

Bayesian non-parametrics and priors over functions

Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

November 22, 2017

<http://www.carlhenrik.com>

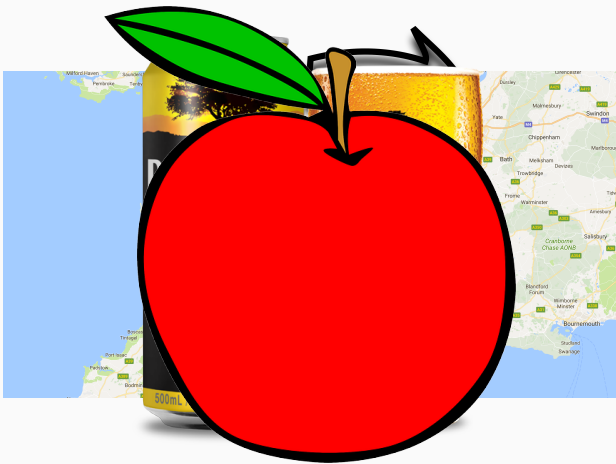
Introductions

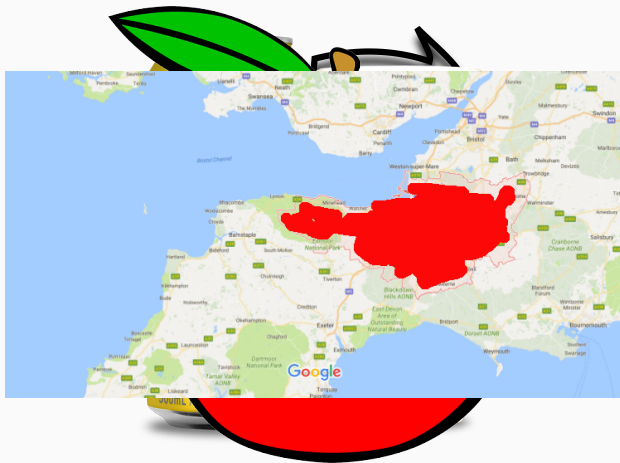


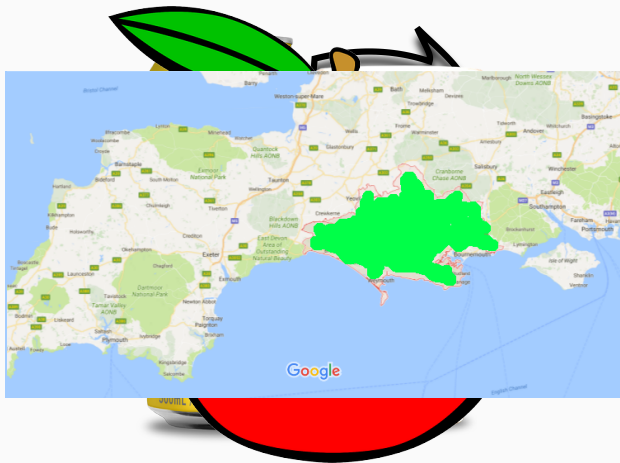


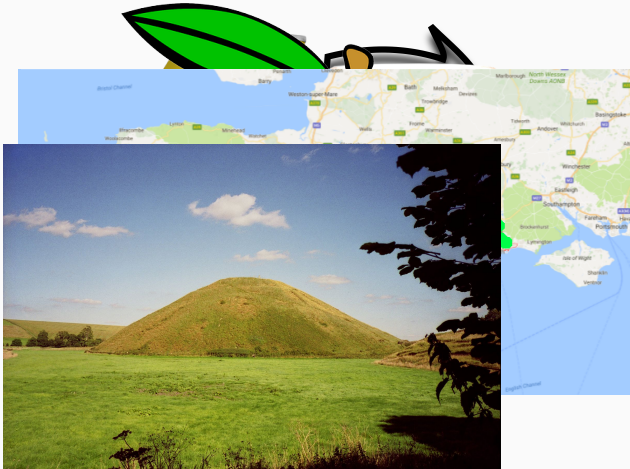


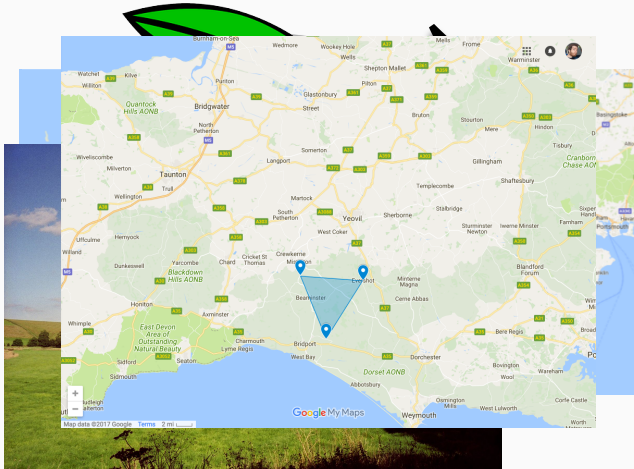




