

# Unsupervised and Composite Gaussian Processes

---

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

September 15, 2020

<http://carlhenrik.com>

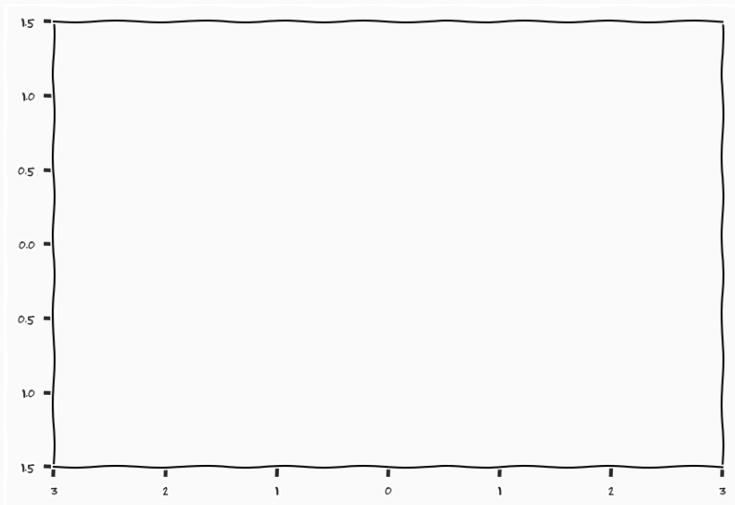
- $\mathcal{F}$  space of functions
- $\mathcal{A}$  learning algorithm
- $\mathcal{S} = \{(x_1, y_1), \dots, (x_N, y_N)\}$
- $\mathcal{S} \sim P(\mathcal{X} \times \mathcal{Y})$
- $\ell(\mathcal{A}_{\mathcal{F}}(\mathcal{S}), x, y)$  loss function

$$e(\mathcal{S}, \mathcal{A}, \mathcal{F}) = \mathbb{E}_{P(\{\mathcal{X}, \mathcal{Y}\})} [\ell(\mathcal{A}_{\mathcal{F}}(\mathcal{S}), x, y)]$$

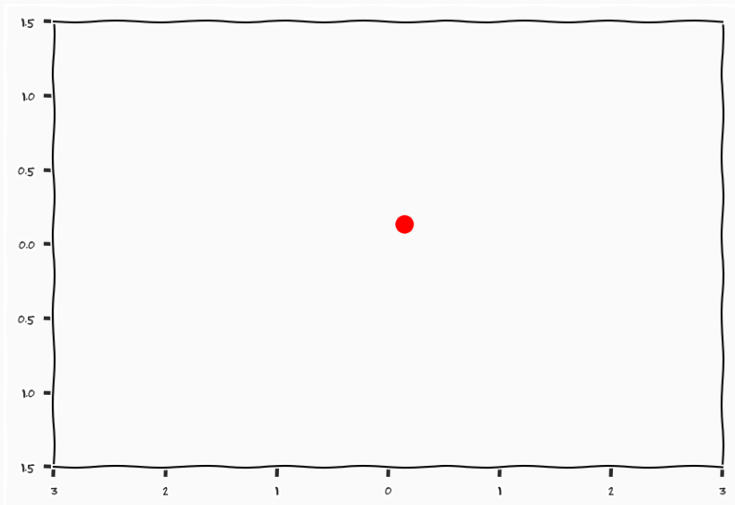
$$\begin{aligned} e(\mathcal{S}, \mathcal{A}, \mathcal{F}) &= \mathbb{E}_{P(\{\mathcal{X}, \mathcal{Y}\})} [\ell(\mathcal{A}_{\mathcal{F}}(\mathcal{S}), x, y)] \\ &\approx \frac{1}{M} \sum_{n=1}^M \ell(\mathcal{A}_{\mathcal{F}}(\mathcal{S}), x_n, y_n) \end{aligned}$$

We can come up with a combination of  $\{\mathcal{S}, \mathcal{A}, \mathcal{F}\}$  that makes  $e(\mathcal{S}, \mathcal{A}, \mathcal{F})$  take an arbitrary value

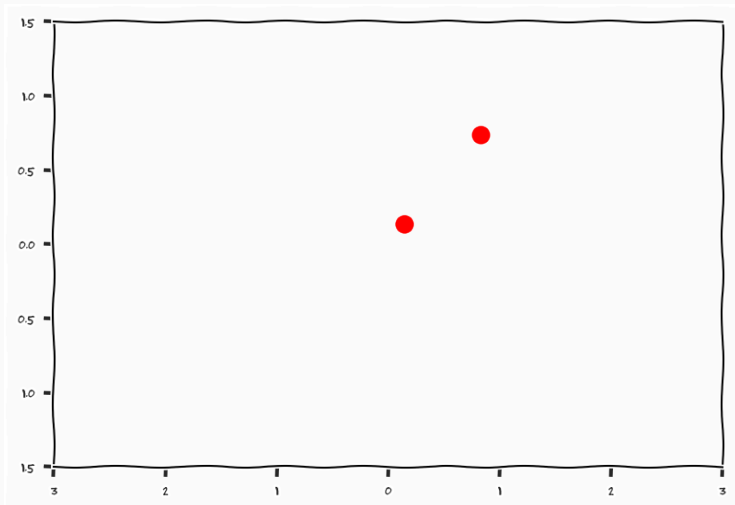
# Example



# Example

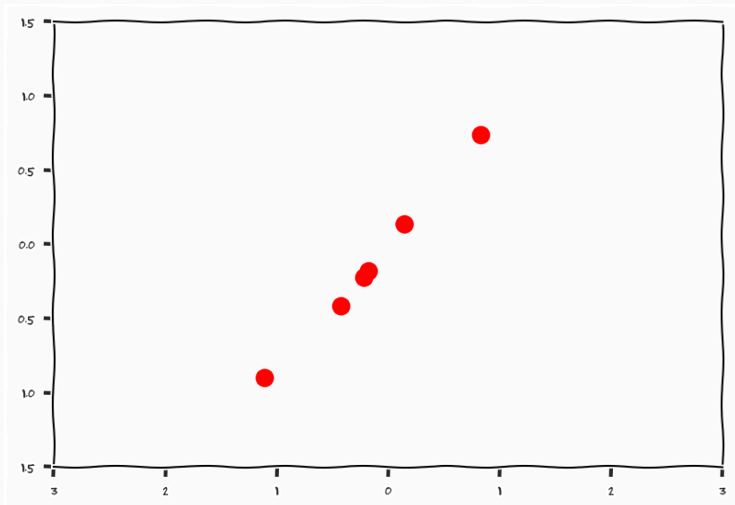


# Example

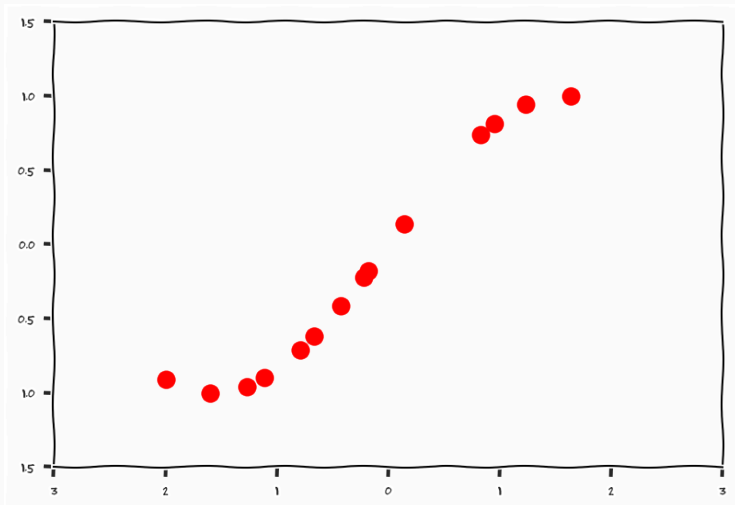




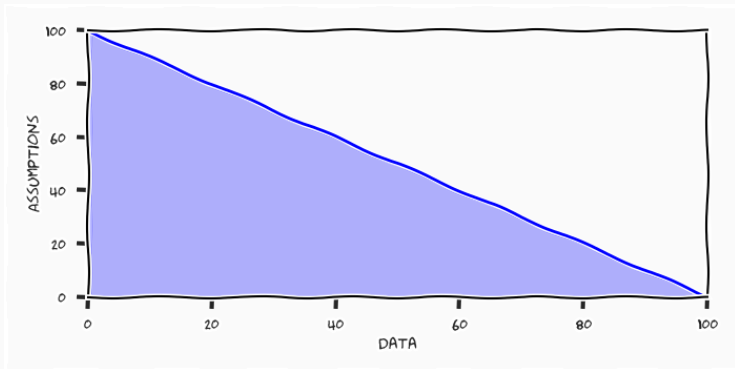
# Example



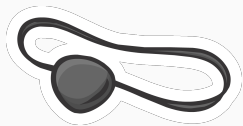
# Example



# Data and Knowledge



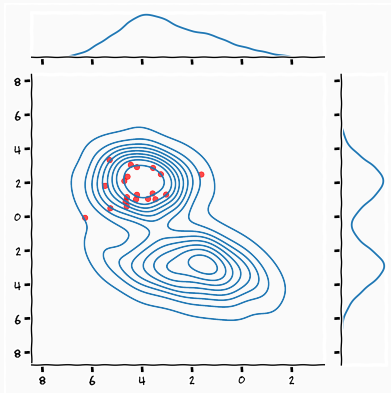
# Assumptions: Algorithms



Statistical Learning

$$A_{\mathcal{F}}(S)$$

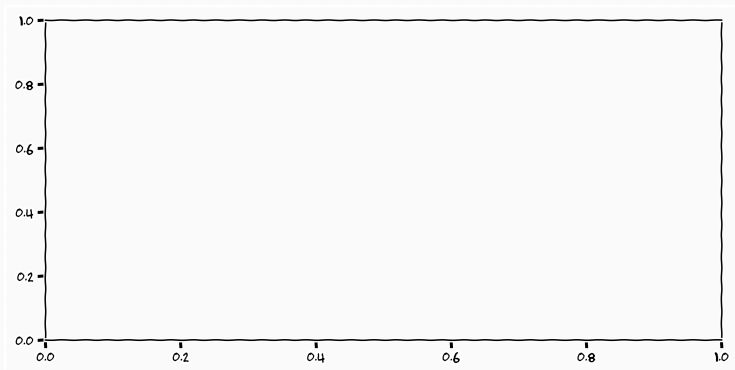
# Assumptions: Biased Sample



Statistical Learning

$$\mathcal{A}_{\mathcal{F}}(\mathcal{S})$$

# Assumptions: Hypothesis space



Statistical Learning

$$\mathcal{A}_{\mathcal{F}}(\mathcal{S})$$

- There seems to be a narrative that the more *flexible* a model is the better it is

# The No Free Lunch

- There seems to be a narrative that the more *flexible* a model is the better it is
  - This is not true



# The No Free Lunch

- There seems to be a narrative that the more *flexible* a model is the better it is
  - This is not true
- The best possible model has infinite support (nothing is excluded) but very focused mass

# The No Free Lunch

- There seems to be a narrative that the more *flexible* a model is the better it is
  - This is not true
- The best possible model has infinite support (nothing is excluded) but very focused mass
- *Your solution can only ever be interpreted in the light of your assumptions*



