)^{-1} \mathbf{y}$$

- both means are linear with respect to the output data

## Regression Models: Mean

$$k(x_*, x)k(x, x)^{-1}\mathbf{y} \quad \mathbf{x}_*^T \mathbf{x} \left( \beta(\alpha\mathbf{I} + \beta\mathbf{x}^T \mathbf{x}) \right)^{-1} \mathbf{y}$$

- both means are linear with respect to the output data
- how about if we take a basis function approach such that

$$\mathbf{x}_*^T \mathbf{x} = \Phi(\mathbf{x}_*)^T \Phi(\mathbf{x})$$

## Regression Models: Mean

$$k(x_*, x)k(x, x)^{-1}\mathbf{y}$$

$$\mathbf{x}_*^T \mathbf{x} \left( \beta(\alpha \mathbf{I} + \beta \mathbf{x}^T \mathbf{x}) \right)^{-1} \mathbf{y}$$

- both means are linear with respect to the output data
- how about if we take a basis function approach such that  
 $\mathbf{x}_*^T \mathbf{x} = \Phi(\mathbf{x}_*)^T \Phi(\mathbf{x})$
- if we think of  $\mathbf{x}$  as the center of each basis function

## Regression Models: Mean

$$k(x_*, x)k(x, x)^{-1}\mathbf{y} \quad \mathbf{x}_*^T \mathbf{x} (\beta(\alpha\mathbf{I} + \beta\mathbf{x}^T \mathbf{x}))^{-1} \mathbf{y}$$

- both means are linear with respect to the output data
- how about if we take a basis function approach such that  
 $\mathbf{x}_*^T \mathbf{x} = \Phi(\mathbf{x}_*)^T \Phi(\mathbf{x})$
- if we think of  $\mathbf{x}$  as the center of each basis function
  - a basis function per data point

## Regression Models: Mean

$$k(x_*, x)k(x, x)^{-1}\mathbf{y} \quad \mathbf{x}_*^T \mathbf{x} \left( \beta(\alpha\mathbf{I} + \beta\mathbf{x}^T \mathbf{x}) \right)^{-1} \mathbf{y}$$

- both means are linear with respect to the output data
- how about if we take a basis function approach such that  
 $\mathbf{x}_*^T \mathbf{x} = \Phi(\mathbf{x}_*)^T \Phi(\mathbf{x})$
- if we think of  $\mathbf{x}$  as the center of each basis function
  - a basis function per data point
- if we could parametrise  $\Phi(\mathbf{x}_*)^T \Phi(\mathbf{x}) = k(\mathbf{x}_*, \mathbf{x})$  as a function

## Regression Models: Mean

$$k(x_*, x)k(x, x)^{-1}\mathbf{y} \quad \mathbf{x}_*^T \mathbf{x} (\beta(\alpha\mathbf{I} + \beta\mathbf{x}^T \mathbf{x}))^{-1} \mathbf{y}$$

- both means are linear with respect to the output data
- how about if we take a basis function approach such that  
 $\mathbf{x}_*^T \mathbf{x} = \Phi(\mathbf{x}_*)^T \Phi(\mathbf{x})$
- if we think of  $\mathbf{x}$  as the center of each basis function
  - a basis function per data point
- if we could parametrise  $\Phi(\mathbf{x}_*)^T \Phi(\mathbf{x}) = k(\mathbf{x}_*, \mathbf{x})$  as a function
- this leads to the interpretation of GPs as infinite basis functions

## The Marginal Likelihood

---

$$p(\theta | y) = \frac{p(y | \theta)p(\theta)}{\int p(y | \theta)p(\theta)d\theta}$$

**Likelihood** How much **evidence** is there in the data for a specific hypothesis

**Prior** What are my beliefs about different hypothesis

**Posterior** What is my **updated** belief after having seen data

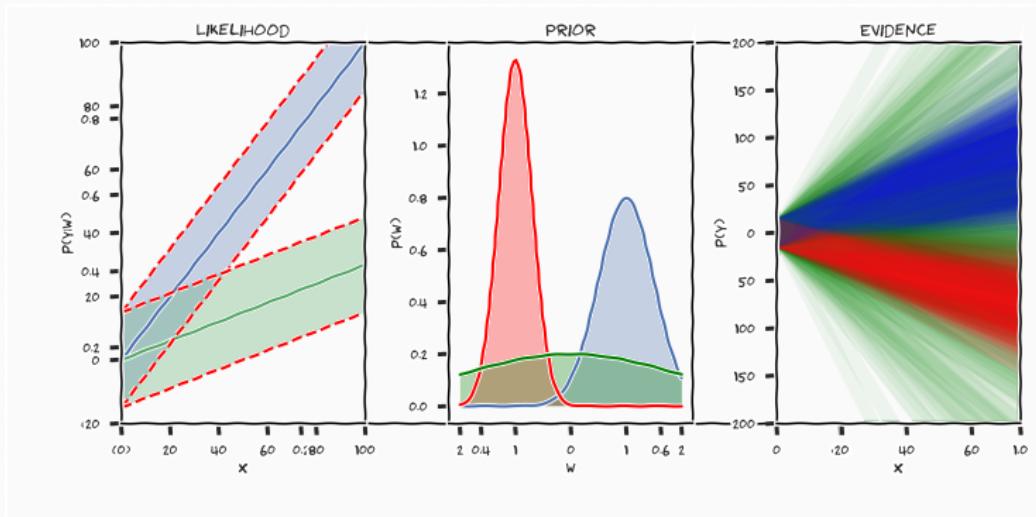
**Evidence** What is my belief about the data

## The Compute: Evidence

---

$$\int p(y \mid \theta)p(\theta)d\theta$$

# Regression Model



# Marginalisation



*Next time you want to give your friends a compliment, tell them that you have completely **marginalised** them from your life*

# Regression Models

---

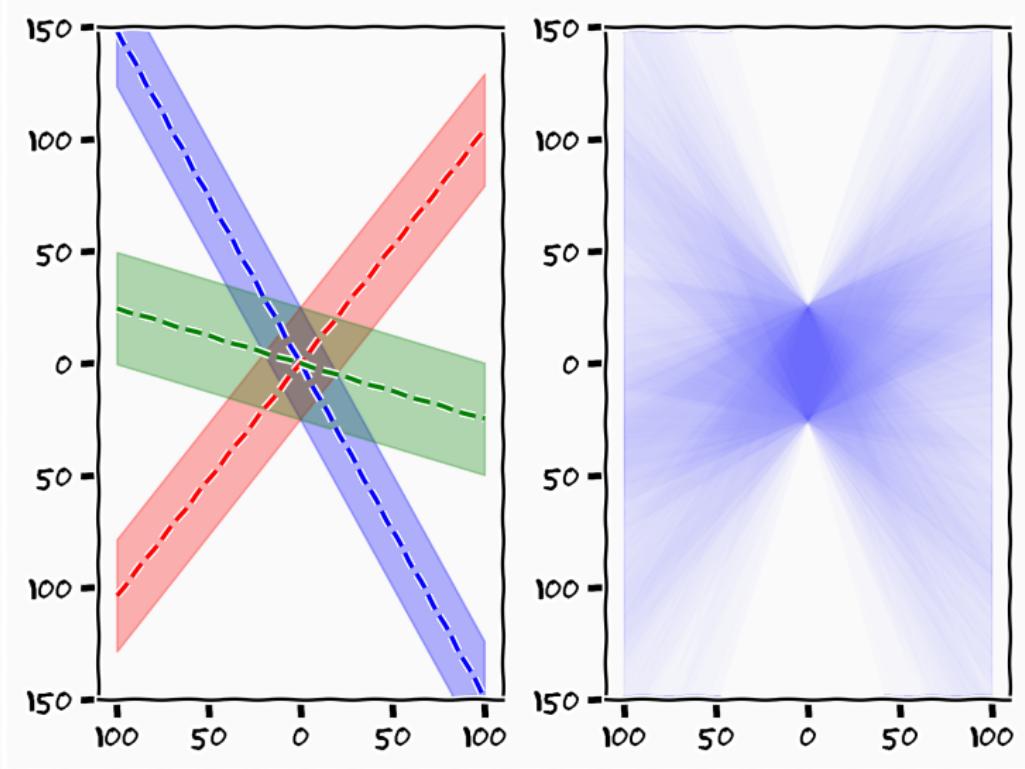
## Linear Model

$$p(y_i|x_i, \mathbf{w}) = \mathcal{N}(w_0 + w_1 \cdot x_i, \beta^{-1})$$

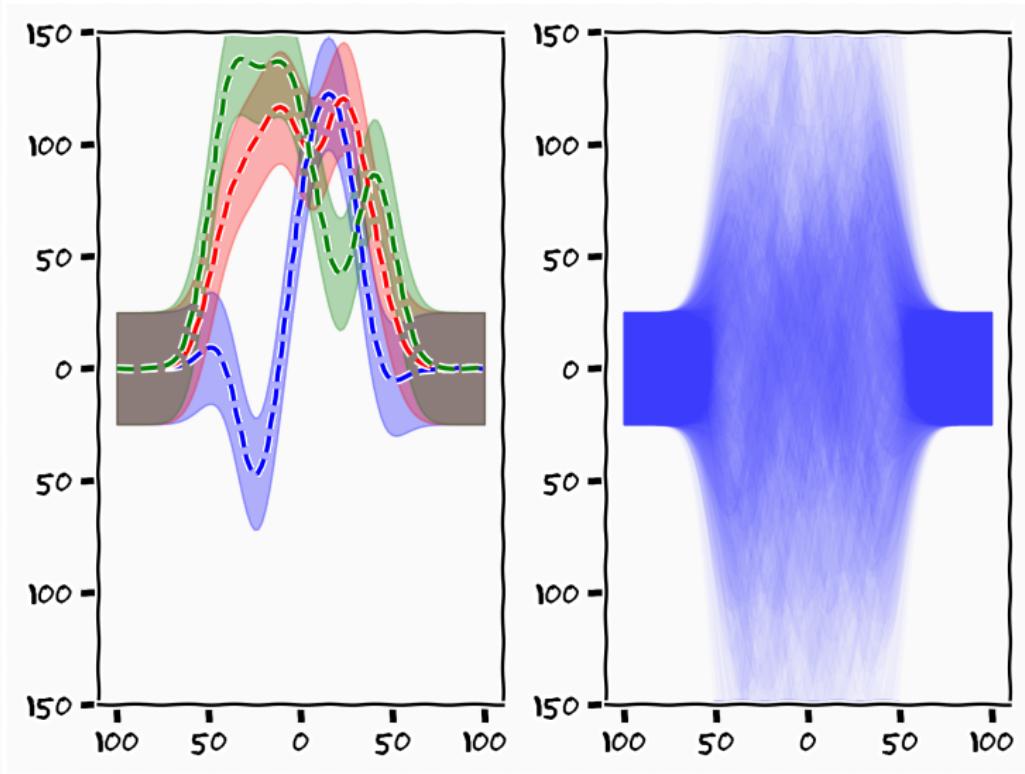
## Basis function

$$p(y_i|x_i, \mathbf{w}) = \mathcal{N}\left(\sum_{i=1}^6 w_i \phi(x_i), \beta^{-1}\right)$$