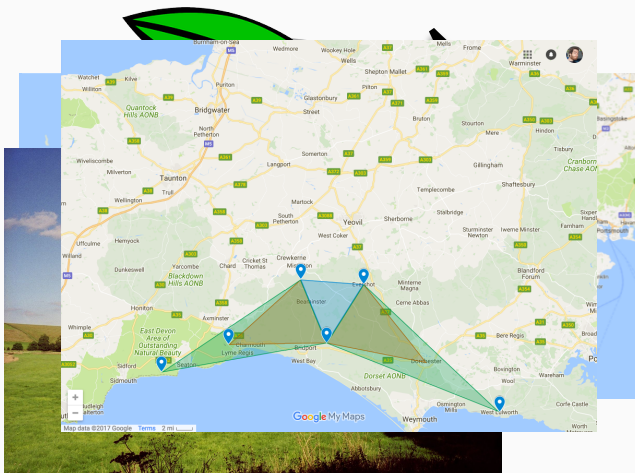


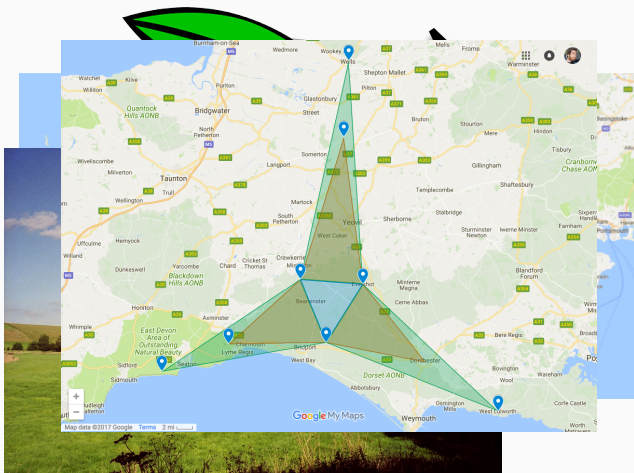


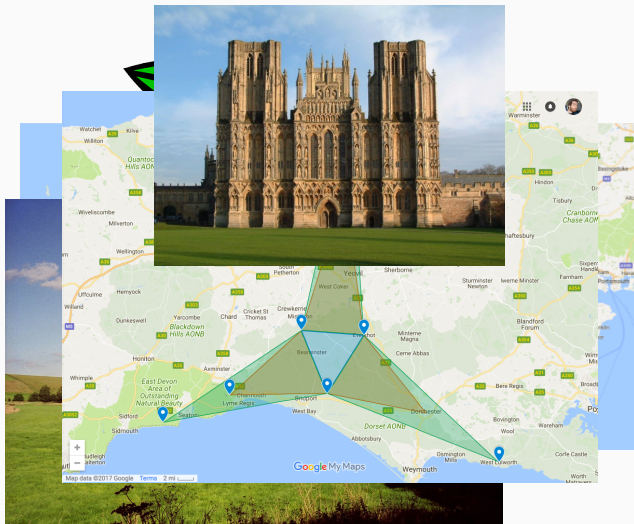


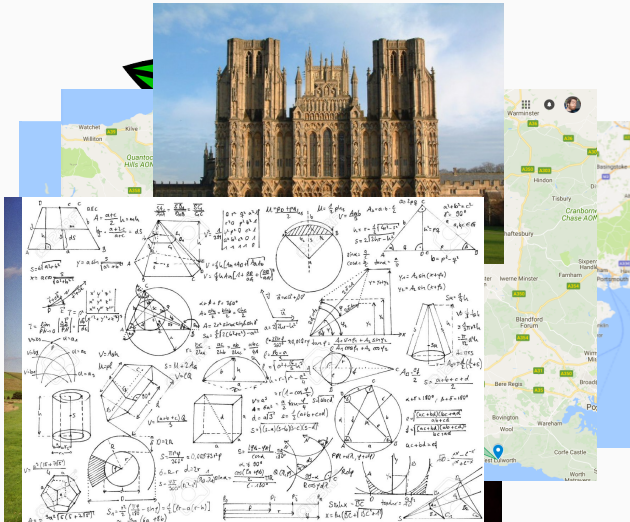


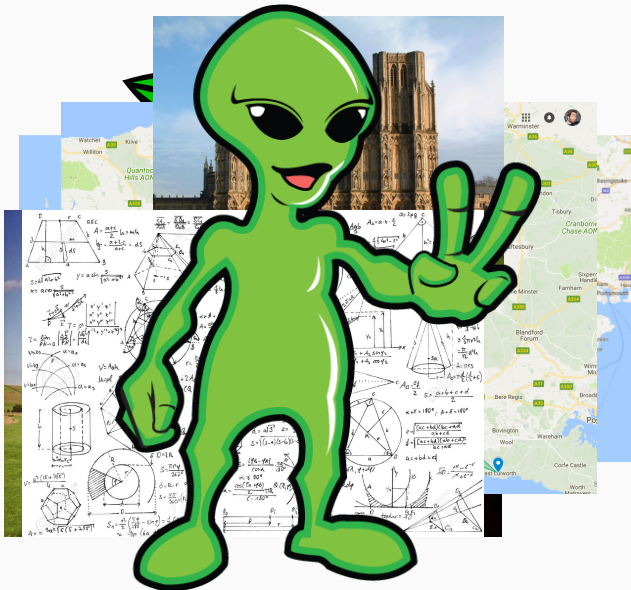






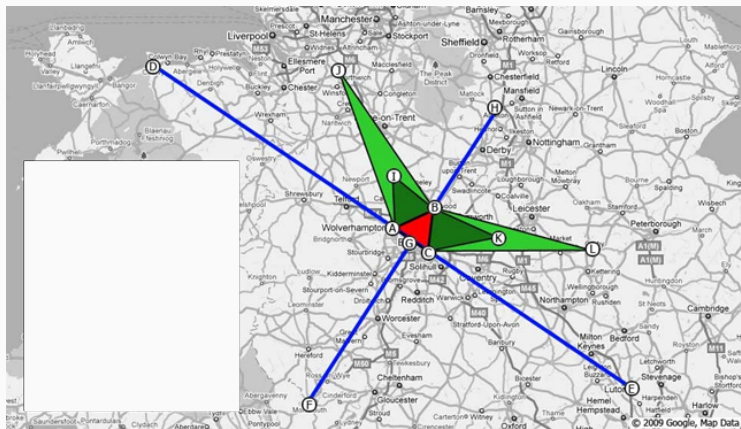








"Brooks has proved, he explains, that there were keen mathematicians here 5,000 years ago, millennia before the Greeks invented geometry: "Such is the mathematical precision, it is inconceivable that this work could have been carried out by the primitive indigenous culture we have always associated with such structures . . . all this suggests a culture existing in these islands in the past quite outside our expectation and experience today." He does not rule out extraterrestrial help." – The Guardian