*Iudicium Posterium Discipulus Est Prioris*<sup>1</sup>

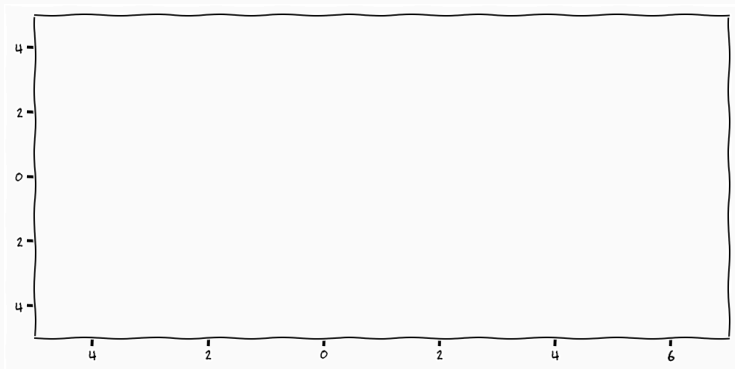
---

<sup>1</sup>The posterior is the student of the prior

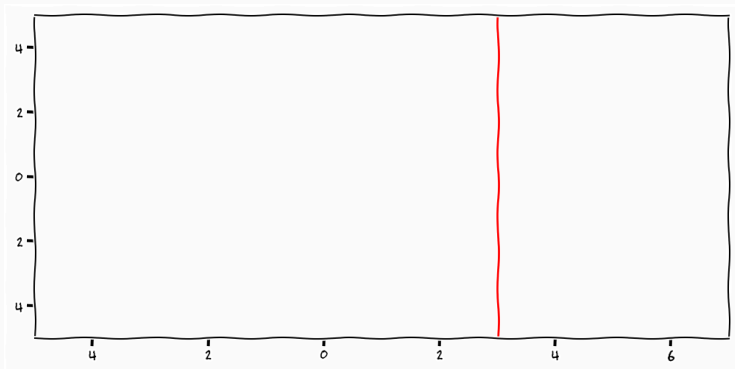
# Gaussian Processes

---

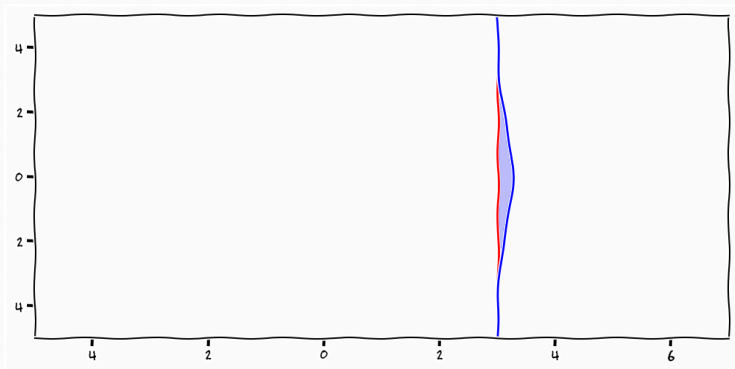
# Gaussian Processes



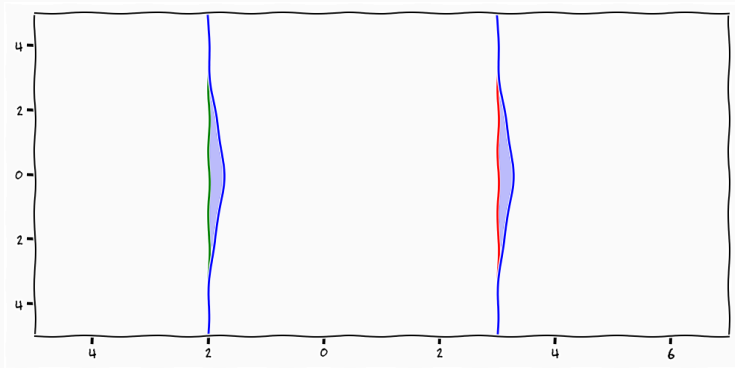
# Gaussian Processes



# Gaussian Processes

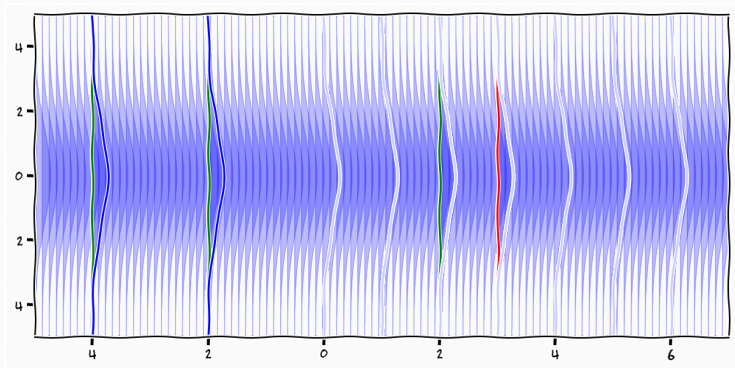


# Gaussian Processes

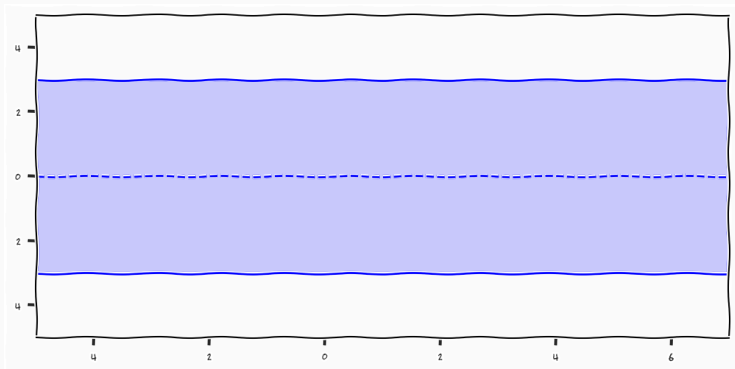




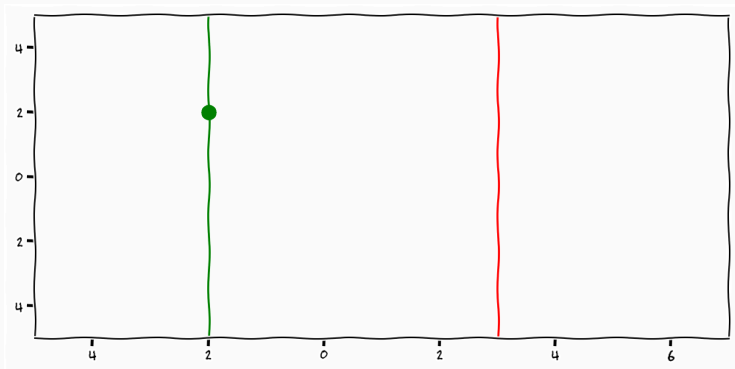
# Gaussian Processes



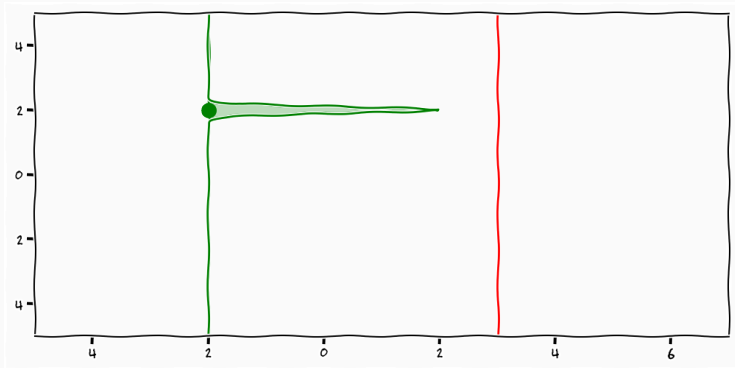
# Gaussian Processes



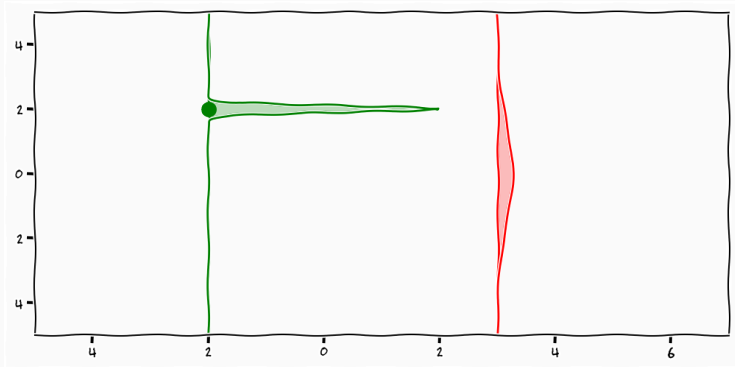
# Gaussian Processes



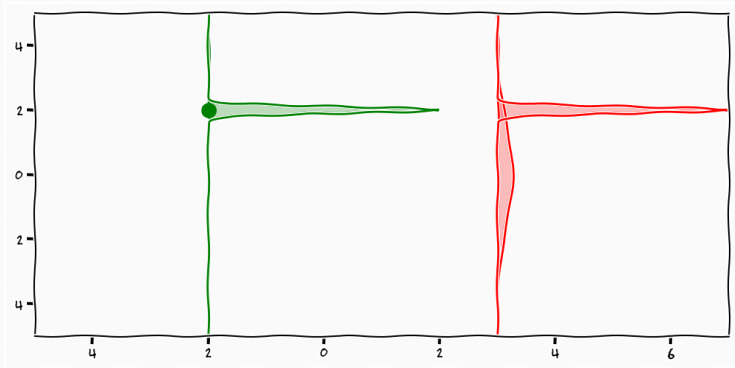
# Gaussian Processes



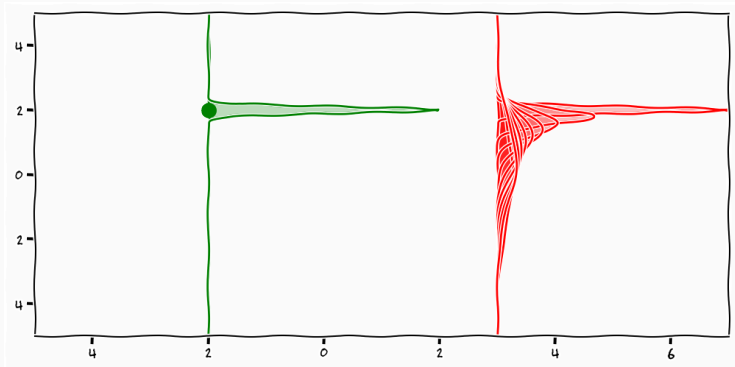
# Gaussian Processes



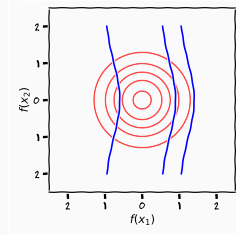
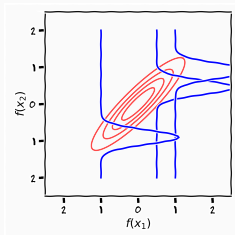
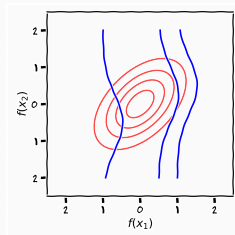
# Gaussian Processes



# Gaussian Processes



# Conditional Gaussians



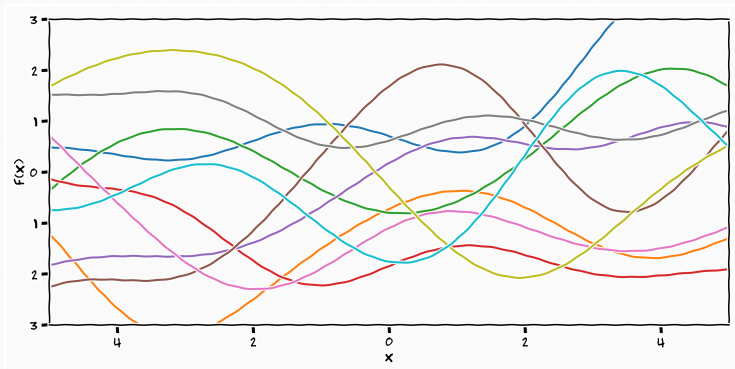
$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$$

$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}\right)$$

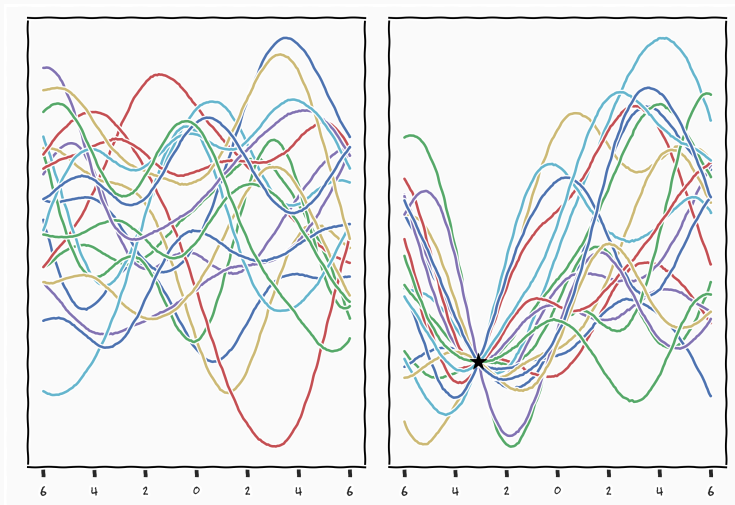
$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$



# Gaussian Processes



# Gaussian Processes



$$p(x_1, x_2) \quad p(x_1) = \int p(x_1, x_2) dx \quad p(x_1|x_2) = \frac{p(x_1, x_2)}{p(x_2)}$$

Gaussian Identities

# Stochastic Processes

---

# Kolmogorov's Existence Theorem

For all permutations  $\pi$ , measurable sets  $F_i \subseteq \mathbb{R}^n$  and probability measure  $\nu$

## 1. Exchangeable

$$\nu_{t_{\pi(1)} \dots t_{\pi(k)}} (F_{\pi(1)} \times \dots \times F_{\pi(k)}) = \nu_{t_1 \dots t_k} (F_1 \times \dots \times F_k)$$

## 2. Marginal

$$\nu_{t_1 \dots t_k} (F_1 \times \dots \times F_k) = \nu_{t_1 \dots t_k, t_{k+1} \dots t_{k+m}} (F_1 \times \dots \times F_k \times \mathbb{R}^n \times \dots \times \mathbb{R}^n)$$

In this case the finite dimensional probability measure is a realisation of an underlying stochastic process

$$p(x_1, x_2) = \mathcal{N} \left( \begin{array}{c|cc} x_1 & \mu_1 & k_{11} & k_{12} \\ x_2 & \mu_2 & k_{21} & k_{22} \end{array} \right)$$