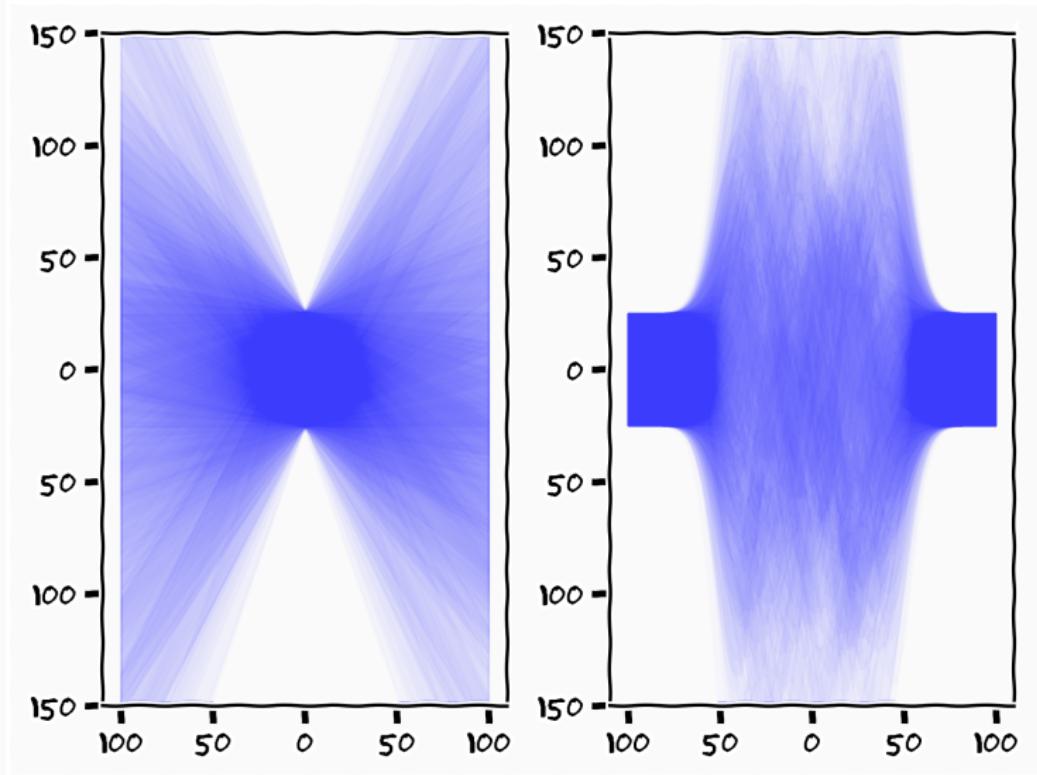
# Model 1



## Model 2



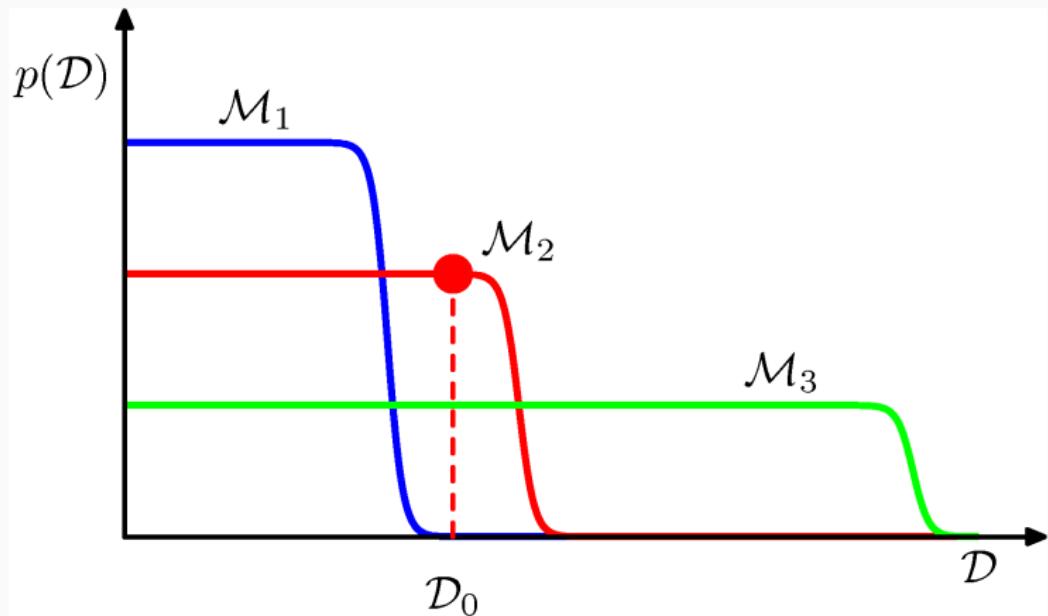
# Evidence



# Probabilities are a zero-sum game



# Model Selection<sup>1</sup>



<sup>1</sup>David MacKay PhD Thesis

# Occams Razor

---



# Occams Razor

---

## **Definition (Occams Razor)**

"All things being equal, the simplest solution tends to be the best one"