"We know so little about the ancient Woolworths stores," he explains, "but we do still know their locations. I thought that if we analysed the sites we could learn more about what life was like in 2008 and how these people went about buying cheap kitchen accessories and discount CDs" – Matt Parker interviewed in The Guardian¹

¹Bad Science Blog



Laplace Demon [1]



Laplace's Demon [1]

An intelligence which at a given instant knew all the forces acting in nature and the position of every object in the universe - if endowed with a brain sufficiently vast to make all necessary calculations - could describe with a single formula the motions of the largest astronomical bodies and those of the smallest atoms. To such an intelligence, nothing would be uncertain; the future, like the past, would be an open book.

All these efforts in the search for truth tend to lead the mind continuously towards the intelligence we have just mentioned, although it will always remain infinitely distant from this intelligence.



Napoleon *"You have written this huge book on the system of the world without once mentioning the author of the universe."*



Napoleon *"You have written this huge book on the system of the world without once mentioning the author of the universe."*

Laplace *"I had no need for that assumption"*



Napoleon *"You have written this huge book on the system of the world without once mentioning the author of the universe."*

Laplace *"I had no need for that assumption"*

Laplace *"Ah, but that is a fine hypothesis. It explains so many things"*

Inductivist Fallacy



2

²Chomsky, N. A., & Fodor, J. A. (1980). The inductivist fallacy. *Language and Learning: The Debate between Jean Piaget and Noam Chomsky*, (). .

22/11/2017

Bayesian non-parametrics

$$p(\mathbf{Y}|\theta)$$

$$\mathbf{Y} \in \mathcal{Y}$$

- Task of machine learning, describe models of data

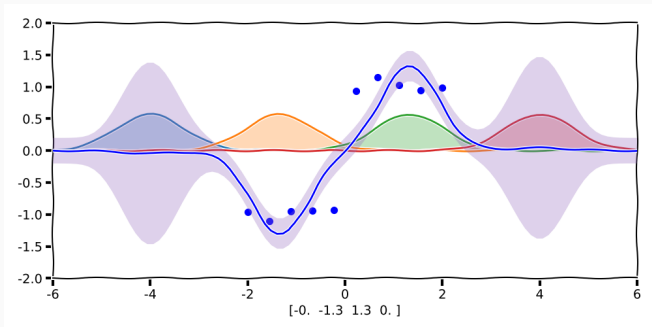
$$M \subset PM(\mathcal{Y})$$

- all probability measures on the sample space \mathcal{Y}

$$M = \{p(\mathbf{Y}|\theta)|\theta \in \mathcal{T}\}$$

- each model is indexed by θ from the parameter space \mathcal{T}

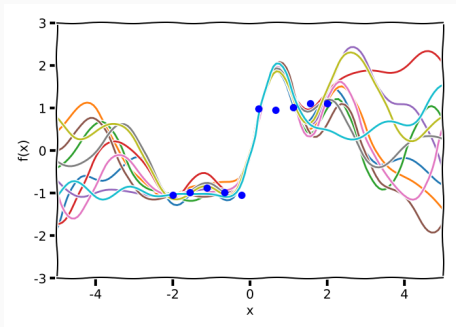
Linear Regression



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

$$\mathcal{T} = \mathbb{R}^4$$

Non-Linear Regression



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

$$\mathcal{T} = \mathbb{R}^\infty$$

Parametric vs. Non-parametric

- If \mathcal{T} is
 - infinite dimensional space we call this a non-parametric
 - finite dimensional space we call this a parametric