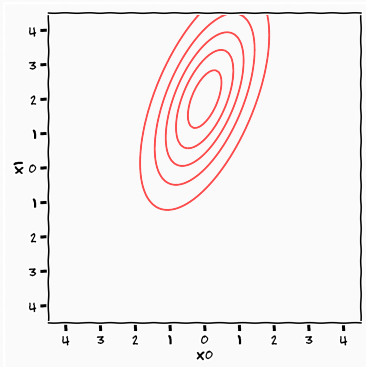
$$\begin{aligned} p(\mathbf{x}_1, \mathbf{x}_2) &= \mathcal{N} \left( \begin{array}{c|cc} \mathbf{x}_1 & \mu_1 & \kappa_{11} & \kappa_{12} \\ \mathbf{x}_2 & \mu_2 & \kappa_{21} & \kappa_{22} \end{array} \right) \\ &= p(\mathbf{x}_2, \mathbf{x}_1) \end{aligned}$$

## Gaussian Distribution - Exchangeable

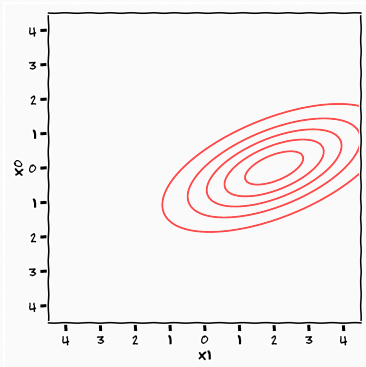
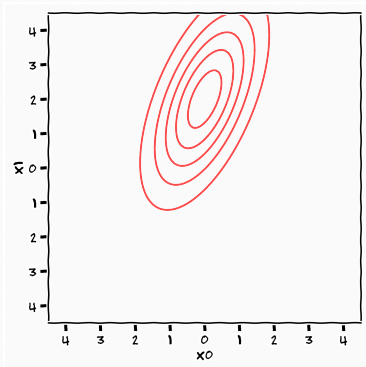
$$\begin{aligned} p(\mathbf{x}_1, \mathbf{x}_2) &= \mathcal{N} \left( \begin{array}{c|cc} \mathbf{x}_1 & \mu_1 & k_{11} & k_{12} \\ \mathbf{x}_2 & \mu_2 & k_{21} & k_{22} \end{array} \right) \\ &= p(\mathbf{x}_2, \mathbf{x}_1) = \mathcal{N} \left( \begin{array}{c|cc} \mathbf{x}_2 & \mu_2 & k_{22} & k_{12} \\ \mathbf{x}_1 & \mu_1 & k_{21} & k_{11} \end{array} \right) \end{aligned}$$



# Gaussian Distribution - Exchangeable



# Gaussian Distribution - Exchangeable



$$p(x_1, x_2) = \mathcal{N} \left( \begin{array}{c|cc} x_1 & \mu_1 & k_{11} & k_{12} \\ x_2 & \mu_2 & k_{21} & k_{22} \end{array} \right)$$

$$p(\mathbf{x}_1, x_2) = \mathcal{N} \left( \begin{array}{c} \mathbf{x}_1 \\ x_2 \end{array} \middle| \begin{array}{cc} \mu_1 & k_{11} & k_{12} \\ \mu_2 & k_{21} & k_{22} \end{array} \right)$$
$$\Rightarrow p(\mathbf{x}_1) = \int_{x_2} p(\mathbf{x}_1, x_2) = \underline{\mathcal{N}(\mathbf{x}_1 \mid \mu_1, k_{11})}$$

## Gaussian Distribution - Marginal

$$p(\mathbf{x}_1, x_2) = \mathcal{N} \left( \begin{array}{c|cc} \mathbf{x}_1 & \mu_1 & k_{11} \\ x_2 & \mu_2 & k_{21} \end{array} \left| \begin{array}{cc} k_{12} & \\ & k_{22} \end{array} \right. \right)$$

$$\Rightarrow p(\mathbf{x}_1) = \int_{x_2} p(\mathbf{x}_1, x_2) = \underline{\mathcal{N}(\mathbf{x}_1 \mid \mu_1, k_{11})}$$

$$p(\mathbf{x}_1, x_2, \dots, x_N) = \mathcal{N} \left( \begin{array}{c|ccc} \mathbf{x}_1 & \mu_1 & k_{11} & k_{12} & \cdots & k_{1N} \\ x_2 & \mu_2 & k_{21} & k_{22} & \cdots & k_{2N} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_N & \mu_N & k_{N1} & k_{N2} & \cdots & k_{NN} \end{array} \right)$$

## Gaussian Distribution - Marginal

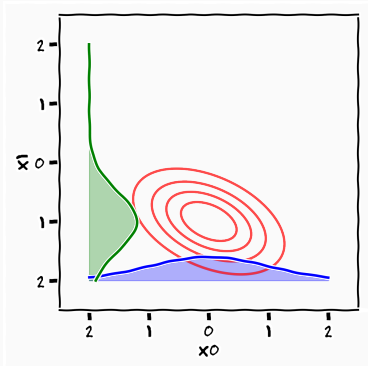
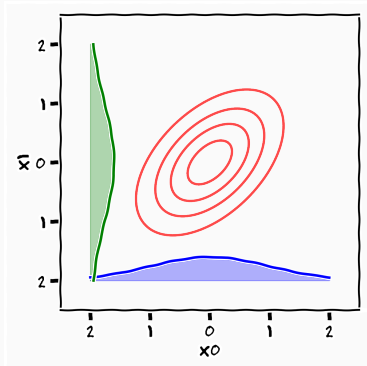
$$p(\mathbf{x}_1, x_2) = \mathcal{N} \left( \begin{array}{c|cc} \mathbf{x}_1 & \mu_1 & k_{11} & k_{12} \\ x_2 & \mu_2 & k_{21} & k_{22} \end{array} \right)$$

$$\Rightarrow p(\mathbf{x}_1) = \int_{x_2} p(\mathbf{x}_1, x_2) = \underline{\mathcal{N}(\mathbf{x}_1 \mid \mu_1, k_{11})}$$

$$p(\mathbf{x}_1, x_2, \dots, x_N) = \mathcal{N} \left( \begin{array}{c|cccc} \mathbf{x}_1 & \mu_1 & k_{11} & k_{12} & \cdots & k_{1N} \\ x_2 & \mu_2 & k_{21} & k_{22} & \cdots & k_{2N} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_N & \mu_N & k_{N1} & k_{N2} & \cdots & k_{NN} \end{array} \right)$$

$$\Rightarrow p(\mathbf{x}_1) = \int_{x_2, \dots, x_N} p(\mathbf{x}_1, x_2, \dots, x_N) = \underline{\mathcal{N}(\mathbf{x}_1 \mid \mu_1, k_{11})}$$

# Gaussian Distribution - Marginal



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

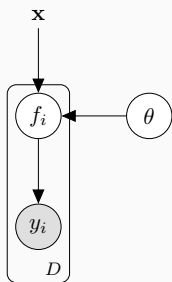
The Gaussian distribution is the projection of the infinite Gaussian process