– William of Ockham

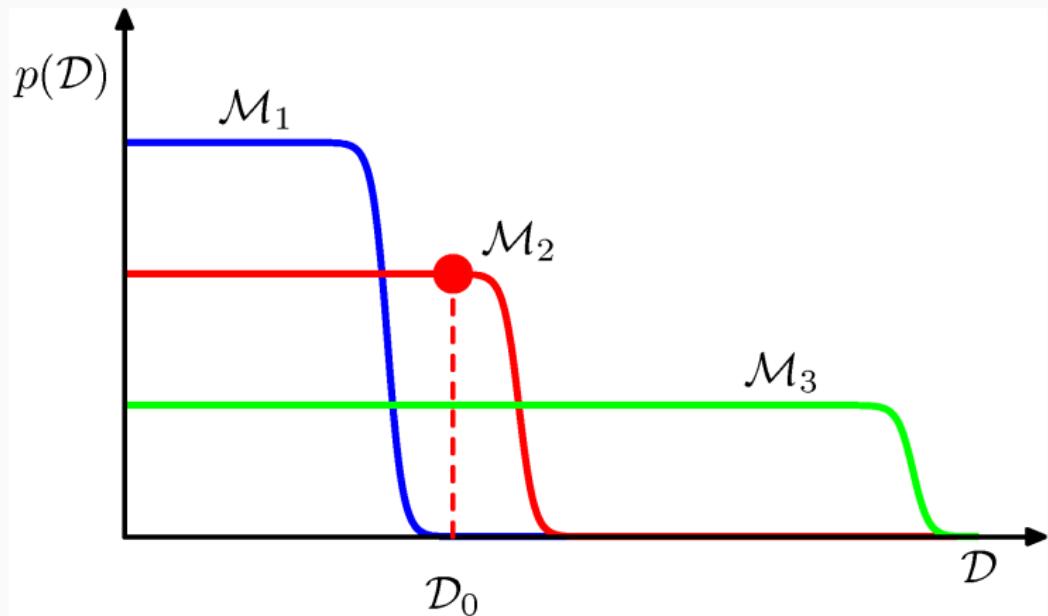
# What is Simple?<sup>2</sup>



---

<sup>2</sup><https://www.imdb.com/title/tt8132700/>

# Model Selection<sup>1</sup>



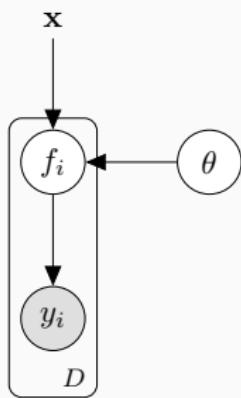
<sup>1</sup>David MacKay PhD Thesis

## The Compute: Evidence

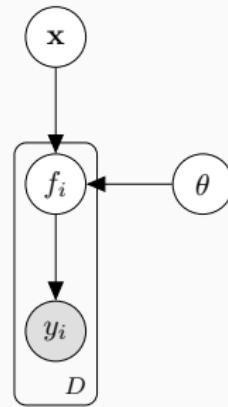
---

$$\int p(y \mid \theta) p(\theta) d\theta$$

# Unsupervised Learning



$$p(y|x)$$



$$p(y)$$

# Gaussian Process Latent Variable Model <sup>3</sup>

- Regression

$$p(y \mid x) = \int p(y \mid f)p(f \mid x)df$$

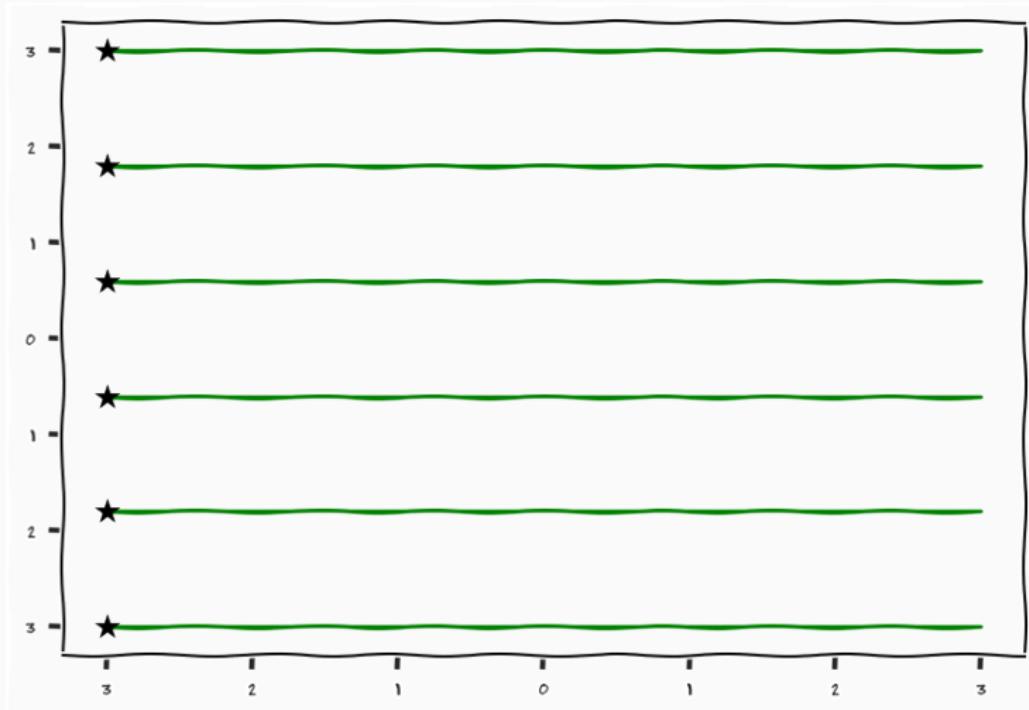
- "Unsupervised" Learning

$$p(y) = \int p(y \mid f)p(f \mid x)p(x)dfdx$$

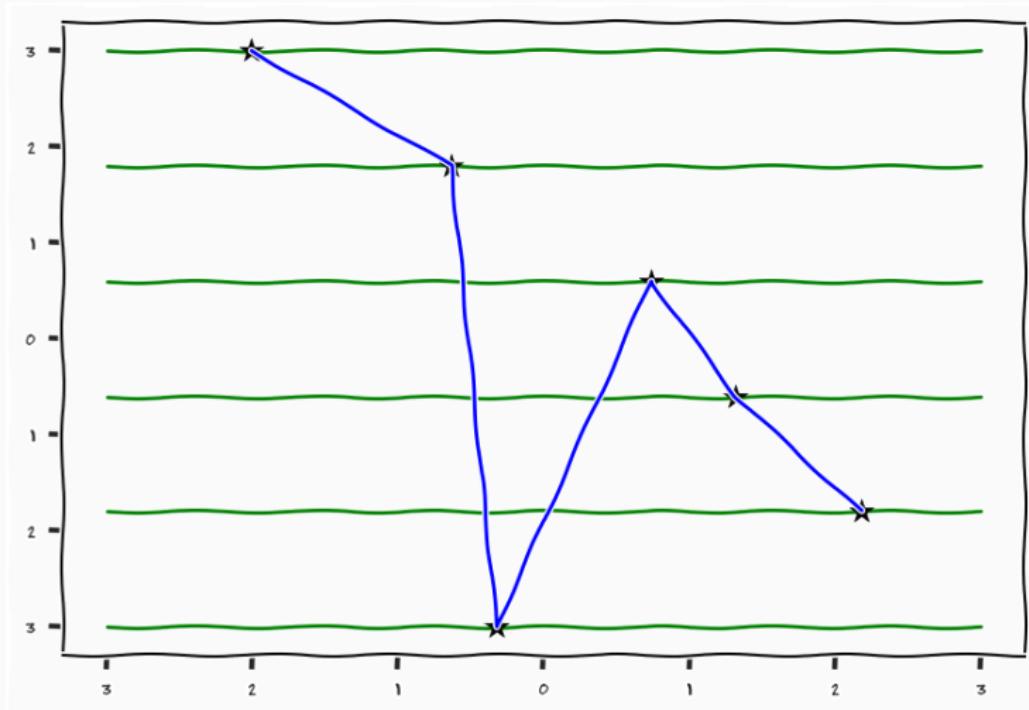
---

<sup>3</sup>Lawrence, 2005

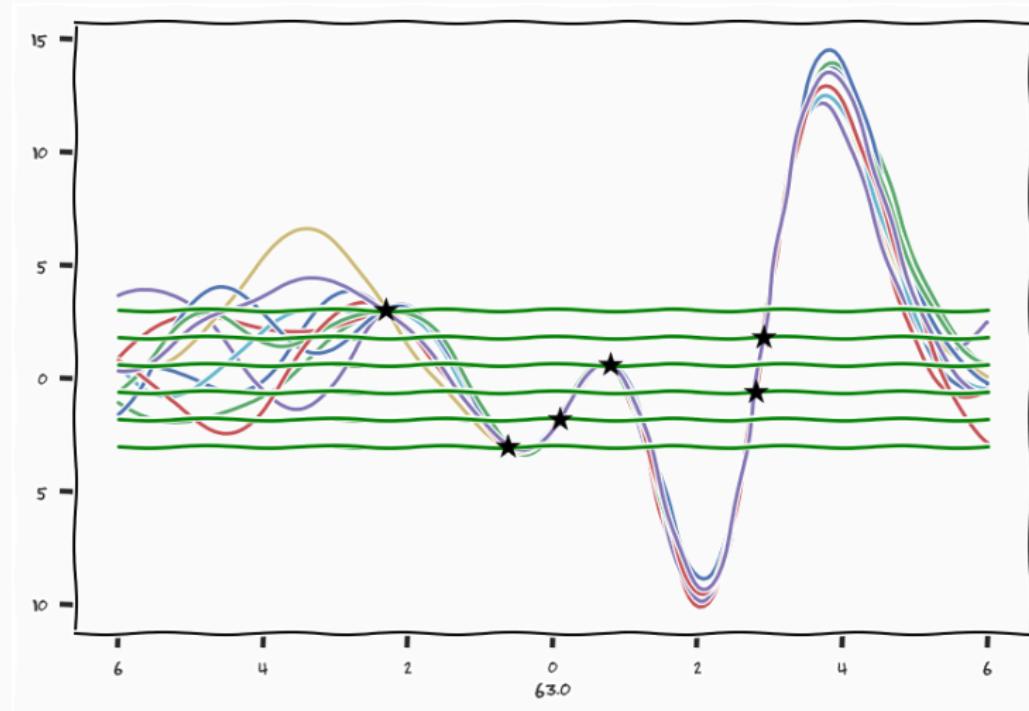
# Unsupervised Learning



# Unsupervised Learning

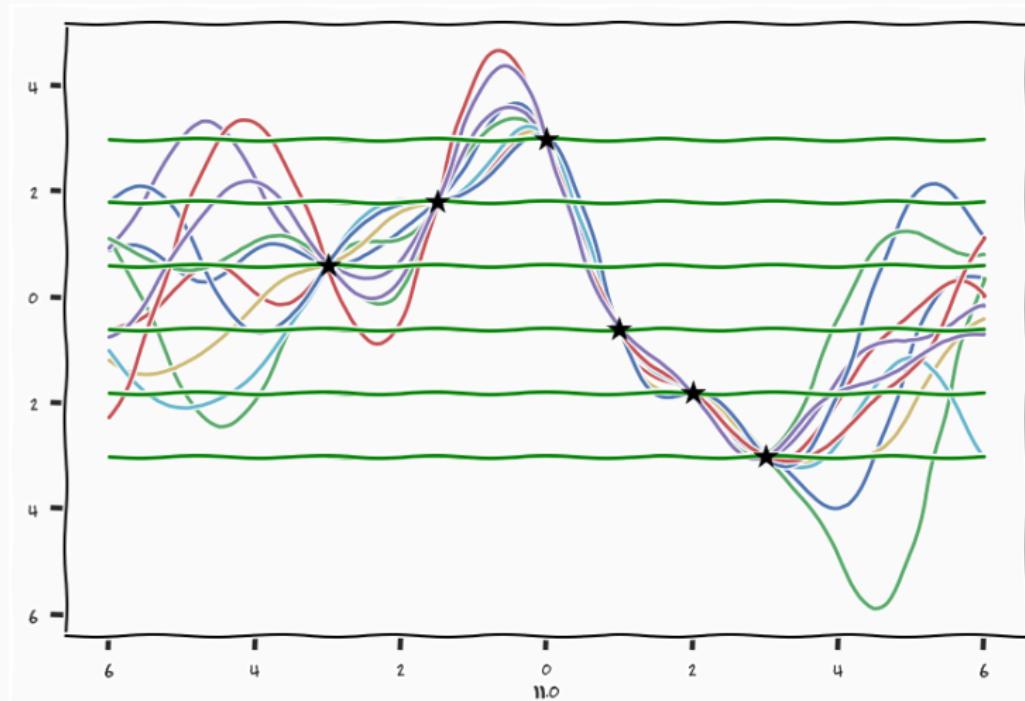


# Gaussian Process Latent Variable Model <sup>4</sup>



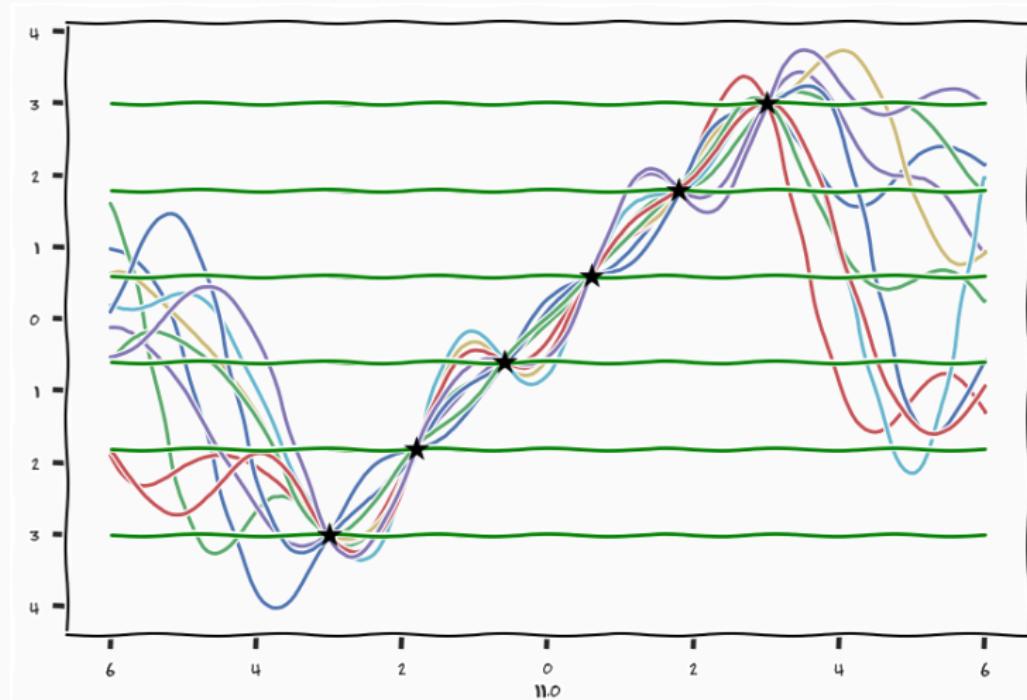
<sup>4</sup>Lawrence, 2005

# Gaussian Process Latent Variable Model <sup>4</sup>



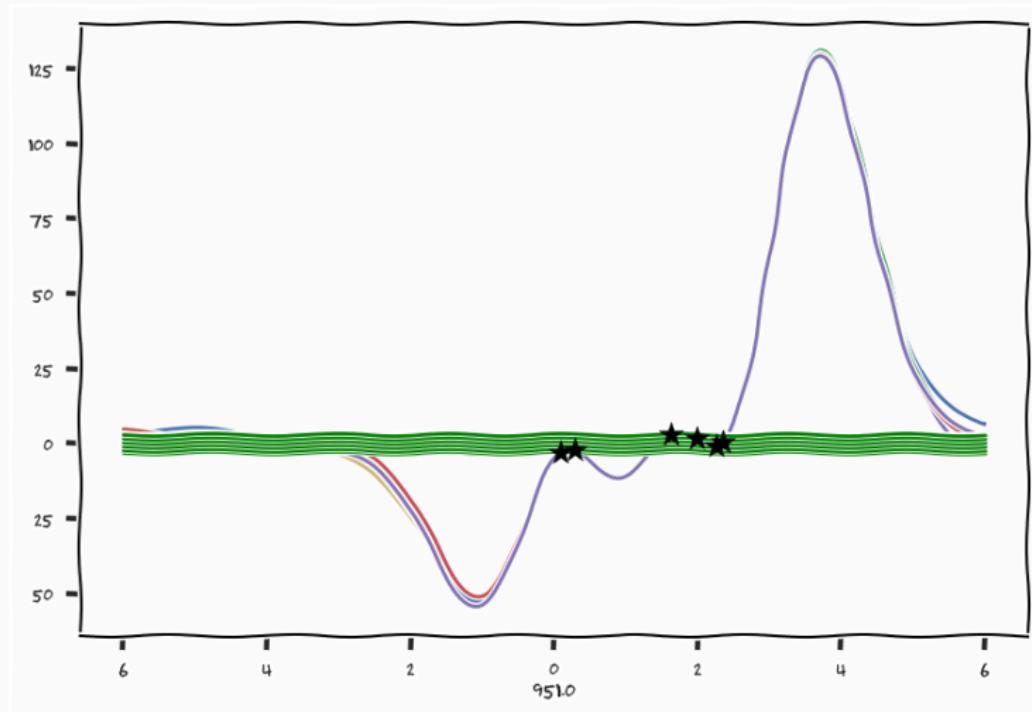
<sup>4</sup>Lawrence, 2005

# Gaussian Process Latent Variable Model <sup>4</sup>



<sup>4</sup>Lawrence, 2005

# Gaussian Process Latent Variable Model <sup>4</sup>



<sup>4</sup>Lawrence, 2005

# Priors

---

$$p(y) = \int p(y|f)p(f|x)p(x)dfdx$$

## 1. Priors that makes sense

$p(f)$  describes our belief/assumptions and defines our notion of complexity in the function

$p(x)$  expresses our belief/assumptions and defines our notion of complexity in the latent space

## 2. Now lets churn the handle

## Relationship between $x$ and data

$$p(y) = \int p(y|f)p(f|x)p(x)dfdx$$

- GP prior

$$p(f|x) \sim \mathcal{N}(0, K) \propto e^{-\frac{1}{2}(f^T K^{-1} f)}$$

$$K_{ij} = e^{-(x_i - x_j)^T M^T M (x_i - x_j)}$$

## Relationship between $x$ and data

$$p(y) = \int p(y|f)p(f|x)p(x)dfdx$$

- GP prior

$$p(f|x) \sim \mathcal{N}(0, K) \propto e^{-\frac{1}{2}(f^T K^{-1} f)}$$

$$K_{ij} = e^{-(x_i - x_j)^T M^T M (x_i - x_j)}$$

- Likelihood

$$p(y|f) \sim N(y|f, \beta) \propto e^{-\frac{1}{2\beta} \text{tr}(y-f)^T (y-f)}$$