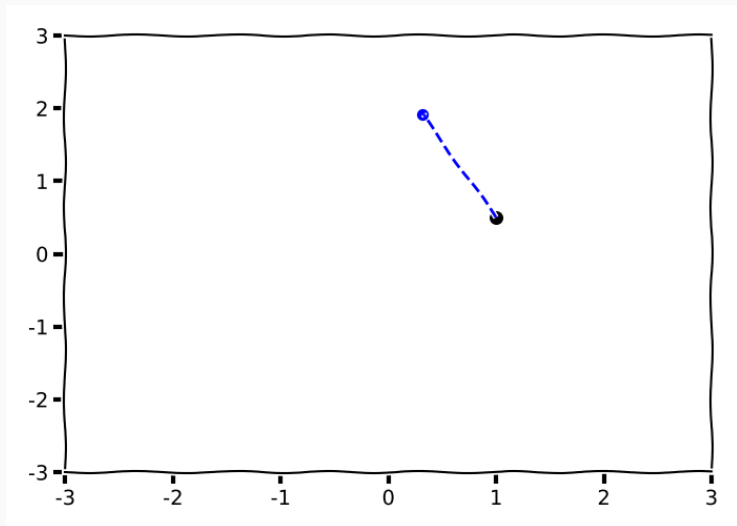
Nearest Neighbour

- Training data: $\{\mathbf{x}_i, y_i\}_{i=1}^N$
- Test data: $\{\mathbf{x}_i\}_{i=1}^M$
- Inference

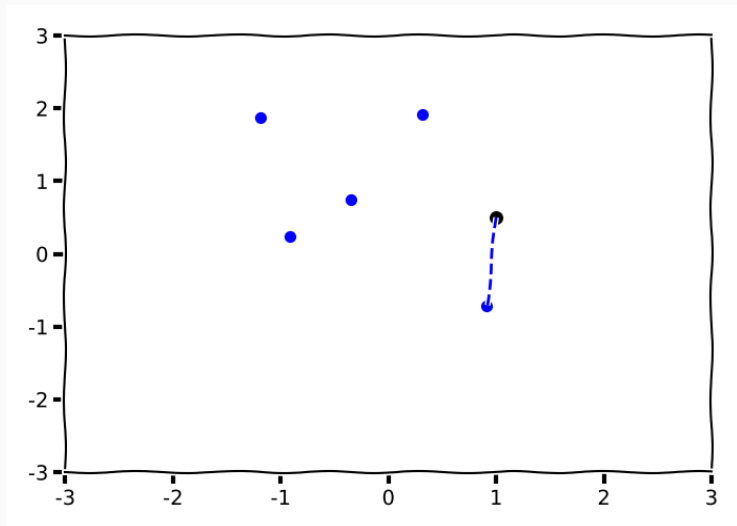
$$\hat{i} = \operatorname{argmin}_i D(\mathbf{x}_*, \mathbf{x}_i)$$