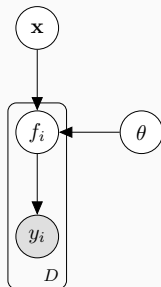
# Unsupervised Gaussian Processes

---

# Unsupervised Learning

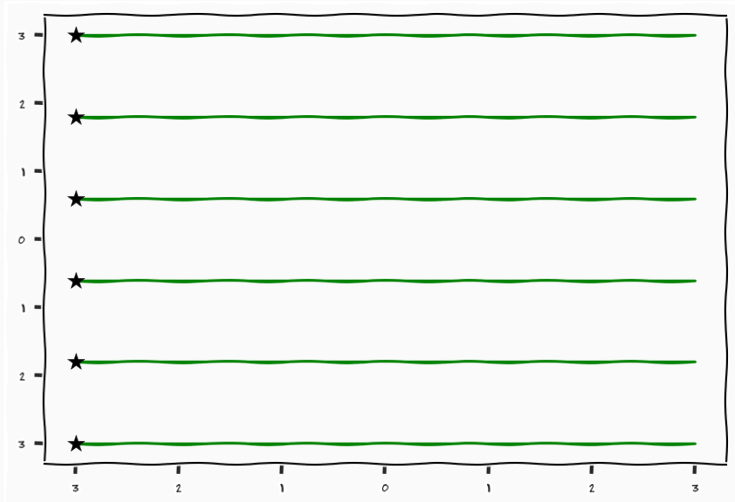


$$p(y|x)$$

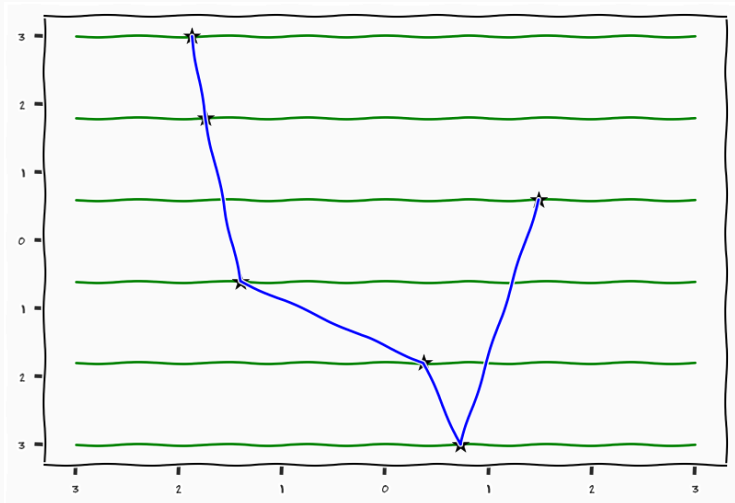


$$p(y)$$

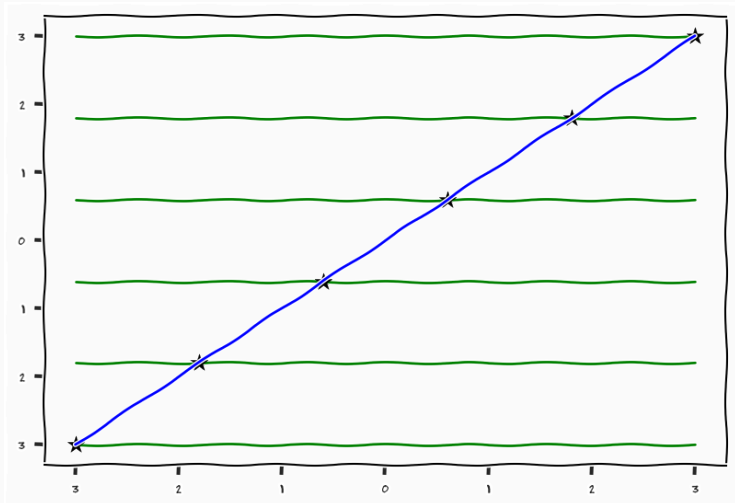
# Unsupervised Learning



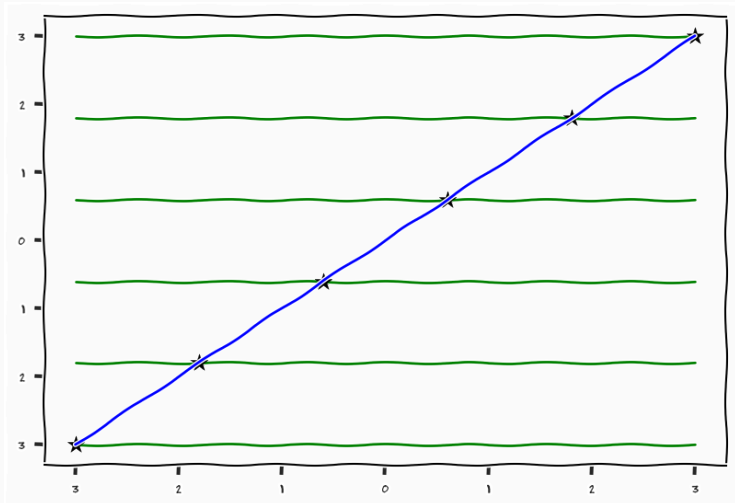
# Unsupervised Learning



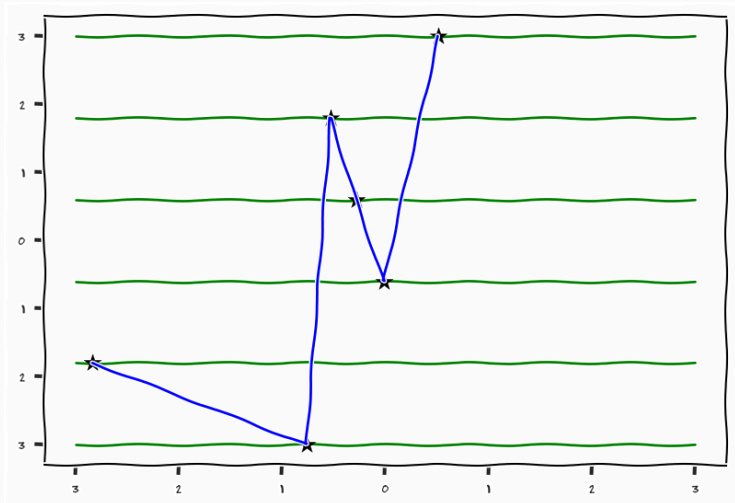
# Unsupervised Learning



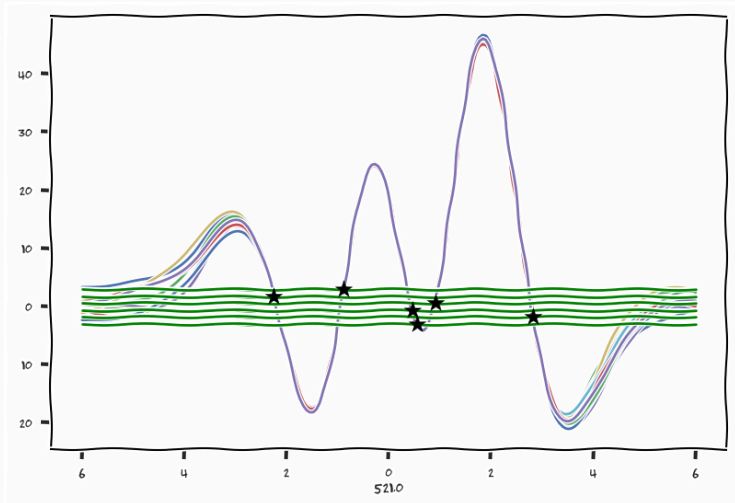
# Unsupervised Learning



# Unsupervised Learning

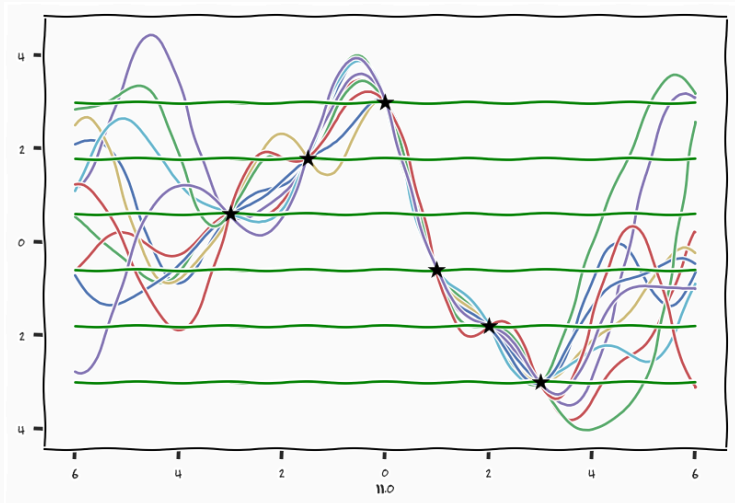


# Unsupervised Learning

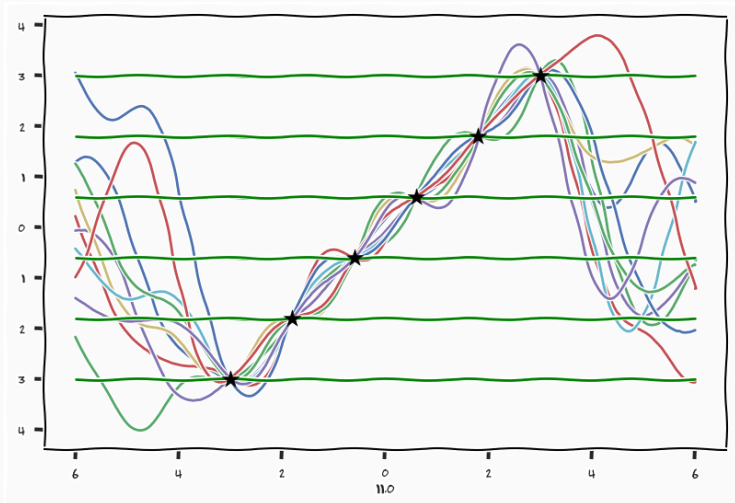




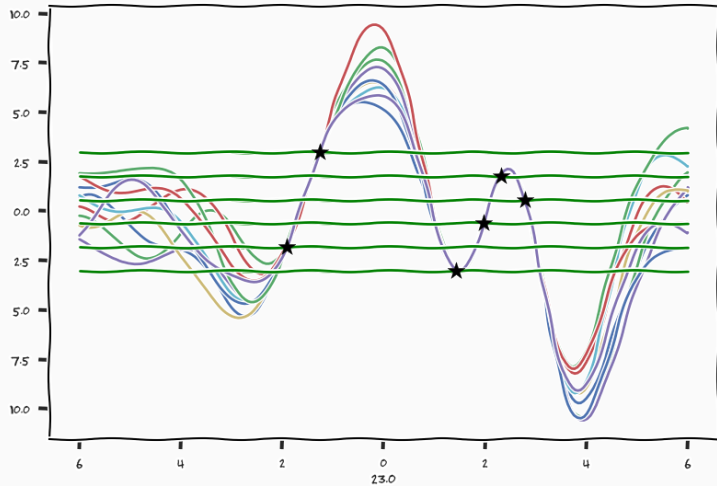
# Unsupervised Learning

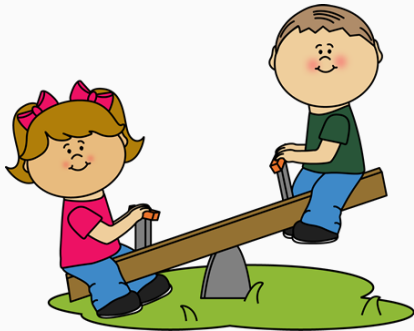


# Unsupervised Learning



# Unsupervised Learning





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

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

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

- GP prior

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

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

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

- 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



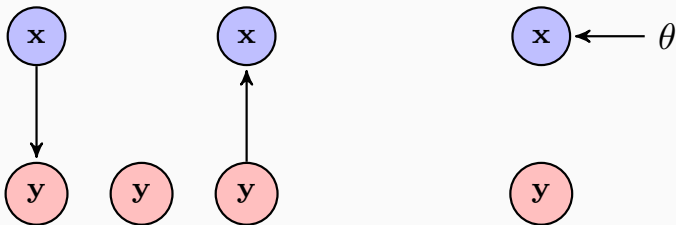
# Approximate Inference

---

$$p(y)$$

- Given some observed data  $y$
- Find a probabilistic model such that the probability of the data is maximised
- Idea: find an approximate model  $q$  that we can integrate

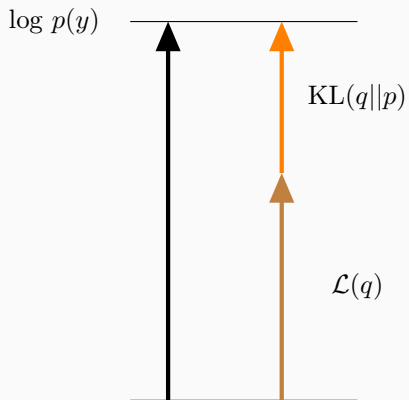
# Lower Bound



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

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

# Deterministic Approximation



$$p(y)$$

$$\log p(y)$$

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

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

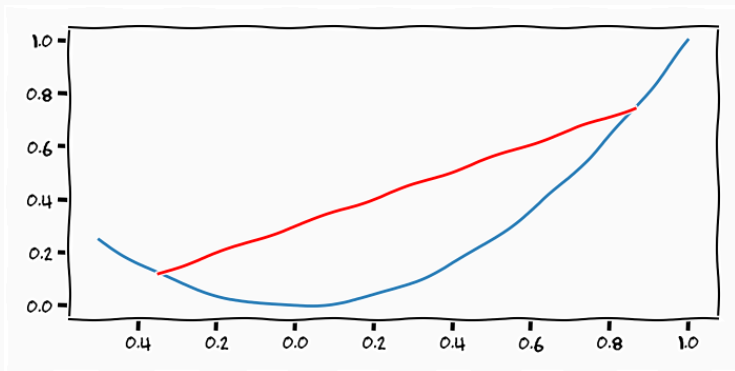


$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\ &= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\ &= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx\end{aligned}$$

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

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

# Jensen Inequality



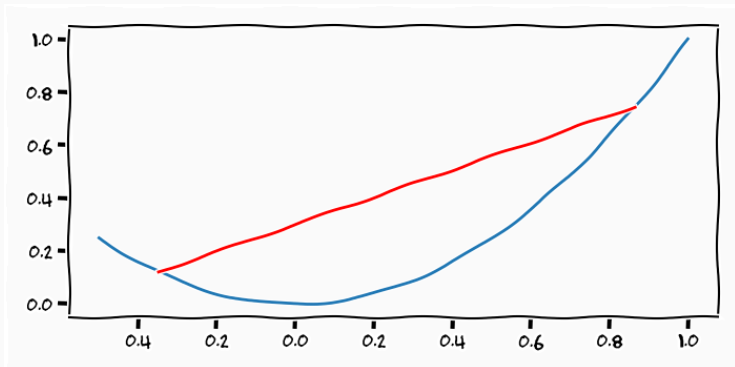
## Convex Function

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

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

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

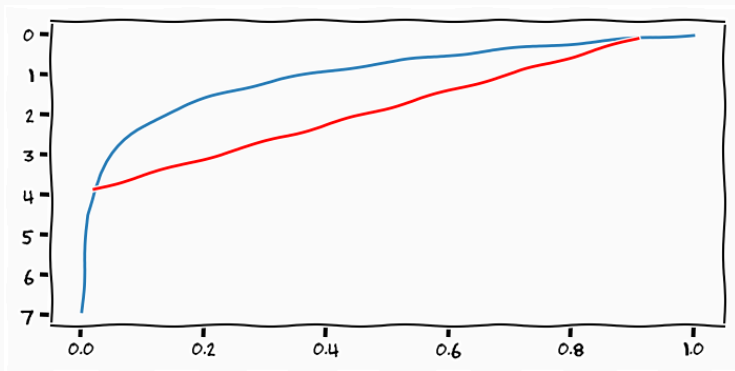
# Jensen Inequality



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

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

## Jensen Inequality in Variational Bayes



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

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

## The "posterior" term

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

## The "posterior" term

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



## The "posterior" term

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

## The "posterior" term

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

## The "posterior" term

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

## The "posterior" term

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

## The "posterior" term

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

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

- 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"

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

## The "other terms"

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

## The "other terms"

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

## The "other terms"

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

## The "other terms"

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

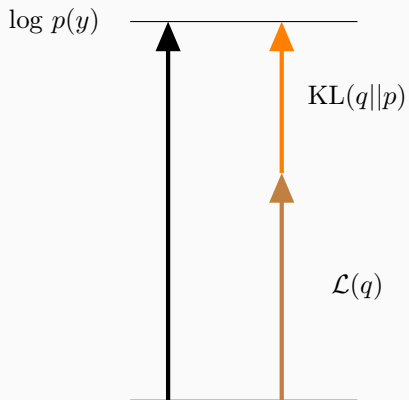
## The "other terms"

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

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

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

# Deterministic Approximation

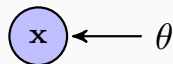




$$\begin{aligned}\log p(y) &\geq - \int q(x) \log q(x) dx + \int q(x) \log p(x, y) dx \\ &= \mathbb{E}_{q(x)} [\log p(x, y)] - H(q(x)) = \mathcal{L}(q(x))\end{aligned}$$

- 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>2</sup>

---

<sup>2</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>2</sup>

---

<sup>2</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>2</sup>

---

<sup>2</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>2</sup>

---

<sup>2</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

$$\mathcal{L} = \int_x q(x) \log \left( \frac{p(y, f, x)}{q(x)} \right)$$

---

<sup>3</sup>Damianou, 2015



$$\begin{aligned}\mathcal{L} &= \int_x q(x) \log \left( \frac{p(y, f, x)}{q(x)} \right) \\ &= \int_x q(x) \log \left( \frac{p(y | f)p(f | x)p(x)}{q(x)} \right)\end{aligned}$$

---

<sup>3</sup>Damianou, 2015

$$\begin{aligned}\mathcal{L} &= \int_x q(x) \log \left( \frac{p(y, f, x)}{q(x)} \right) \\ &= \int_x q(x) \log \left( \frac{p(y | f)p(f | x)p(x)}{q(x)} \right) \\ &= \int_x q(x) \log p(y | f)p(f | x) - \int_x q(x) \log \frac{q(x)}{p(x)}\end{aligned}$$

---

<sup>3</sup>Damianou, 2015

$$\begin{aligned}\mathcal{L} &= \int_x q(x) \log \left( \frac{p(y, f, x)}{q(x)} \right) \\ &= \int_x q(x) \log \left( \frac{p(y | f)p(f | x)p(x)}{q(x)} \right) \\ &= \int_x q(x) \log p(y | f)p(f | x) - \int_x q(x) \log \frac{q(x)}{p(x)} \\ &= \tilde{\mathcal{L}} - \text{KL}(q(x) \parallel p(x))\end{aligned}$$

---

<sup>3</sup>Damianou, 2015

$$\tilde{\mathcal{L}} = \int q(x) \log p(y|f)p(f|x)dfdx$$

- Has not eliviate the problem at all,  $x$  still needs to go through  $f$  to reach the data
- Idea of sparse approximations<sup>4</sup>