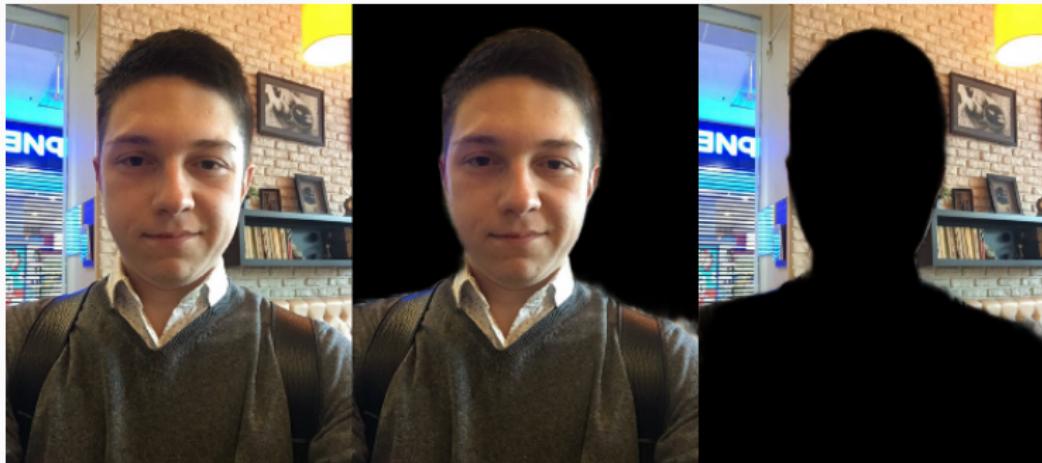
# Laplace Integration

---

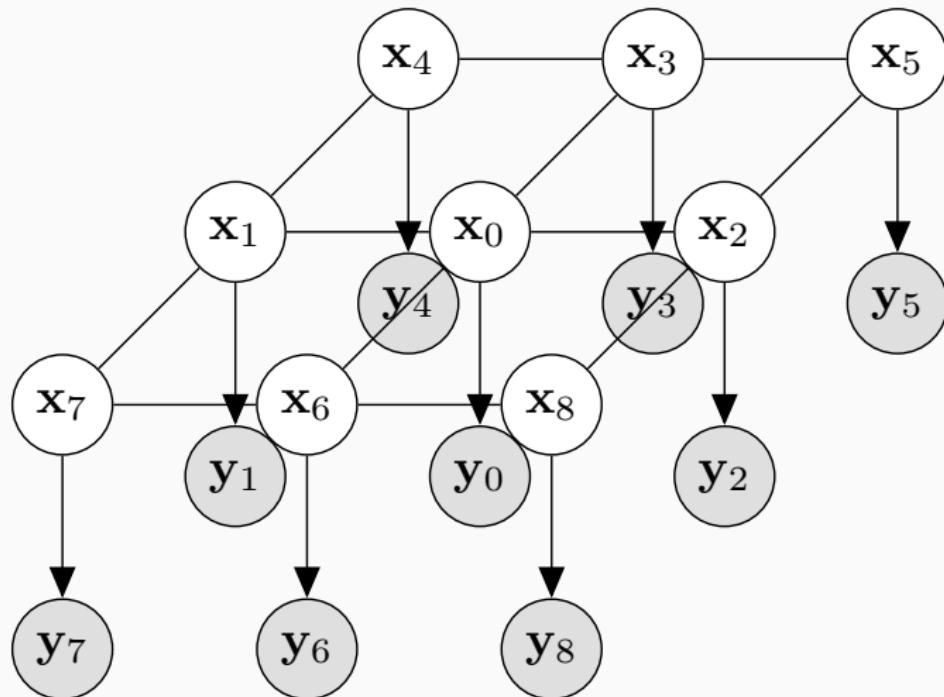


*"Nature laughs at the difficulties of integrations"*  
– Simon Laplace

# Image Segmentation



# Markov Random Field



# Markov Random Field

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) = \sum_i^N p(\mathbf{y}|\mathbf{x}_i)p(\mathbf{x}_i)$$

- $\mathbf{x}_i$  is a specific binary images

# Numbers

---

- Possible black and white 3 Megapixel images

$$2^{3145728}$$

# Numbers

---

- Possible black and white 3 Megapixel images

$$2^{3145728}$$

- Number of atoms in the universe

$$10^{80} \approx (2^{\frac{10}{3}})^{80} \approx 2^{267}$$

# Numbers

---

- Possible black and white 3 Megapixel images

$$2^{3145728}$$

- Number of atoms in the universe

$$10^{80} \approx (2^{\frac{10}{3}})^{80} \approx 2^{267}$$

- Age of the universe in seconds

$$4.35 \cdot 10^{17} \approx 2^{59}$$

# Intractability

---

- For most interesting problems computing the evidence is intractable
- Computational intractability: there are too many states to sum over
- Analytic: no closed form exists for the distribution

- We have been living in a fairy land so far, the world is not conjugate

# Today

---

- We have been living in a fairy land so far, the world is not conjugate
- We cannot give one lecture on approximate inference

- We have been living in a fairy land so far, the world is not conjugate
- We cannot give one lecture on approximate inference
- We hope to give you a flavour of approaches

- We have been living in a fairy land so far, the world is not conjugate
- We cannot give one lecture on approximate inference
- We hope to give you a flavour of approaches
  - Stochastic

## Stochastic Inference

---

## Basic Sampling

---

$$\mathbb{E}_{p(x)}[f] = \int f(x)p(x)dx \approx \frac{1}{L} \sum_{l=1}^L f(x^{(l)}) = \hat{f}$$

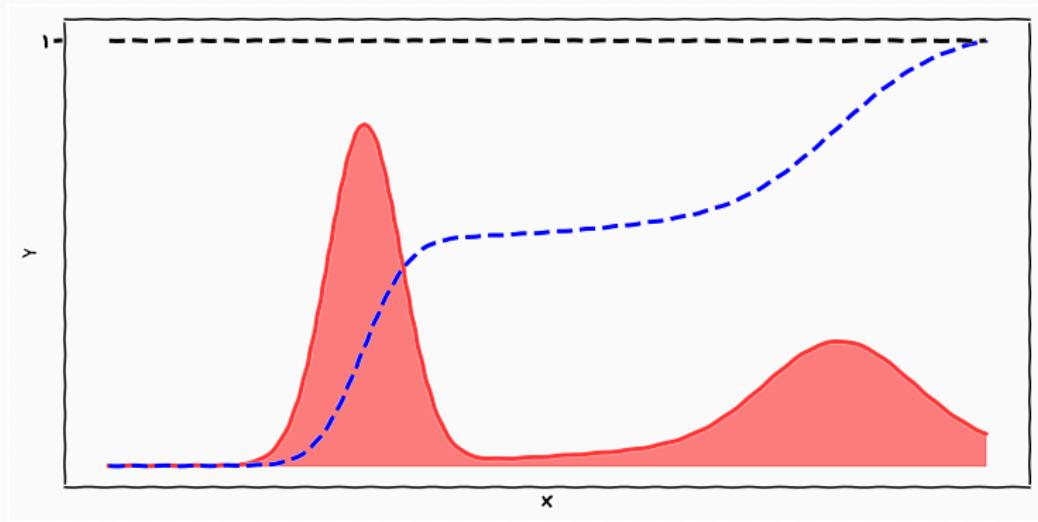
$$x^{(l)} \sim p(x)$$

$$\mathbb{E}[\hat{f}] = \mathbb{E}[f]$$

$$\text{var}[\hat{f}] = \frac{1}{L}\mathbb{E}[(f(x) - \mathbb{E}[f])^2]$$

- Approximation not dependent on dimensionality of  $x$
- Variance of estimator shrinks with number of samples