- Complexity grows with number of training data
- Does not generalise at all

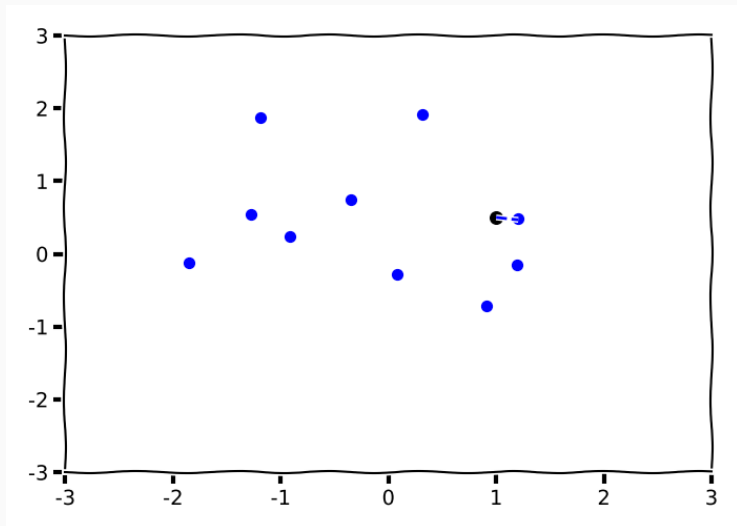
Nearest Neighbour



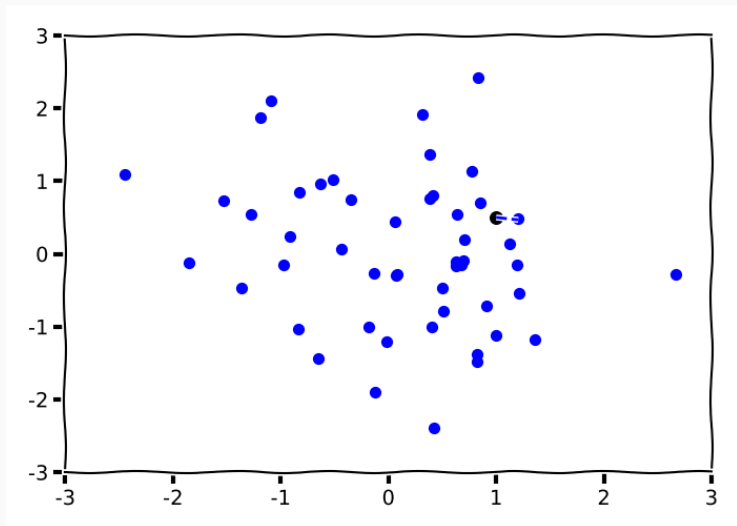
Nearest Neighbour



Nearest Neighbour



Nearest Neighbour



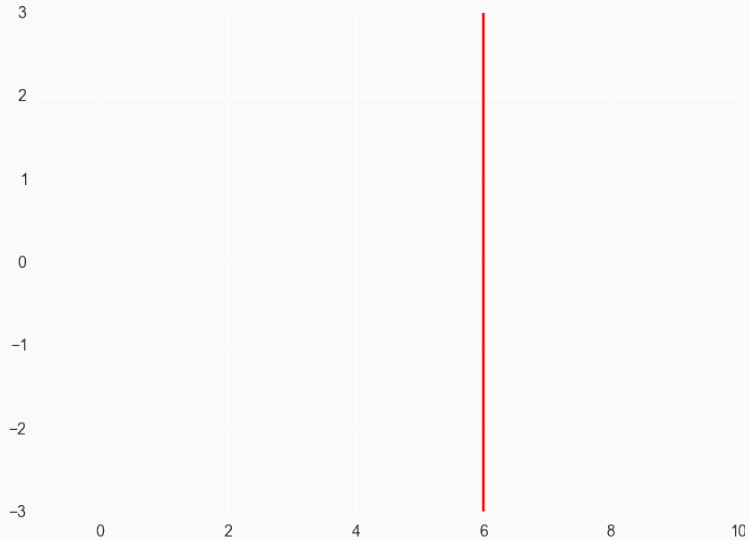
$$\theta \sim Q$$