---

<sup>4</sup>Candela et al., [2005](#)

$$p(f, u \mid x, z)$$

- Add another set of samples from the same prior
- Conditional distribution

---

<sup>5</sup>Titsias et al., [2010](#)

$$p(f, u | x, z) = p(f | u, x, z)p(u | z)$$

- Add another set of samples from the same prior
- Conditional distribution

---

<sup>5</sup>Titsias et al., [2010](#)

$$\begin{aligned} p(f, u \mid x, z) &= p(f \mid u, x, z)p(u \mid z) \\ &= \mathcal{N}(f \mid K_{fu}K_{uu}^{-1}u, K_{ff} - K_{fu}K_{uu}^{-1}K_{uf})\mathcal{N}(u \mid \mathbf{0}, K_{uu}) \end{aligned}$$

- Add another set of samples from the same prior
- Conditional distribution

---

<sup>5</sup>Titsias et al., 2010

$$p(y, f, u, x | z) = p(y | f)p(f | u, x)p(u | z)p(x)$$

- we have done nothing to the model, just project an additional set of marginals from the GP
- *however* we will now **interpret**  $u$  and  $z$  not as **random** variables but **variational** parameters
- i.e. the variational distribution  $q(\cdot)$  is parametrised by these



- Variational distributions are approximations to intractable posteriors,

$$q(u) \approx p(u \mid y, x, z, f)$$

$$q(f) \approx p(f \mid u, x, z, y)$$

$$q(x) \approx p(x \mid y)$$

- Variational distributions are approximations to intractable posteriors,

$$q(u) \approx p(u \mid y, x, z, f)$$

$$q(f) \approx p(f \mid u, x, z, y)$$

$$q(x) \approx p(x \mid y)$$

- Bound is **tight** if  $u$  completely represents  $f$  i.e.  $u$  is sufficient statistics for  $f$

$$q(f) \approx p(f \mid u, x, z, y) = p(f \mid u, x, z)$$

$$\tilde{\mathcal{L}} = \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y, f, y | x, z)}{q(f)q(u)}$$

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y, f, y | x, z)}{q(f)q(u)} \\ &= \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y | f)p(f | u, x, z)p(u | z)}{q(f)q(u)}\end{aligned}$$

- Assume that  $u$  is sufficient statistics of  $f$

$$q(f) = p(f | u, x, z)$$

$$\tilde{\mathcal{L}} = \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y | f)p(f | u, x, z)p(u | z)}{q(f)q(u)}$$

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{q(f)q(u)} \\ &= \int_{x,f,u} p(f|u,x,z)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{p(f|u,x,z)q(u)}\end{aligned}$$

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{q(f)q(u)} \\ &= \int_{x,f,u} p(f|u,x,z)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{p(f|u,x,z)q(u)}\end{aligned}$$

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{q(f)q(u)} \\ &= \int_{x,f,u} p(f|u,x,z)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{p(f|u,x,z)q(u)} \\ &= \int_{x,f,u} p(f|u,x,z)q(u)q(x) \log \frac{p(y|f)p(u|z)}{q(u)}\end{aligned}$$

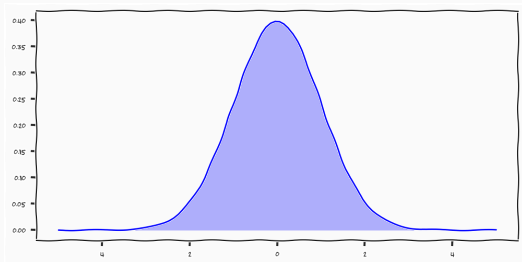


$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{x,f,u} q(f)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{q(f)q(u)} \\ &= \int_{x,f,u} p(f|u,x,z)q(u)q(x) \log \frac{p(y|f)p(f|u,x,z)p(u|z)}{p(f|u,x,z)q(u)} \\ &= \int_{x,f,u} p(f|u,x,z)q(u)q(x) \log \frac{p(y|f)p(u|z)}{q(u)} \\ &= \mathbb{E}_{p(f|u,x,z)} [p(y|f)] - \text{KL}(q(u) \parallel p(u|z))\end{aligned}$$

$$\mathcal{L} = \mathbb{E}_{p(f|u,x,z)}[p(y | f)] - \text{KL}(q(u) \parallel p(u | z)) - \text{KL}(q(x) \parallel p(x))$$

- Expectation tractable (for some co-variances)
- Allows us to place priors and not "regularisers" over the latent representation
- Stochastic inference Hensman et al., [2013](#)
- Importantly  $p(x)$  only appears in  $\text{KL}(\cdot \parallel \cdot)$  term!

# Latent Space Priors



$$p(x) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

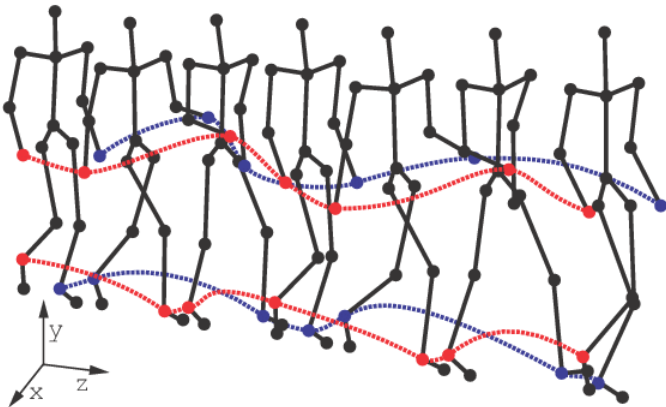
$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma e^{-\sum_d^D \alpha_d \cdot (x_{i,d} - x_{j,d})^2}$$

GPy

Code

```
[ ]python RBF(...,ARD=True) Matern32(...,ARD=True)
```

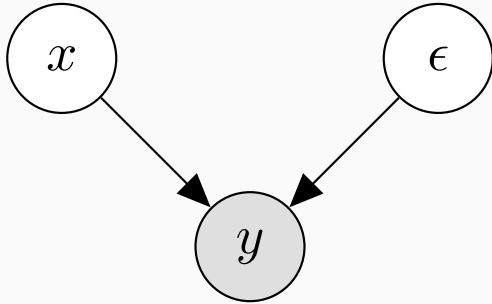
# Dynamic Prior



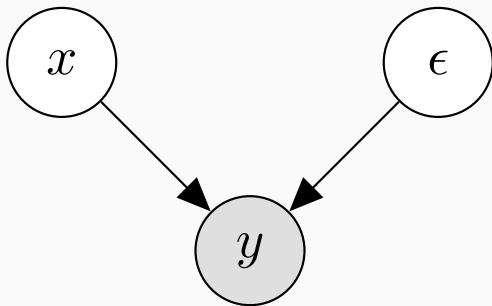
$$p(x | t) = \mathcal{N}(\mu_t, K_t)$$

# Structured Latent Spaces

---



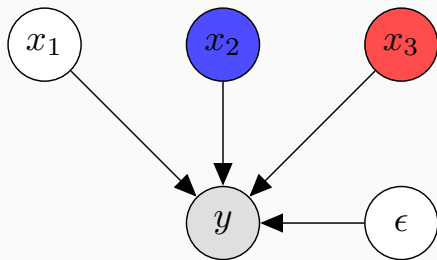
$$y = f(x) + \epsilon$$



$$y - \epsilon = f(x)$$



# Factor Analysis

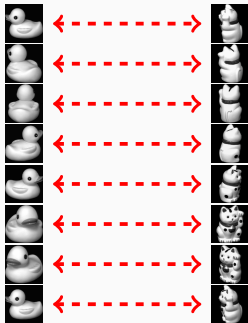


$$y = f(x_1, x_2, x_3) + \epsilon$$

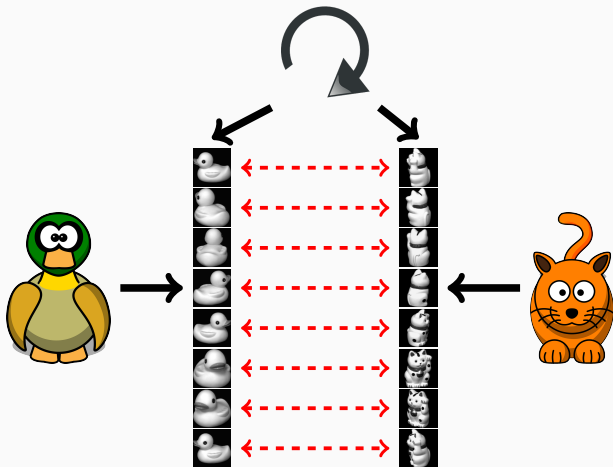
# Alignments



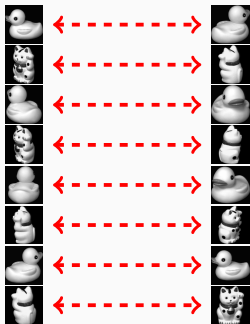
# Alignments



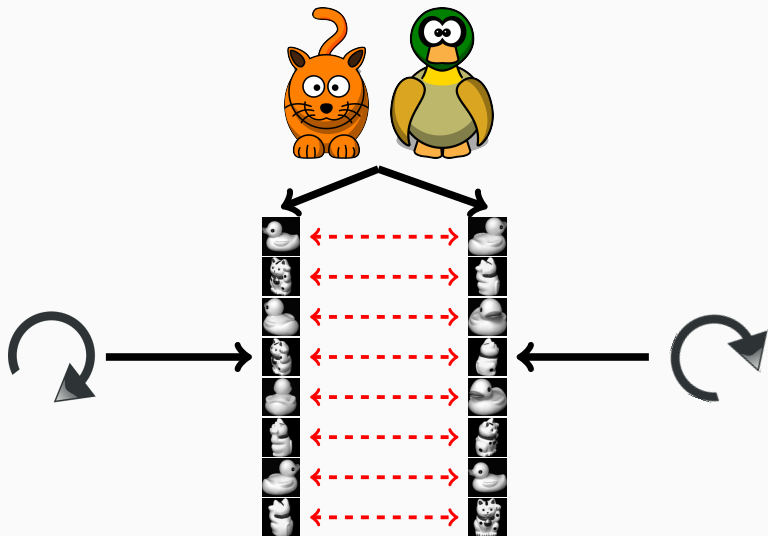
# Alignments



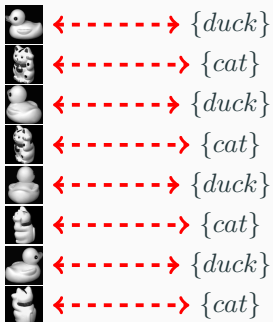
# Alignments



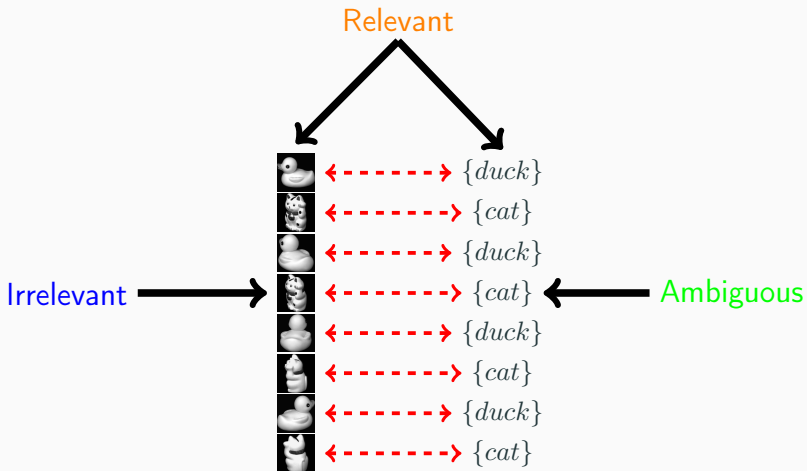
# Alignments



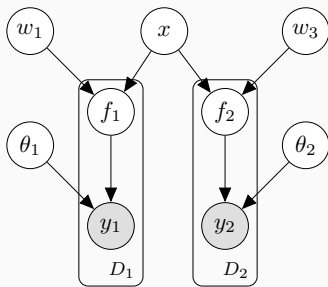
# Alignments



# Alignments



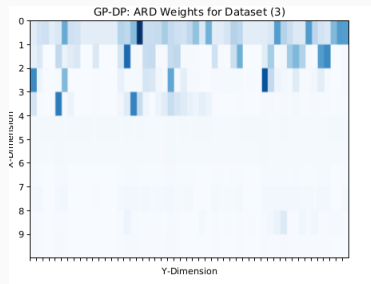
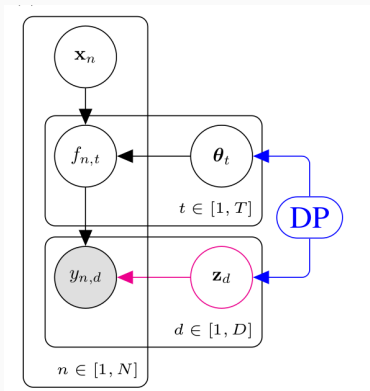