# Basic Sampling

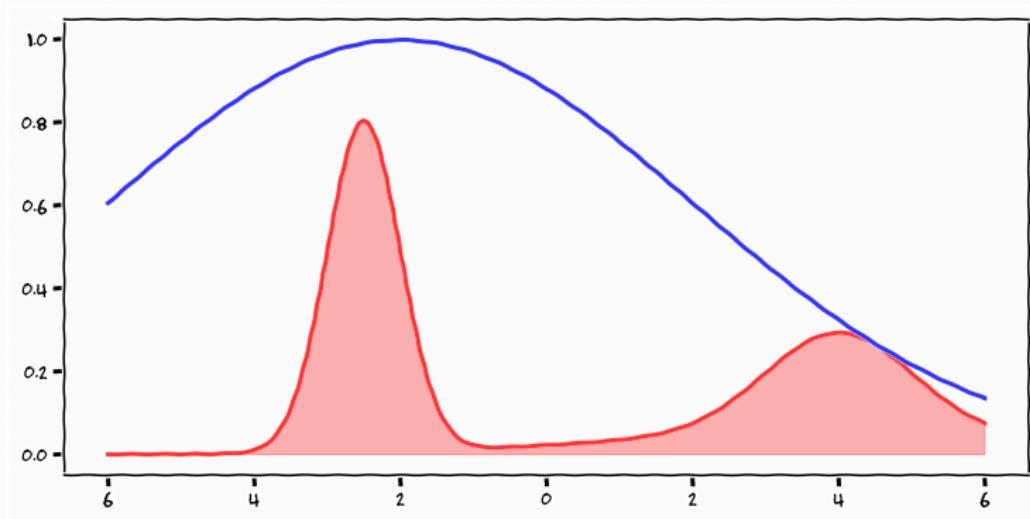


# Sampling

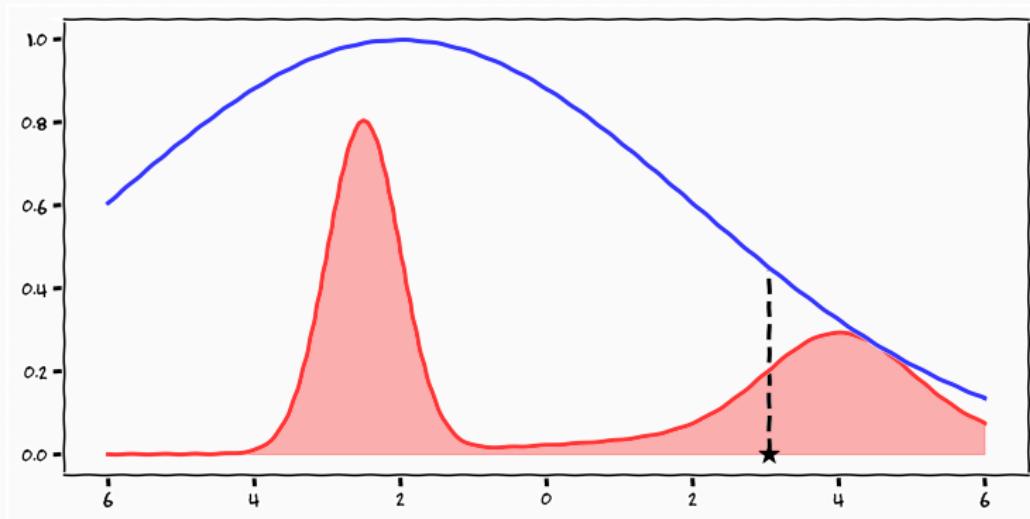
---

- We know how to transform samples from uniform to any distribution we can formulate the cumulative distribution
- Can we sample from distributions we do not know the form of?
  1. Rejection Sampling
  2. Importance Sampling

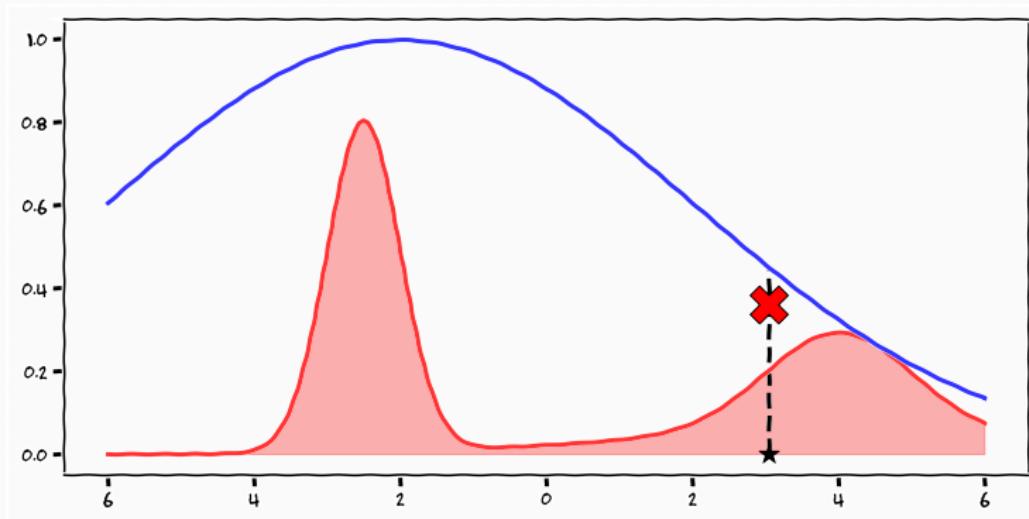
# Rejection Sampling



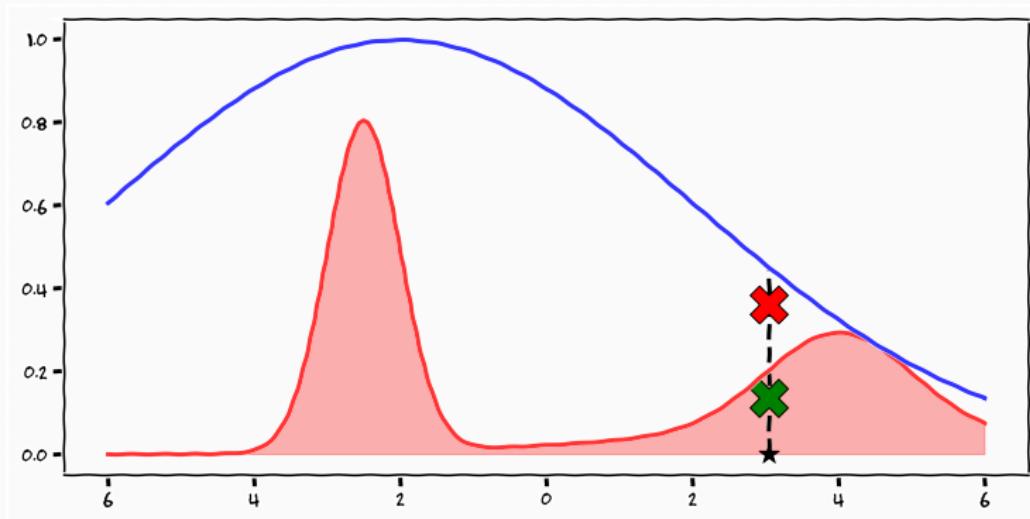
# Rejection Sampling



# Rejection Sampling



# Rejection Sampling



# Importance Sampling

---

$$\mathbb{E}_{p(x)}[f] = \int f(x)p(x)dx$$

# Importance Sampling

---

$$\mathbb{E}_{p(x)}[f] = \int f(x)p(x)dx = \int f(x)p(x)\frac{q(x)}{q(x)}dx$$

# Importance Sampling

---

$$\begin{aligned}\mathbb{E}_{p(x)}[f] &= \int f(x)p(x)dx = \int f(x)p(x)\frac{q(x)}{q(x)}dx \\ &= \int f(x)\frac{p(x)}{q(x)}q(x)dx\end{aligned}$$

# Importance Sampling

---

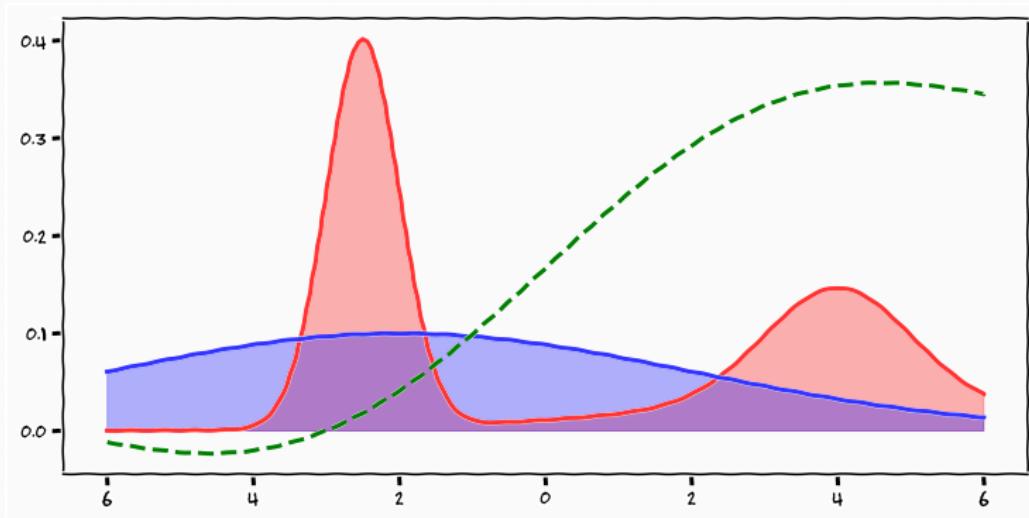
$$\begin{aligned}\mathbb{E}_{p(x)}[f] &= \int f(x)p(x)dx = \int f(x)p(x)\frac{q(x)}{q(x)}dx \\ &= \int f(x)\frac{p(x)}{q(x)}q(x)dx = \mathbb{E}_{q(x)}\left[f(x)\frac{p(x)}{q(x)}\right]\end{aligned}$$

# Importance Sampling

---

$$\begin{aligned}\mathbb{E}_{p(x)}[f] &= \int f(x)p(x)dx = \int f(x)p(x)\frac{q(x)}{q(x)}dx \\ &= \int f(x)\frac{p(x)}{q(x)}q(x)dx = \mathbb{E}_{q(x)}\left[f(x)\frac{p(x)}{q(x)}\right] \\ &\approx \frac{1}{L} \sum_{l=1}^L f(x^{(l)})\frac{p(x^{(l)})}{q(x^{(l)})}\end{aligned}$$

# Importance Sampling



# Markov Chain Monte Carlo

---

- Sample from a proposal distribution
- Remembers the state and samples from a conditional
- Can lead to much better exploration of the space

# Metropolis Sampling

---

1. start with state  $x^{(0)}$

# Metropolis Sampling

---

1. start with state  $x^{(0)}$
2. sample from conditional proposal distribution  $q(x^* \mid x^{(0)})$

# Metropolis Sampling

---

1. start with state  $x^{(0)}$
2. sample from conditional proposal distribution  $q(x^* \mid x^{(0)})$
3. compute acceptance probability

$$A(x^*, x^{(0)}) = \min \left( 1, \frac{\tilde{p}(x^*)}{\tilde{p}(x^{(0)})} \right)$$