Treating the index into the parameter space as a random variable

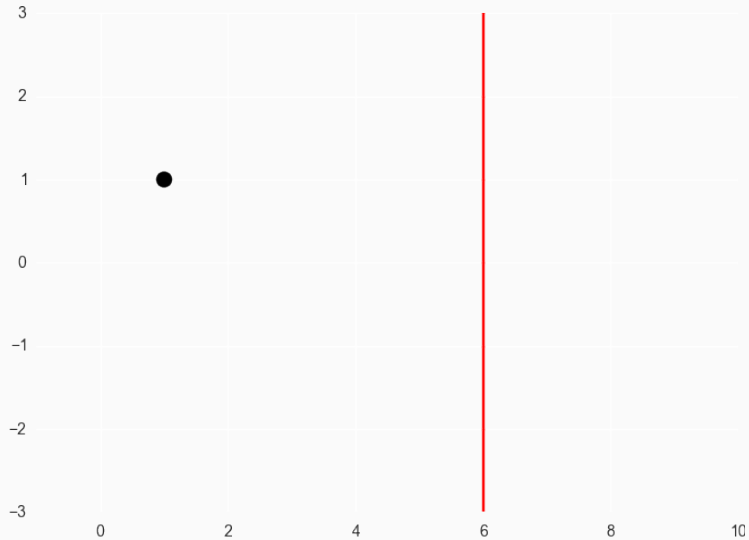
Functions



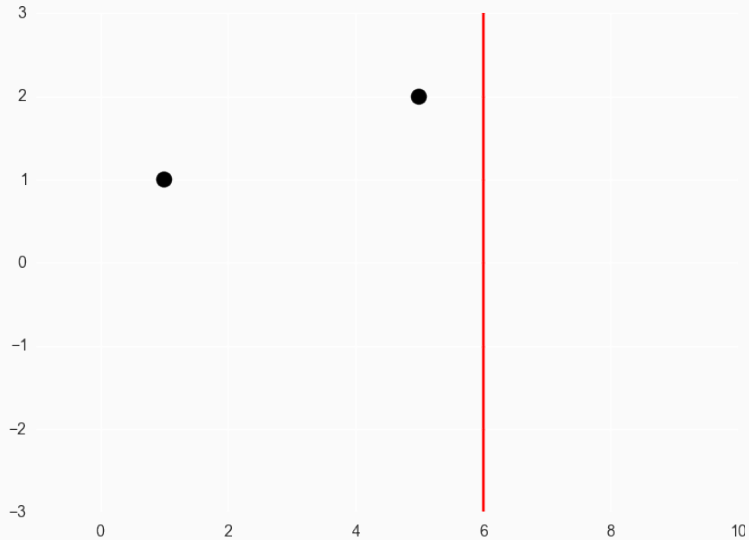
Functions



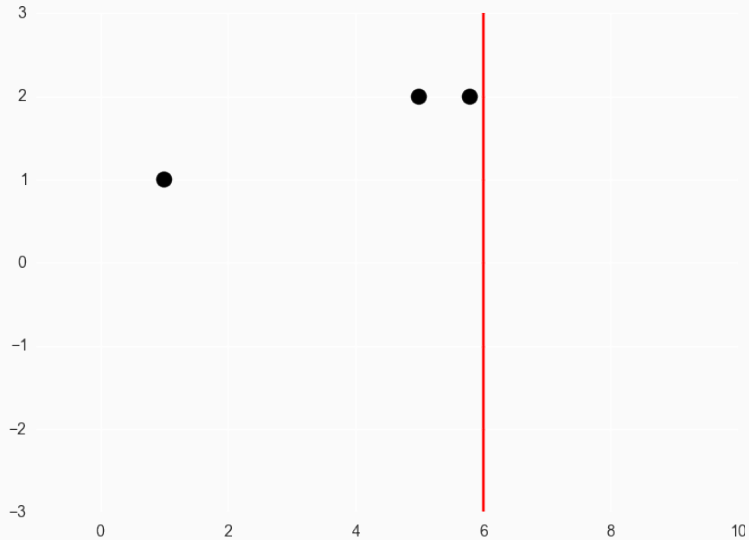
Functions



Functions

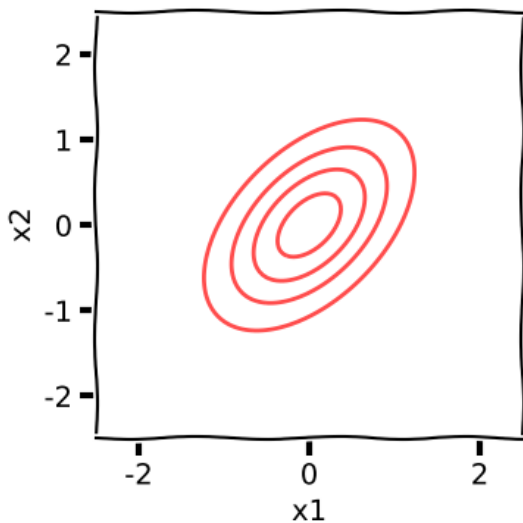


Functions