$$y_1 = f(w_1^T x) \quad y_2 = f(w_2^T x)$$

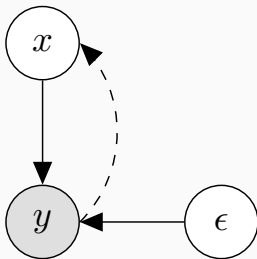
---

<sup>6</sup>Damianou et al., 2016



<sup>7</sup>Lawrence et al., 2019

## Constrained Latent Space<sup>8</sup>

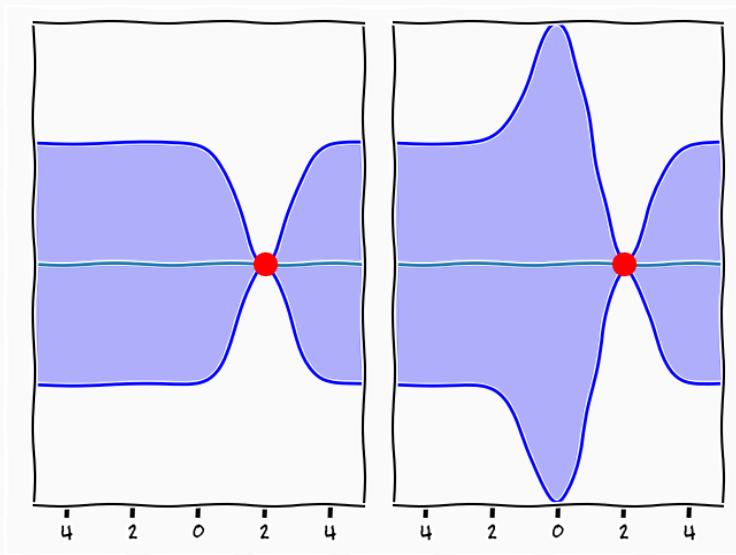


$$y = f(g(y)) + \epsilon$$

---

<sup>8</sup>Lawrence et al., 2006

# Geometry

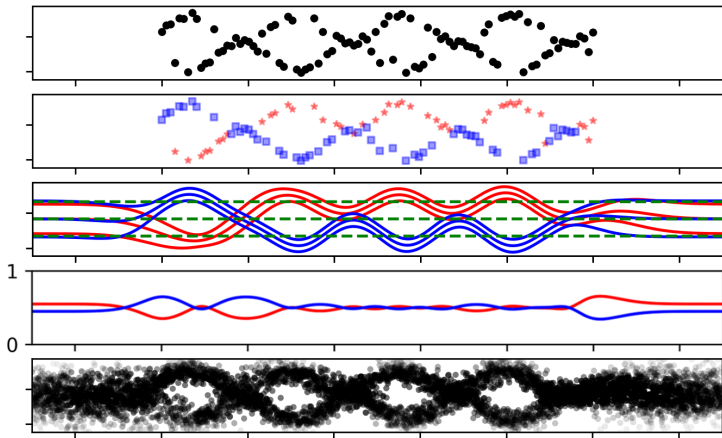


$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{F}) p(\mathbf{F}|\mathbf{X}, \mathbf{X}^{(C)}) p(\mathbf{X}^{(C)}) d\mathbf{F} d\mathbf{X}^{(C)}.$$

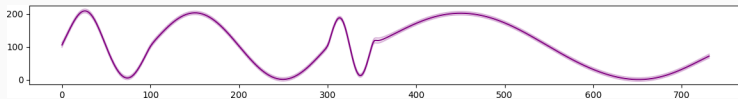
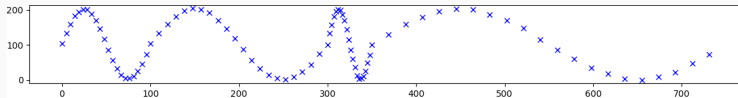
---

<sup>9</sup>Bodin et al., [2017](#), Yousefi et al., [2016](#)

# Discrete



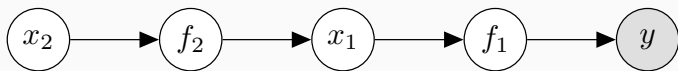
# Continuous



# Composite Gaussian Processes

---

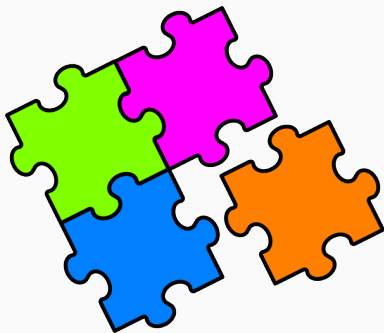




---

<sup>10</sup>Damianou et al., [2013](#)

# Composite Functions



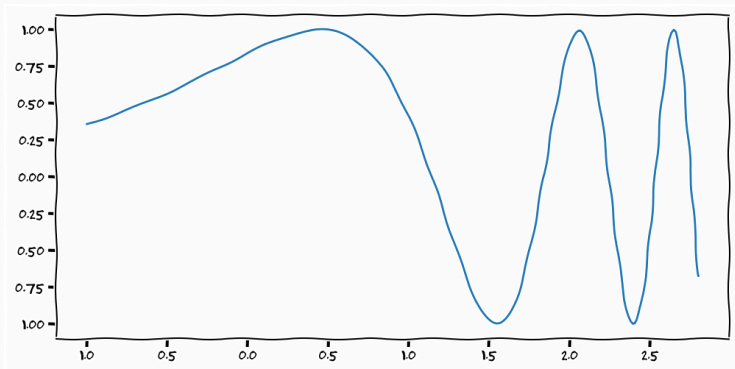
$$y = f_k(f_{k-1}(\dots f_0(x))) = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$$

# When do I want Composite Functions

$$y = f_k \circ f_{k-1} \circ \cdots \circ f_1(x)$$

1. My generative process is composite
  - my prior knowledge is composite
2. I want to "re-parametrise" my kernel in a learning setting
  - i have knowledge of the re-parametrisation

## Because we lack "models"?



## Diff Levels of Abstraction

- Hierarchical Learning
  - Natural progression from low level to high level structure as seen in natural complexity
  - Easier to monitor what is being learnt and to guide the machine to better subspaces
  - A good lower level representation can be used for many distinct tasks

Feature representation



3rd layer  
"Objects"



2nd layer  
"Object parts"

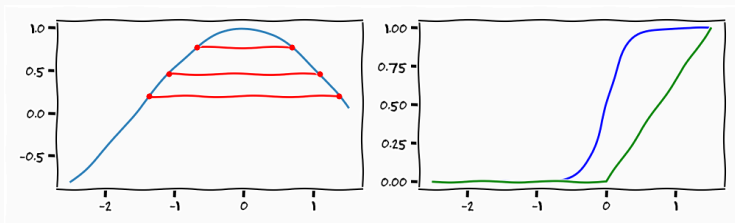


1st layer  
"Edges"



Pixels

# Composite functions

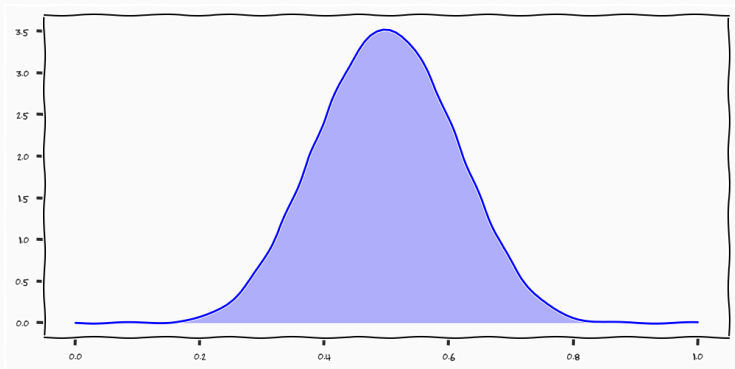


$$y = f_k(f_{k-1}(\dots f_0(x))) = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$$

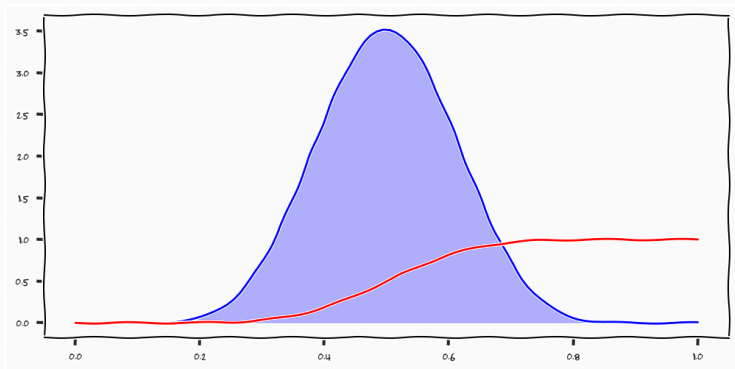
$$\text{Kern}(f_1) \subseteq \text{Kern}(f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Kern}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1)$$

$$\text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2) \subseteq \dots \subseteq \text{Im}(f_k)$$