# Metropolis Sampling

---

1. start with state  $x^{(0)}$
2. sample from conditional proposal distribution  $q(x^* \mid x^{(0)})$
3. compute acceptance probability

$$A(x^*, x^{(0)}) = \min \left( 1, \frac{\tilde{p}(x^*)}{\tilde{p}(x^{(0)})} \right)$$

4. Draw uniform random number  $u \sim \text{Uniform}(0, 1)$

# Metropolis Sampling

1. start with state  $x^{(0)}$
2. sample from conditional proposal distribution  $q(x^* \mid x^{(0)})$
3. compute acceptance probability

$$A(x^*, x^{(0)}) = \min \left( 1, \frac{\tilde{p}(x^*)}{\tilde{p}(x^{(0)})} \right)$$

4. Draw uniform random number  $u \sim \text{Uniform}(0, 1)$ 
  - if  $A(x^*, x^{(0)}) > u \rightarrow x^{(1)} = x^*$

# Metropolis Sampling

1. start with state  $x^{(0)}$
2. sample from conditional proposal distribution  $q(x^* \mid x^{(0)})$
3. compute acceptance probability

$$A(x^*, x^{(0)}) = \min \left( 1, \frac{\tilde{p}(x^*)}{\tilde{p}(x^{(0)})} \right)$$

4. Draw uniform random number  $u \sim \text{Uniform}(0, 1)$ 
  - if  $A(x^*, x^{(0)}) > u \rightarrow x^{(1)} = x^*$
  - otherwise reject  $x^*$  and start over

## Gibbs Sampling

---

- Often 1D samples are easy to get
- Gibbs sampling exploits this to create a very simple Markov Chain
- Sample each variable in turn conditioned on the others and cycle through

# Gibbs Sampling

---

1. Initialise  $x^{(0)}$

# Gibbs Sampling

---

1. Initialise  $x^{(0)}$
2. Pick single variable  $x_i \in x$

# Gibbs Sampling

---

1. Initialise  $x^{(0)}$
2. Pick single variable  $x_i \in x$
3. Formulate posterior  $p(x_i|x_{-i})$

# Gibbs Sampling

---

1. Initialise  $x^{(0)}$
2. Pick single variable  $x_i \in x$
3. Formulate posterior  $p(x_i|x_{-i})$
4. Sample from posterior

$$x_i^{(1)} \sim p(x_i|\mathbf{x}_{\neg i})$$

# Gibbs Sampling

---

1. Initialise  $x^{(0)}$
2. Pick single variable  $x_i \in x$
3. Formulate posterior  $p(x_i|x_{-i})$
4. Sample from posterior  
$$x_i^{(1)} \sim p(x_i|\mathbf{x}_{\neg i})$$
5. cycle through variables

# Why is this easier?

---

## Multivariate case

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

$$p(\mathbf{y}) = \sum_i p(\mathbf{y}|\mathbf{x}^{(i)})p(\mathbf{x}^{(i)})$$

# Why is this easier?

---

## Multivariate case

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{x})}{p(\mathbf{y})}$$
$$p(\mathbf{y}) = \sum_i p(\mathbf{y}|\mathbf{x}^{(i)})p(\mathbf{x}^{(i)})$$

## 1D case

$$p(x_i|\mathbf{x}_{\neg i}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x}_{\neg i}, \mathbf{y})}$$

# Why is this easier?

---

## Multivariate case

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

$$p(\mathbf{y}) = \sum_i p(\mathbf{y}|\mathbf{x}^{(i)})p(\mathbf{x}^{(i)})$$

## 1D case

$$p(x_i|\mathbf{x}_{\neg i}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x}_{\neg i}, \mathbf{y})}$$

$$p(\mathbf{x}_{\neg i}, \mathbf{y}) = \int p(\mathbf{x}, \mathbf{y}) dx_i = \sum_{x_i \in [1, -1]} p(x_i, \mathbf{x}_{\neg i}, \mathbf{y})$$

$$= p(x_i = 1, \mathbf{x}_{\neg i}, \mathbf{y}) + p(x_i = -1, \mathbf{x}_{\neg i}, \mathbf{y})$$

## Summary

---

- There are lots and lots of sampling methods

## Summary

---

- There are lots and lots of sampling methods
- Most provide guarantees in the limit

## Summary

---

- There are lots and lots of sampling methods
- Most provide guarantees in the limit
- Often requires a lot of problem knowledge to make efficient

## Summary

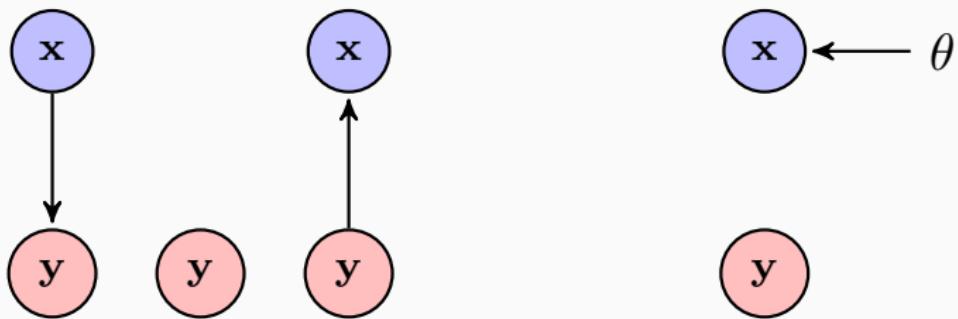
---

- There are lots and lots of sampling methods
- Most provide guarantees in the limit
- Often requires a lot of problem knowledge to make efficient
- Hard to get a measure of how well you are doing

## Deterministic Inference

---

## Lower Bound



$$p(y) = \int_x p(y|x)p(x) = \frac{p(y|x)p(x)}{p(x|y)}$$

$$q_{\theta}(x) \approx p(x|y)$$

# Variational Bayes

---

$$p(y)$$

# Variational Bayes

---

$$\log p(y)$$

# Variational Bayes

---

$$\log p(y) = \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx$$

# Variational Bayes

---

$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\ &= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx\end{aligned}$$

# Variational Bayes

---

$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\ &= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\ &= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx\end{aligned}$$

# Variational Bayes

---

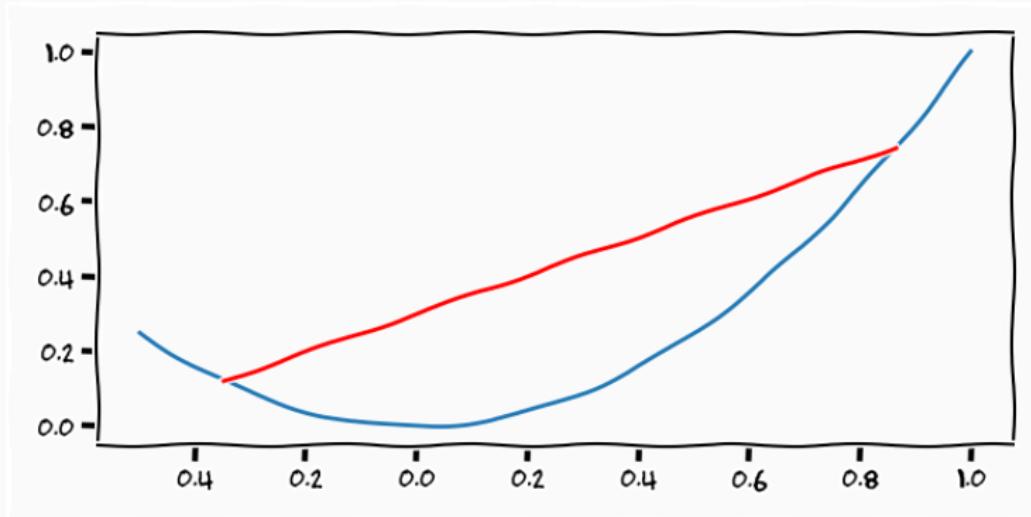
$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\&= \int q(x) \log \frac{q(x)}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{1}{p(x|y)} dx\end{aligned}$$

## Variational Bayes

---

$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\&= \int q(x) \log \frac{q(x)}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{1}{p(x|y)} dx \\&= \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{q(x)}{p(x|y)} dx\end{aligned}$$

# Jensen Inequality



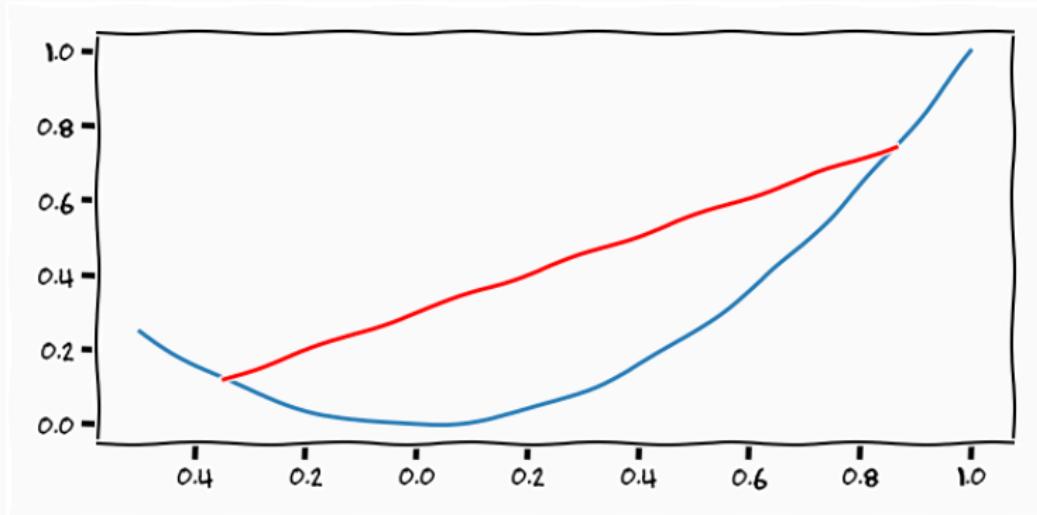
## Convex Function

$$\lambda f(x_0) + (1 - \lambda)f(x_1) \geq f(\lambda x_0 + (1 - \lambda)x_1)$$

$$x \in [x_{min}, x_{max}]$$

$$\lambda \in [0, 1]$$

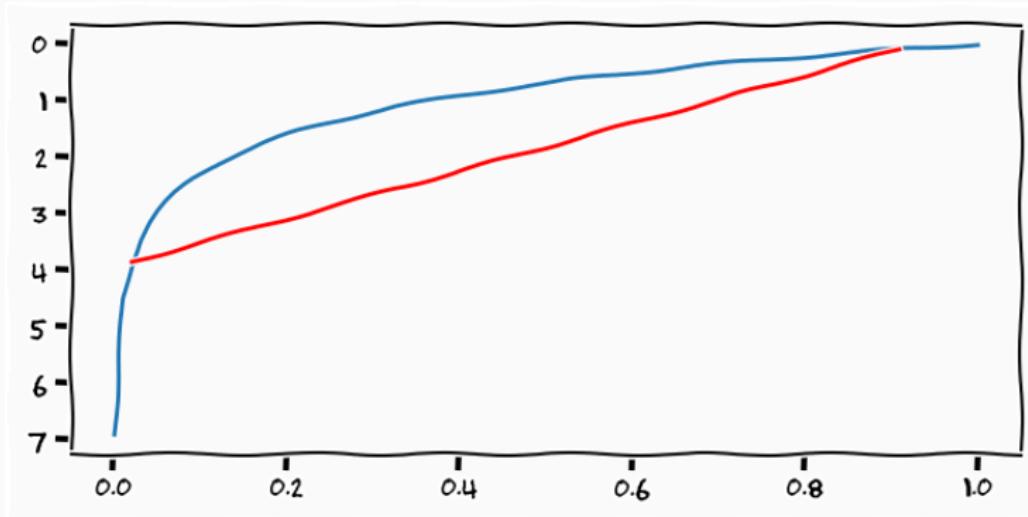
## Jensen Inequality



$$\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$$

$$\int f(x)p(x)dx \geq f\left(\int xp(x)dx\right)$$

## Jensen Inequality in Variational Bayes



$$\int \log(x)p(x)dx \leq \log \left( \int xp(x)dx \right)$$

*moving the log inside the integral is a lower-bound on the integral*

## The "posterior" term

---

$$\int q(x) \log \frac{q(x)}{p(x|y)} dx$$

## The "posterior" term

---

$$\int q(x) \log \frac{q(x)}{p(x|y)} dx$$

$$= - \int q(x) \log \frac{p(x|y)}{q(x)} dx$$

## The "posterior" term

---

$$\begin{aligned} & \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= - \int q(x) \log \frac{p(x|y)}{q(x)} dx \\ &\geq -\log \int p(x|y) dx = -\log 1 = 0 \end{aligned}$$

## The "posterior" term

---

$$\int q(x) \log \frac{q(x)}{p(x|y)} dx$$

## The "posterior" term

---

$$\int q(x) \log \frac{q(x)}{p(x|y)} dx$$

= {Lets assume that  $q(x) = p(x|y)$ }

## The "posterior" term

---

$$\begin{aligned} & \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= \{\text{Let's assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx \end{aligned}$$

## The "posterior" term

---

$$\begin{aligned} & \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= \{\text{Lets assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx \\ &= 0 \end{aligned}$$

# Kullback-Leibler Divergence

---

$$KL(q(x)||p(x|y)) = \int q(x) \log \frac{q(x)}{p(x|y)} dx$$

- Measure of divergence between distributions
- Not a metric (not symmetric)
- $KL(q(x)||p(x|y)) = 0 \Leftrightarrow q(x) = p(x|y)$
- $KL(q(x)||p(x|y)) \geq 0$

## The "other terms"

---

$$\int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx =$$

## The "other terms"

---

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \end{aligned}$$

## The "other terms"

---

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{\text{Let's assume that } q(x) = p(x|y)\} \end{aligned}$$

## The "other terms"

---

$$\begin{aligned}& \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\&= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\&= \{\text{Let's assume that } q(x) = p(x|y)\} \\&= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx\end{aligned}$$

## The "other terms"

---

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{\text{Let's assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \end{aligned}$$

## The "other terms"

---

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{\text{Let's assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx + \int p(x|y) \log p(y) dx \end{aligned}$$

## The "other terms"

---

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{\text{Let's assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx + \int p(x|y) \log p(y) dx \\ &= \underbrace{\int p(x|y) dx}_{=1} \log p(y) \end{aligned}$$

## The "other terms"

---

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{\text{Let's assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx + \int p(x|y) \log p(y) dx \\ &= \underbrace{\int p(x|y) dx}_{=1} \log p(y) = \log p(y) \end{aligned}$$

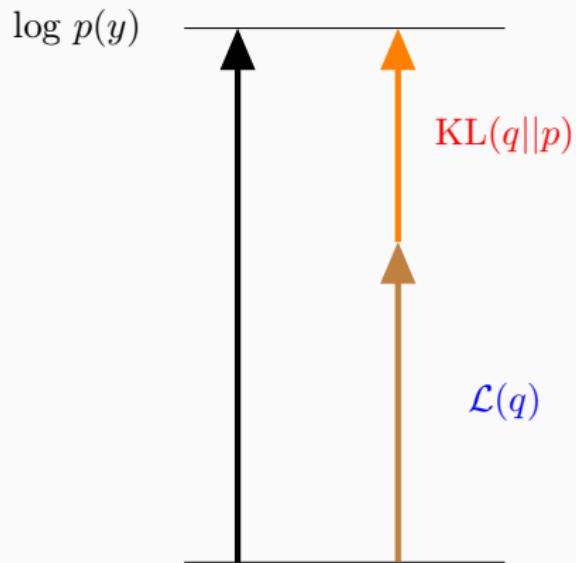
# Variational Bayes

---

$$\begin{aligned}\log p(y) &= \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &\geq - \int q(x) \log q(x) dx + \int q(x) \log p(x, y) dx\end{aligned}$$

- The Evidence Lower BOnd
- Tight if  $q(x) = p(x|y)$

# Deterministic Approximation



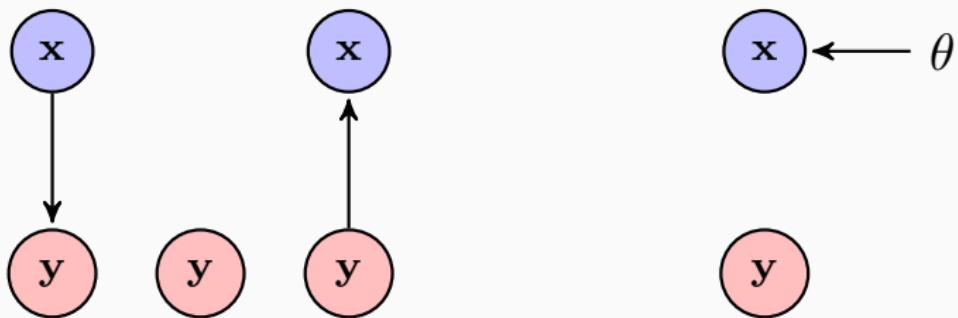
# ELBO

---

$$\begin{aligned}\log p(y) &\geq - \int q(x) \log q(x) dx + \int q(x) \log p(x, y) dx \\ &= \mathbb{E}_{q(x)} [\log p(x, y)] - H(q(x)) = \mathcal{L}(q(x))\end{aligned}$$

- if we maximise the ELBO we,
  - find an approximate posterior
  - lower bound the marginal likelihood
- *maximising  $p(y)$*  is learning
- finding  $q(x) \approx p(x|y)$  is prediction

## Lower Bound



$$p(y) = \int_x p(y|x)p(x) = \frac{p(y|x)p(x)}{p(x|y)}$$

$$q_{\theta}(x) \approx p(x|y)$$

## Why is this useful?

---

Why is this a sensible thing to do?

– Ryan Adams<sup>5</sup>

---

<sup>5</sup>Talking Machines Podcast

## Why is this useful?

---

### Why is this a sensible thing to do?

- If we can't formulate the joint distribution there isn't much we can do

– Ryan Adams<sup>5</sup>

---

<sup>5</sup>Talking Machines Podcast

## Why is this useful?

---

### Why is this a sensible thing to do?

- If we can't formulate the joint distribution there isn't much we can do
- Taking the expectation of a log is usually easier than the expectation

– Ryan Adams<sup>5</sup>

---

<sup>5</sup>Talking Machines Podcast

## Why is this useful?

---

### Why is this a sensible thing to do?

- If we can't formulate the joint distribution there isn't much we can do
- Taking the expectation of a log is usually easier than the expectation
- We are allowed to choose the distribution to take the expectation over
  - Ryan Adams<sup>5</sup>

---

<sup>5</sup>Talking Machines Podcast

## How to choose Q?

---

$$\mathcal{L}(q(x)) = \mathbb{E}_{q(x)} [\log p(x, y)] - H(q(x))$$

- We have to be able to compute an expectation over the joint distribution
- The second term should be trivial

# Mean Field Approximation

---

$$q(\mathbf{X}) = \prod_i q_i(\mathbf{x}_i)$$

$$\mathcal{L}(q_j) = \mathcal{L}_j(q_j) + \mathcal{L}_{\neg j}(q_{\neg j}),$$

- Model originating in Physics
- We model marginals rather than the full distribution
- We can update each distribution in turn and cycle

## Summary

---

- Variational Bayes is **never** exact
- It can be made very fast
- There are exciting works using *implicit* variational distributions

## Summary

---

# Summary

---

- The evidence allows you to do model selection
  - "simple" is a relative concept that needs to be defined
- Most interesting things are intractable
- This leads to the field of approximative inference
  - you have gotten a "hint" of what is out there and some scaffolding to build upon

## But what...

---

- *Now how about if we would treat the whole thing as a machine learning/inference problem instead?*

## But what...

---

- *Now how about if we would treat the whole thing as a machine learning/inference problem instead?*
- you can sample the integrand → data

## But what...

---

- *Now how about if we would treat the whole thing as a machine learning/inference problem instead?*
- you can sample the integrand → data
- you kinda have an idea of the integral → prior

## But what...

---

- Now how about if we would treat the whole thing as a machine learning/inference problem instead?
- you can sample the integrand → data
- you kinda have an idea of the integral → prior
- you would be able to match the integral to data → likelihood

## But what... ---

- *Now how about if we would treat the whole thing as a machine learning/inference problem instead?*
- you can sample the integrand → data
- you kinda have an idea of the integral → prior
- you would be able to match the integral to data → likelihood
- and this we might talk about later in the course :-)

eof

## References

---

- 
- Lawrence, Neil D (2005). "Probabilistic non-linear principal component analysis with Gaussian process latent variable models". In: *Journal of Machine Learning Research* 6, pp. 1783–1816.