# Sampling

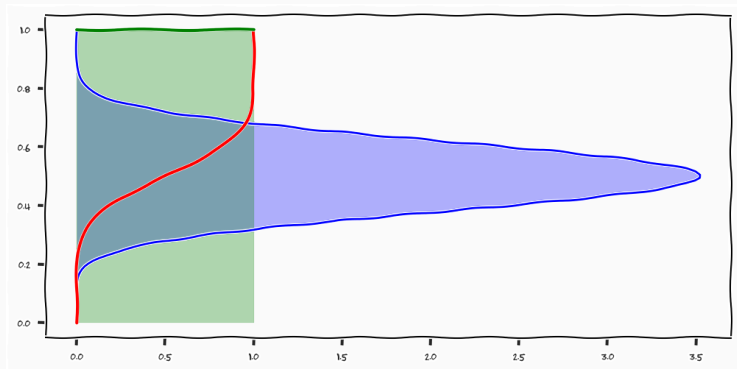


# Sampling

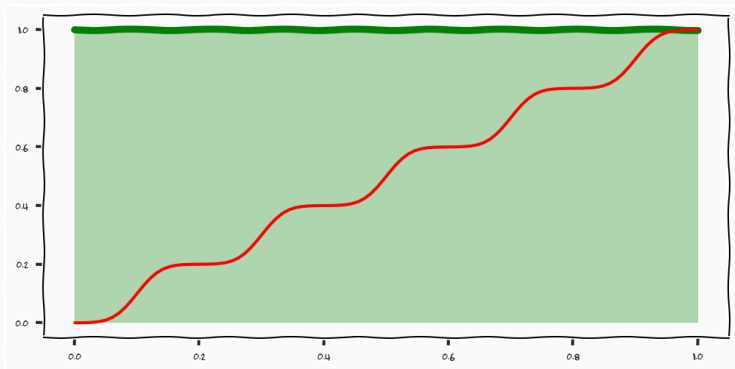




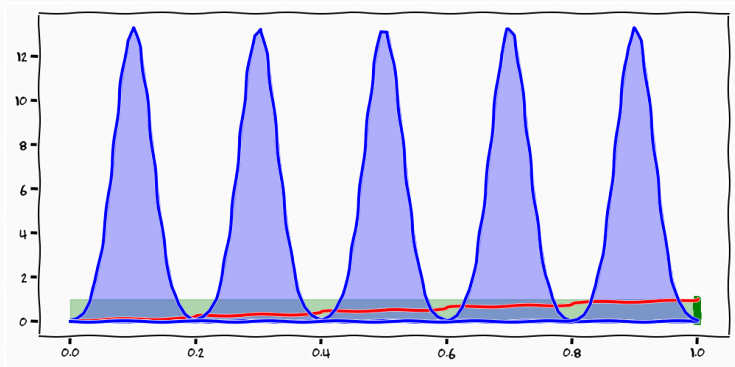
# Sampling



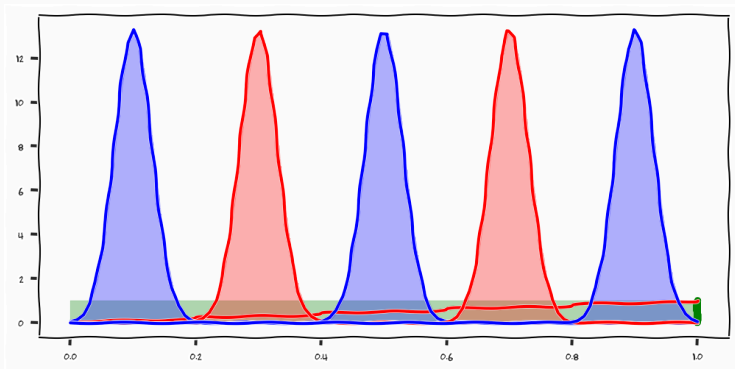
# Change of Variables



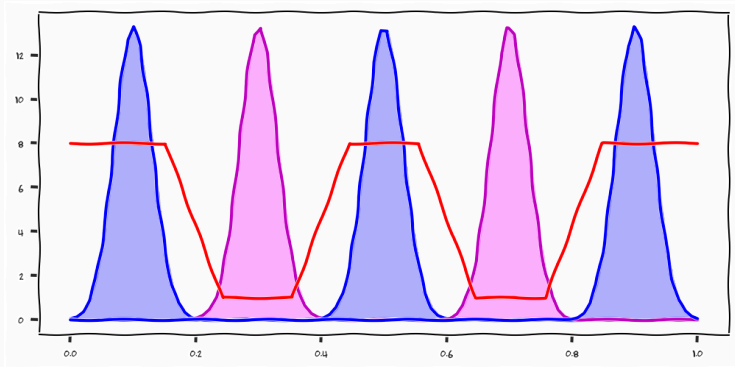
# Change of Variables



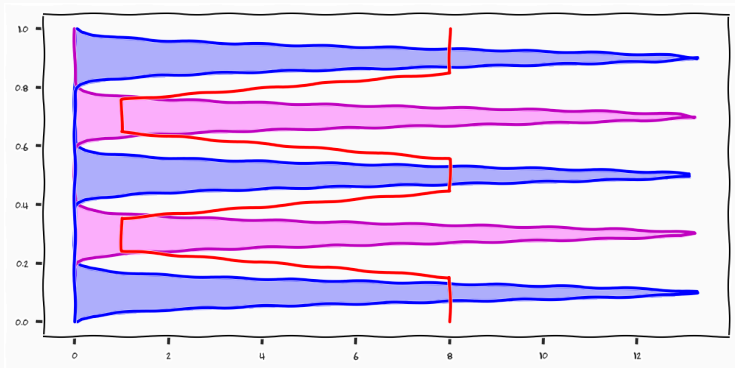
# Change of Variables



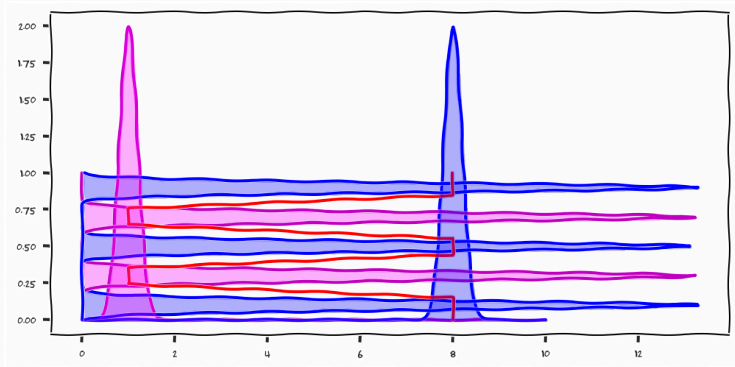
# Change of Variables



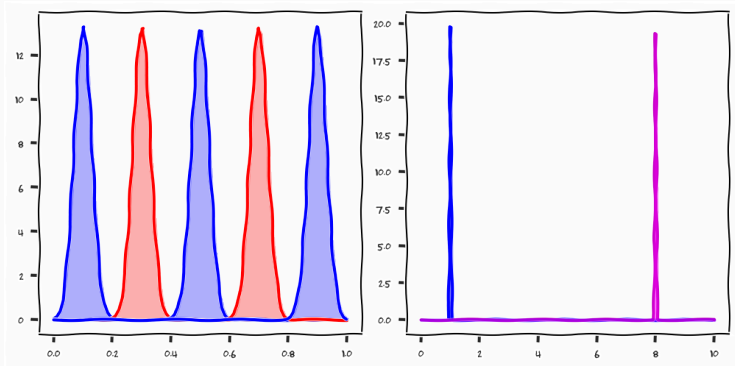
# Change of Variables



# Change of Variables



# Change of Variables





"I'm Bayesian therefore I am superior"



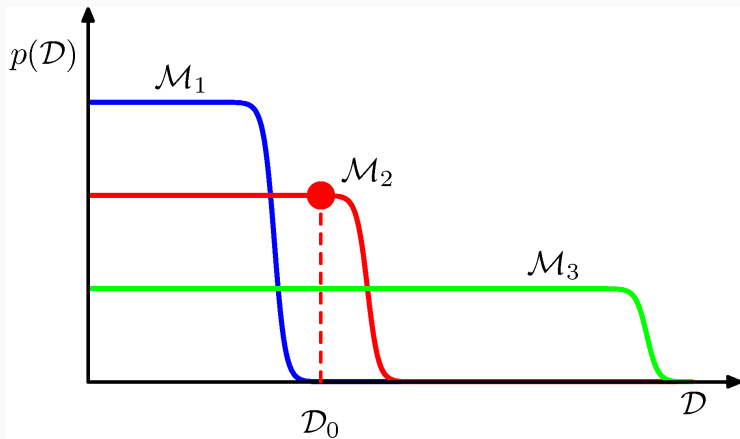
## Because we want to hang out with the cool kids



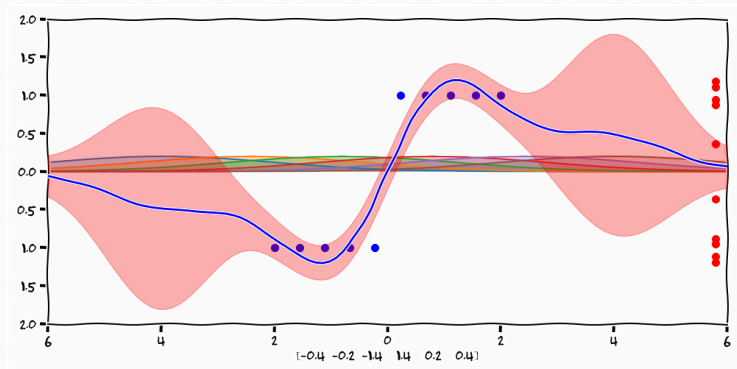
*Deep Learning is a bit like smoking, you know that its wrong but you do it anyway because you want to look cool.*