# Mean Field Approximation

$$\mathcal{L}(q) = \int q(\mathbf{X}) \log \frac{p(\mathbf{Y}, \mathbf{X})}{q(\mathbf{X})} d\mathbf{X}$$

# Mean Field Approximation

$$\mathcal{L}(q) = \int q(\mathbf{X}) \log \frac{p(\mathbf{Y}, \mathbf{X})}{q(\mathbf{X})} d\mathbf{X} = \int \prod_i q_i(\mathbf{x}_i) \log \frac{p(\mathbf{Y}, \mathbf{X})}{\prod_k q_k(\mathbf{x}_k)} d\mathbf{X}$$

# Mean Field Approximation

$$\begin{aligned}\mathcal{L}(q) &= \int q(\mathbf{X}) \log \frac{p(\mathbf{Y}, \mathbf{X})}{q(\mathbf{X})} d\mathbf{X} = \int \prod_i q_i(\mathbf{x}_i) \log \frac{p(\mathbf{Y}, \mathbf{X})}{\prod_k q_k(\mathbf{x}_k)} d\mathbf{X} \\ &= \int \prod_i q_i(\mathbf{x}_i) \left( \log p(\mathbf{Y}, \mathbf{X}) - \sum_k \log q_k(\mathbf{x}_k) \right)\end{aligned}$$

# Mean Field Approximation

$$\mathcal{L}(q_j) = \int \prod_i q_i(\mathbf{x}_i) \left( \log p(\mathbf{Y}, \mathbf{X}) - \sum_k \log q_k(\mathbf{x}_k) \right) d\mathbf{x}$$

# Mean Field Approximation

$$\begin{aligned}\mathcal{L}(q_j) &= \int \prod_i q_i(\mathbf{x}_i) \left( \log p(\mathbf{Y}, \mathbf{X}) - \sum_k \log q_k(\mathbf{x}_k) \right) d\mathbf{x} \\ &= \int_j \int_{\neg j} q_j(\mathbf{x}_j) \prod_{i \neq j} q_i(\mathbf{x}_i) \left( \log p(\mathbf{X}, \mathbf{Y}) - \sum_k \log q_k(\mathbf{x}_k) \right) d\mathbf{x}_{\neg j} d\mathbf{x}_j\end{aligned}$$

# Mean Field Approximation

$$\begin{aligned}\mathcal{L}(q_j) &= \int \prod_i q_i(\mathbf{x}_i) \left( \log p(\mathbf{Y}, \mathbf{X}) - \sum_k \log q_k(\mathbf{x}_k) \right) d\mathbf{x} \\ &= \int_j \int_{\neg j} q_j(\mathbf{x}_j) \prod_{i \neq j} q_i(\mathbf{x}_i) \left( \log p(\mathbf{X}, \mathbf{Y}) - \sum_k \log q_k(\mathbf{x}_k) \right) d\mathbf{x}_{\neg j} d\mathbf{x}_j \\ &= \int_j q_j(\mathbf{x}_j) \underbrace{\int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) \log p(\mathbf{Y}, \mathbf{X}) d\mathbf{x}_{\neg j}}_{\log f_j(\mathbf{x}_j)} d\mathbf{x}_j \\ &\quad - \int_j q_j(\mathbf{x}_j) \int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) \left( \log q_j(\mathbf{x}_j) + \sum_{k \neq j} \log q_k(\mathbf{x}_k) \right) d\mathbf{x}_{\neg j} d\mathbf{x}_j\end{aligned}$$

# Mean Field Approximation

$$\begin{aligned} &= \int_j q_j(\mathbf{x}_j) \underbrace{\int_{\neg j} \prod_{i \neq j} q(\mathbf{x}_i) \log p(\mathbf{Y}, \mathbf{X}) d\mathbf{x}_{\neg j} d\mathbf{x}_j}_{\log f_j(\mathbf{x}_j)} \\ &- \int_j q_j(\mathbf{x}_j) \int_{\neg j} \prod_{i \neq j} q(\mathbf{x}_i) \left( \log q_j(\mathbf{x}_j) + \sum_{k \neq j} \log q_k(\mathbf{x}_k) \right) d\mathbf{x}_{\neg j} d\mathbf{x}_j \end{aligned}$$

# Mean Field Approximation

$$\begin{aligned} &= \int_j q_j(\mathbf{x}_j) \underbrace{\int_{\neg j} \prod_{i \neq j} q(\mathbf{x}_i) \log p(\mathbf{Y}, \mathbf{X}) d\mathbf{x}_{\neg j}}_{\log f_j(\mathbf{x}_j)} d\mathbf{x}_j \\ &\quad - \int_j q_j(\mathbf{x}_j) \int_{\neg j} \prod_{i \neq j} q(\mathbf{x}_i) \left( \log q_j(\mathbf{x}_j) + \sum_{k \neq j} \log q_k(\mathbf{x}_k) \right) d\mathbf{x}_{\neg j} d\mathbf{x}_j \\ &= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j \\ &\quad - \int_j q_j(\mathbf{x}_j) \left( \underbrace{\log q_j(\mathbf{x}_j) \int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) d\mathbf{x}_{\neg j}}_{=1} + \underbrace{\int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) \sum_{k \neq j} \log q_k(\mathbf{x}_k) d\mathbf{x}_{\neg j}}_{\text{constant w.r.t. } q_j} \right) \end{aligned}$$

# Mean Field Approximation

$$= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j$$
$$- \int_j q_j(\mathbf{x}_j) \left( \underbrace{\log q_j(\mathbf{x}_j) \int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) d\mathbf{x}_{\neg j}}_{=1} + \underbrace{\int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) \sum_{k \neq j} \log q_k(\mathbf{x}_k) d\mathbf{x}_{\neg j}}_{\text{constant w.r.t. } q_j} \right) d\mathbf{x}_j$$

# Mean Field Approximation

$$\begin{aligned} &= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j \\ &- \int_j q_j(\mathbf{x}_j) \left( \underbrace{\log q_j(\mathbf{x}_j) \int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) d\mathbf{x}_{\neg j}}_{=1} + \underbrace{\int_{\neg j} \prod_{i \neq j} q_i(\mathbf{x}_i) \sum_{k \neq j} \log q_k(\mathbf{x}_k) d\mathbf{x}_{\neg j}}_{\text{constant w.r.t. } q_j} \right) d\mathbf{x}_j \\ &= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j - \int_j q_j(\mathbf{x}_j) \log q_j(\mathbf{x}_j) d\mathbf{x}_j + \text{const.} \cdot \underbrace{\int_j q_j(\mathbf{x}_j) d\mathbf{x}_j}_{=1} \end{aligned}$$

# Mean Field Approximation

$$= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j - \int_j q_j(\mathbf{x}_j) \log q_j(\mathbf{x}_j) d\mathbf{x}_j + \text{const.} \cdot \underbrace{\int_j q_j(\mathbf{x}_j) d\mathbf{x}_j}_{=1}$$

# Mean Field Approximation

$$= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j - \int_j q_j(\mathbf{x}_j) \log q_j(\mathbf{x}_j) d\mathbf{x}_j + \text{const.} \cdot \underbrace{\int_j q_j(\mathbf{x}_j) d\mathbf{x}_j}_{=1}$$

$$= \int_j q_j(\mathbf{x}_j) \log \frac{f_j(\mathbf{x}_j)}{q_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.}$$

# Mean Field Approximation

$$= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j - \int_j q_j(\mathbf{x}_j) \log q_j(\mathbf{x}_j) d\mathbf{x}_j + \text{const.} \cdot \underbrace{\int_j q_j(\mathbf{x}_j) d\mathbf{x}_j}_{=1}$$

$$\begin{aligned} &= \int_j q_j(\mathbf{x}_j) \log \frac{f_j(\mathbf{x}_j)}{q_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.} \\ &= - \int_j q_j(\mathbf{x}_j) \log \frac{q_j(\mathbf{x}_j)}{f_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.} \end{aligned}$$

# Mean Field Approximation

$$= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j - \int_j q_j(\mathbf{x}_j) \log q_j(\mathbf{x}_j) d\mathbf{x}_j + \text{const.} \cdot \underbrace{\int_j q_j(\mathbf{x}_j) d\mathbf{x}_j}_{=1}$$

$$\begin{aligned} &= \int_j q_j(\mathbf{x}_j) \log \frac{f_j(\mathbf{x}_j)}{q_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.} \\ &= - \int_j q_j(\mathbf{x}_j) \log \frac{q_j(\mathbf{x}_j)}{f_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.} \\ &= -\text{KL}(q_j(\mathbf{x}_j) || f_j(\mathbf{x}_j)) + \text{const.} \end{aligned}$$

# Mean Field Approximation

$$= \int_j q_j(\mathbf{x}_j) \log f_j(\mathbf{x}_j) d\mathbf{x}_j - \int_j q_j(\mathbf{x}_j) \log q_j(\mathbf{x}_j) d\mathbf{x}_j + \text{const.} \cdot \underbrace{\int_j q_j(\mathbf{x}_j) d\mathbf{x}_j}_{=1}$$

$$= \int_j q_j(\mathbf{x}_j) \log \frac{f_j(\mathbf{x}_j)}{q_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.}$$

$$= - \int_j q_j(\mathbf{x}_j) \log \frac{q_j(\mathbf{x}_j)}{f_j(\mathbf{x}_j)} d\mathbf{x}_j + \text{const.}$$

$$= -\text{KL}(q_j(\mathbf{x}_j) || f_j(\mathbf{x}_j)) + \text{const.} = \mathcal{L}(q_j)$$

# Mean Field Approximation

$$\mathcal{L}(q_j) = -\text{KL}(q_j(\mathbf{x}_j) || f_j(\mathbf{x}_j)) + \text{const.}$$

- Want to maximise lower bound
- Negative KL  $\rightarrow$  minimise KL term
- *we are free to choose the form of the distribution*

## Mean Field Approximation

$$\begin{aligned}\log q_j(\mathbf{x}_j) &= \log f_j(\mathbf{x}_j) = \int_{\neg j} \underbrace{\prod_{i \neq j} q_i(\mathbf{x}_i)}_{q_{\neg j}(\mathbf{x}_{\neg j})} \log p(\mathbf{Y}, \mathbf{X}) d\mathbf{x}_{\neg j} \\ &= \mathbb{E}_{q_{\neg j}(\mathbf{x}_{\neg j})} [\log p(\mathbf{Y}, \mathbf{X})]\end{aligned}$$

- Choose the marginal distribution that makes the bound tight
- Will not make the bound tight in general though

# Mean Field Variational Bayes

1. Formulate joint distribution over data and latent parameters

# Mean Field Variational Bayes

1. Formulate joint distribution over data and latent parameters
2. Formulate fully factorised approximative posterior over latent variables

# Mean Field Variational Bayes

1. Formulate joint distribution over data and latent parameters
2. Formulate fully factorised approximative posterior over latent variables
3. Fit marginal approximation by making bound tight

# Mean Field Variational Bayes

1. Formulate joint distribution over data and latent parameters
2. Formulate fully factorised approximative posterior over latent variables
3. Fit marginal approximation by making bound tight
4. Iterate through variables

## Taking stock

- We have derived a general update for the mean-field approximation

## Taking stock

- We have derived a general update for the mean-field approximation
- Looks an awful lot like Expectation-Maximisation
  - EM is VB (E-local, M-global)

## Taking stock

- We have derived a general update for the mean-field approximation
- Looks an awful lot like Expectation-Maximisation
  - EM is VB (E-local, M-global)
- Looks an awful lot like Gibbs sampling
  - Gibbs: sample from complete conditional
  - VB: variational expectations of natural parameters of the complete conditional