*– Fantomens Djungelordspråk*

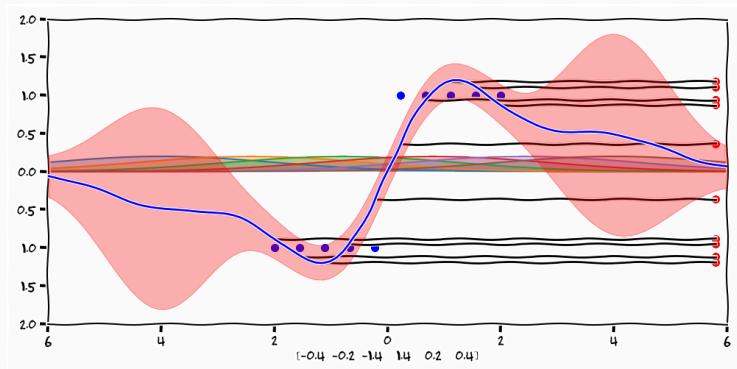
# MacKay plot



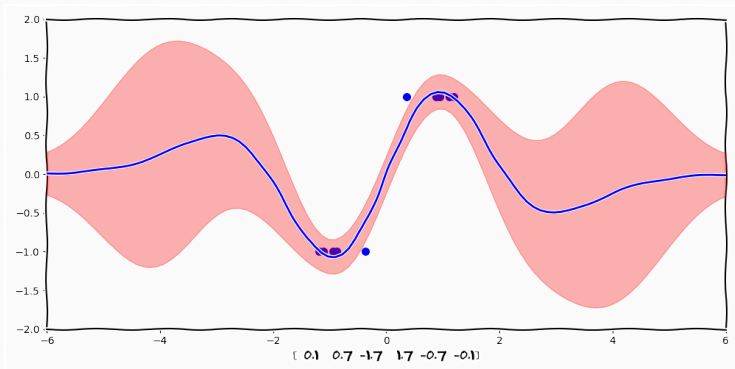
# Composite Functions



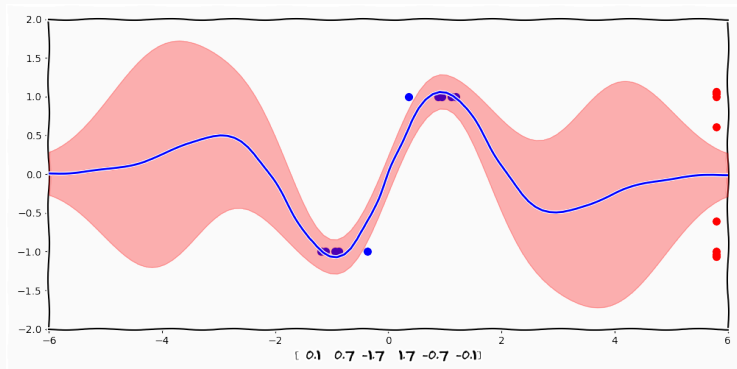
# Composite Functions



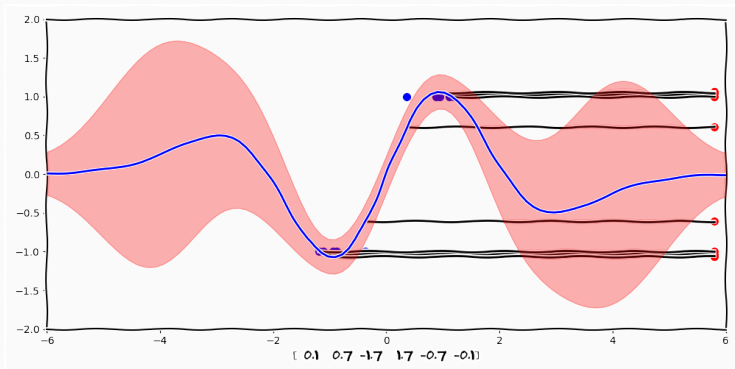
# Composite Functions



# Composite Functions

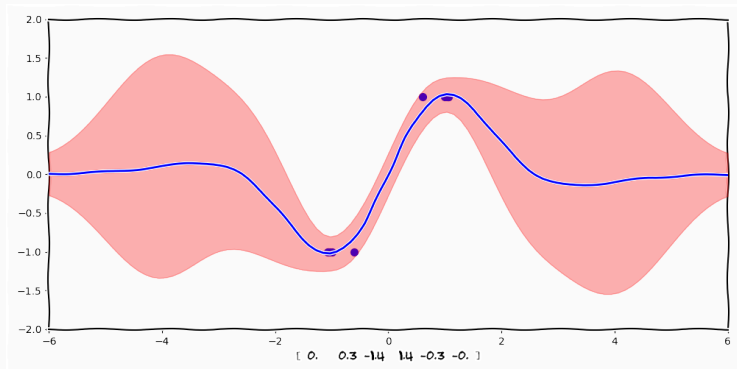


# Composite Functions

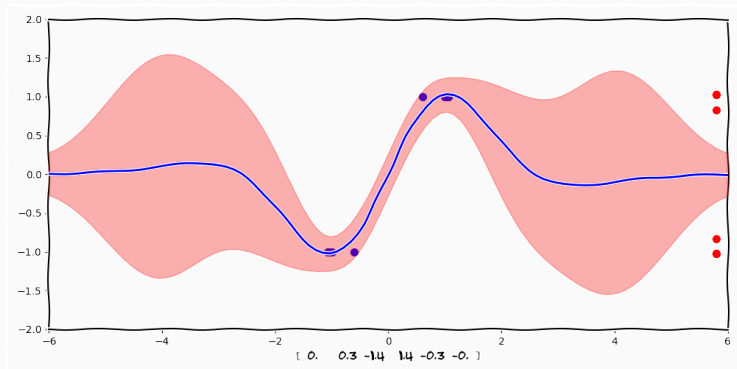




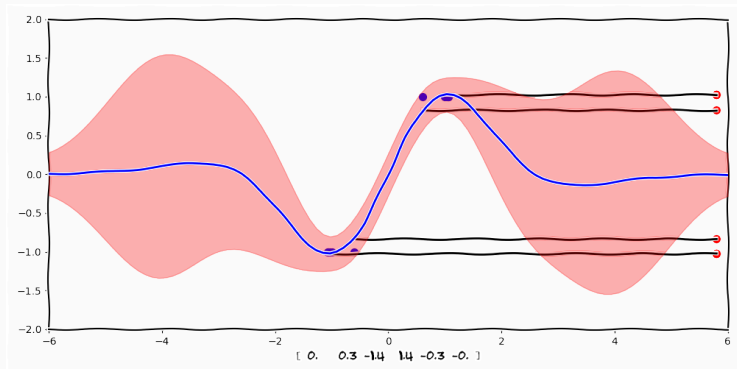
# Composite Functions



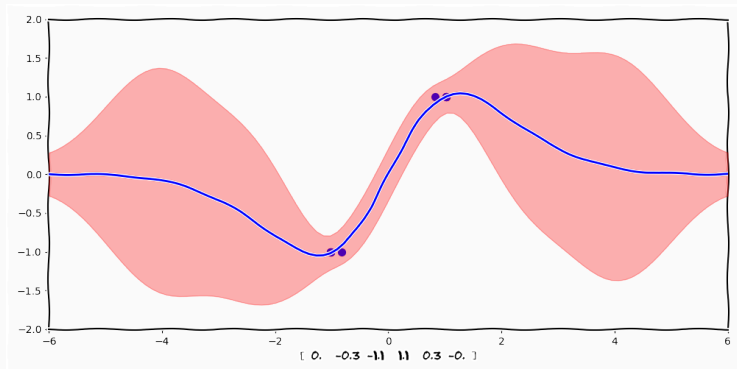
# Composite Functions



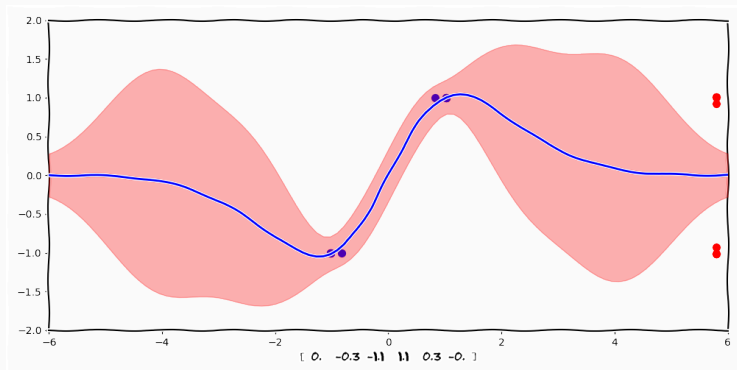
# Composite Functions



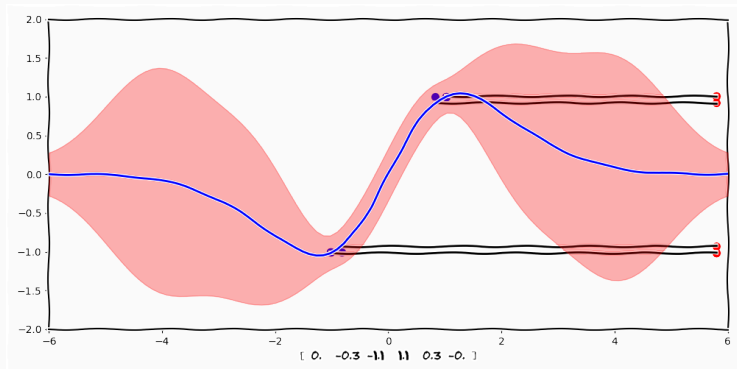
# Composite Functions



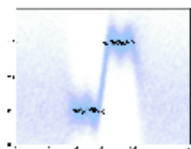
# Composite Functions



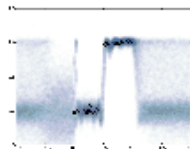
# Composite Functions



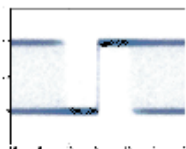
# The Final Composition



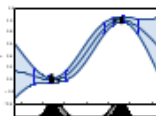
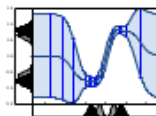
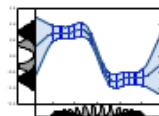
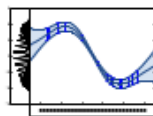
(a) GP



(b) 2 layers

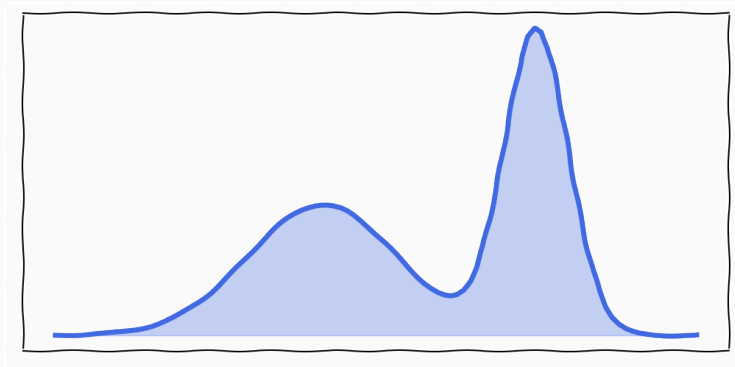


(c) 4 layers



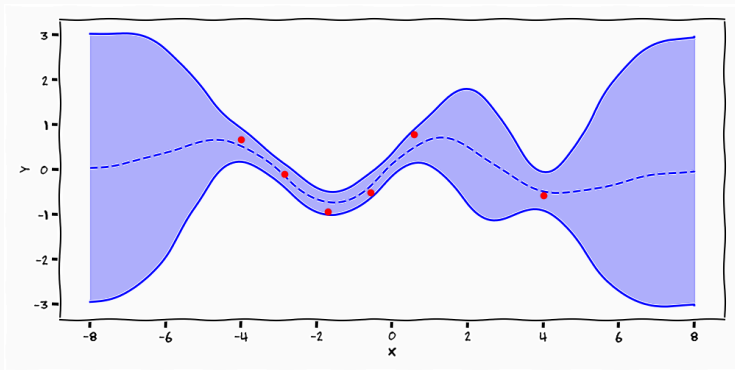
(d) Hidden spaces for 4 layer model

Remember why we did this in the first place

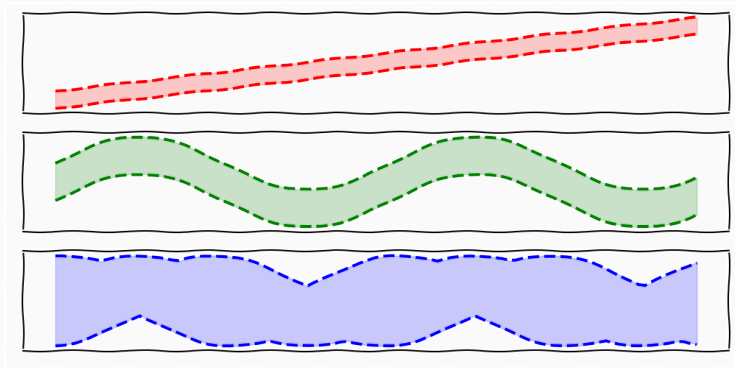




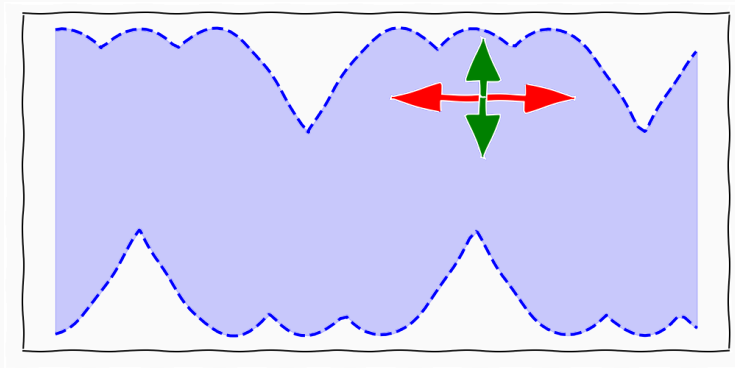
## These damn plots



# It gets worse



It gets even worse



# Approximate Inference

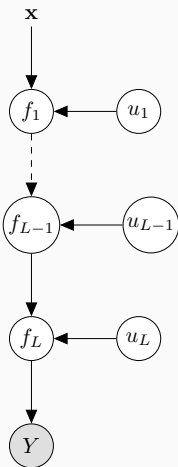
- Sufficient statistics

$$q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) = p(\mathbf{F}|\mathbf{Y}, \mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})q(\mathbf{X})$$

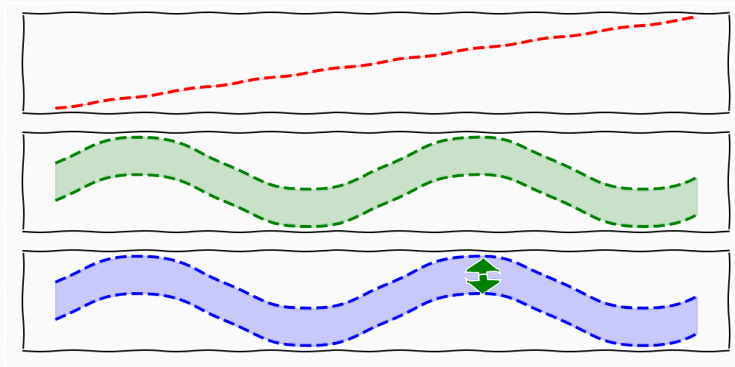
$$= p(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})q(\mathbf{X})$$

- Mean-Field

$$q(\mathbf{U}) = \prod_i^L q(\mathbf{U}_i)$$



# The effect



## What have we lost

- Our priors are not reflected correctly
  - → we cannot interpret the results
- No intermediate uncertainties
  - → we cannot do sequential decision making

# What have we lost

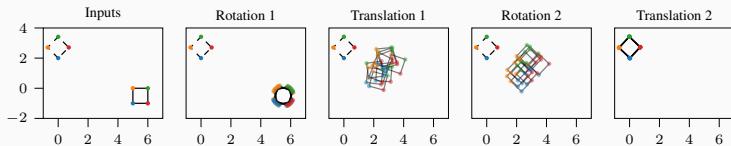
- Our priors are not reflected correctly
  - → we cannot interpret the results
- No intermediate uncertainties
  - → we cannot do sequential decision making
- We are performing a massive computational overhead for very little use

## What have we lost

- Our priors are not reflected correctly
  - → we cannot interpret the results
- No intermediate uncertainties
  - → we cannot do sequential decision making
- We are performing a massive computational overhead for very little use
- *"... throwing out the baby with the bathwater... "*



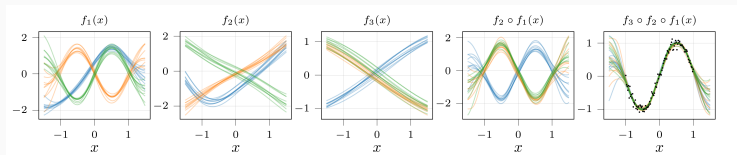
# What we really want<sup>11</sup>



---

<sup>11</sup>Ustyuzhaninov et al., 2019

# What we really want<sup>12</sup>



---

<sup>12</sup>Ustyuzhaninov et al., 2019

## Summary

---

- Unsupervised learning<sup>13</sup> is **very** hard.

---

<sup>13</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning<sup>13</sup> is **very** hard.
  - *Its actually not, its really really easy.*

---

<sup>13</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning<sup>13</sup> is **very** hard.
  - *Its actually not, its really really easy.*
- Relevant assumptions needed to learn anything useful

---

<sup>13</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning<sup>13</sup> is **very** hard.
  - *Its actually not, its really really easy.*
- Relevant assumptions needed to learn anything useful
- Strong assumptions needed to learn anything from "sensible" amounts of data

---

<sup>13</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning<sup>13</sup> is **very** hard.
  - *Its actually not, its really really easy.*
- Relevant assumptions needed to learn anything useful
- Strong assumptions needed to learn anything from "sensible" amounts of data
- Stochastic processes such as GPs provide strong, interpretative assumptions that aligns well to our intuitions allowing us to make **relevant** assumptions

---

<sup>13</sup>I would argue that there is no such thing



- Composite functions **cannot** model more things

- Composite functions **cannot** model more things
- However, they can easily warp the input space to model **less** things

- Composite functions **cannot** model more things
- However, they can easily warp the input space to model **less** things
- This leads to high requirements on data

## Summary II

- Composite functions **cannot** model more things
- However, they can easily warp the input space to model **less** things
- This leads to high requirements on data
- Even bigger need for uncertainty propagation, we cannot assume noiseless data

## Summary II

- Composite functions **cannot** model more things
- However, they can easily warp the input space to model **less** things
- This leads to high requirements on data
- Even bigger need for uncertainty propagation, we cannot assume noiseless data
- We need to think about correlated uncertainty, not marginals

eof

## Reference

---



## References




---

-  Bodin, Erik, Neill D. F. Campbell, and Carl Henrik Ek (2017). *Latent Gaussian Process Regression*.
-  Candela, Joaquin Quiñonero and Carl Edward Rasmussen (2005). “A Unifying View of Sparse Approximate Gaussian Process Regression”. In: *Journal of Machine Learning Research* 6, pp. 1939–1959.
-  Damianou, Andreas, Neil D Lawrence, and Carl Henrik Ek (2016). “Multi-view Learning as a Nonparametric Nonlinear Inter-Battery Factor Analysis”. In: *arXiv preprint arXiv:1604.04939*.



-  Damianou, Andreas C (Feb. 2015). “Deep Gaussian Processes and Variational Propagation of Uncertainty”. PhD thesis. University of Sheffield.
-  Damianou, Andreas C and Neil D Lawrence (2013). “Deep Gaussian Processes”. In: *International Conference on Artificial Intelligence and Statistical Learning*, pp. 207–215.
-  Hensman, James, N Fusi, and Neil D Lawrence (2013). “Gaussian Processes for Big Data”. In: *Uncertainty in Artificial Intelligence*.

-  Lawrence, Andrew R., Carl Henrik Ek, and Neill W. Campbell (2019). “DP-GP-LVM: A Bayesian Non-Parametric Model for Learning Multivariate Dependency Structures”. In: *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pp. 3682–3691.
-  Lawrence, Neil D. and Joaquin Quiñonero Candela (2006). “Local Distance Preservation in the GP-LVM Through Back Constraints”. In: *Proceedings of the 23rd International Conference on Machine Learning. ICML '06*. Pittsburgh, Pennsylvania, USA: ACM, pp. 513–520.

-  Titsias, Michalis and Neil D Lawrence (2010). “Bayesian Gaussian Process Latent Variable Model”. In: *International Conference on Artificial Intelligence and Statistical Learning*, pp. 844–851.
-  Ustyuzhaninov, Ivan et al. (2019). “Compositional Uncertainty in Deep Gaussian Processes”. In: *CoRR*.
-  Yousefi, Fariba, Zhenwen Dai, Carl Henrik Ek, and Neil Lawrence (2016). “Unsupervised Learning With Imbalanced Data Via Structure Consolidation Latent Variable Model”. In: *CoRR*.