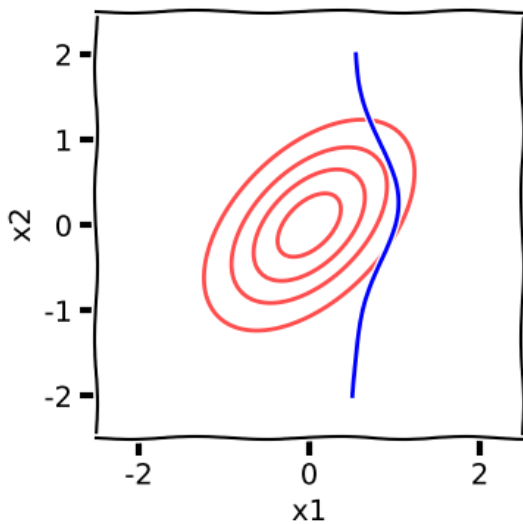


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$$

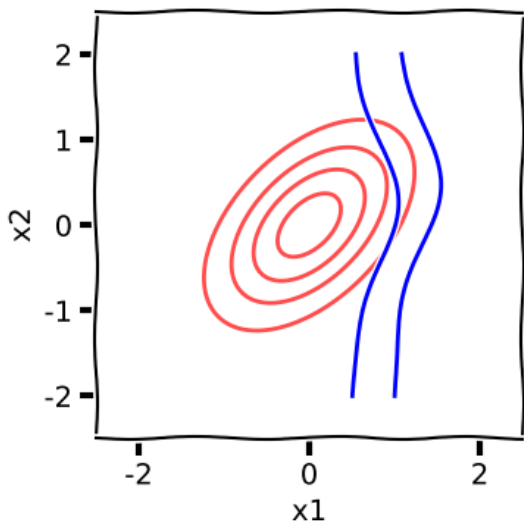
Conditional Gaussians



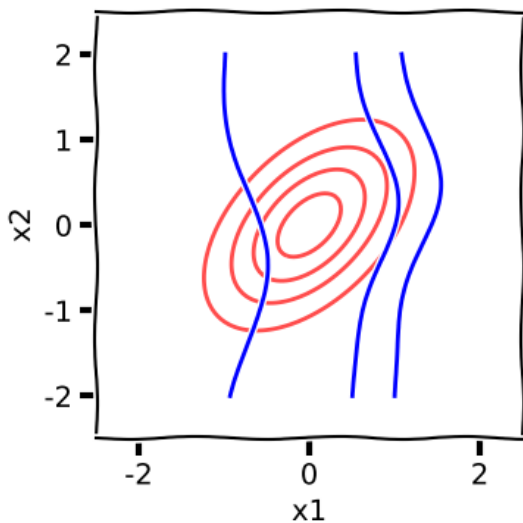
Conditional Gaussians



Conditional Gaussians

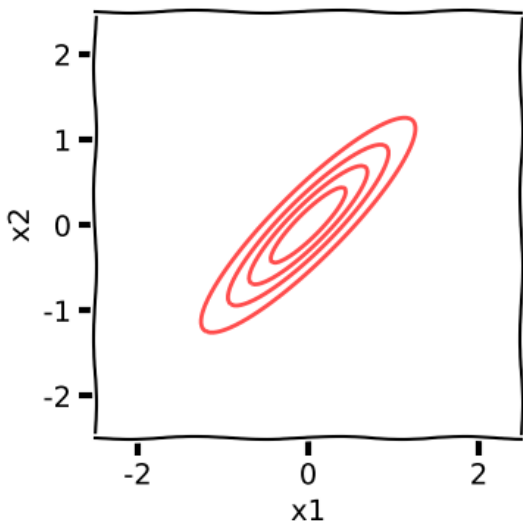


Conditional Gaussians

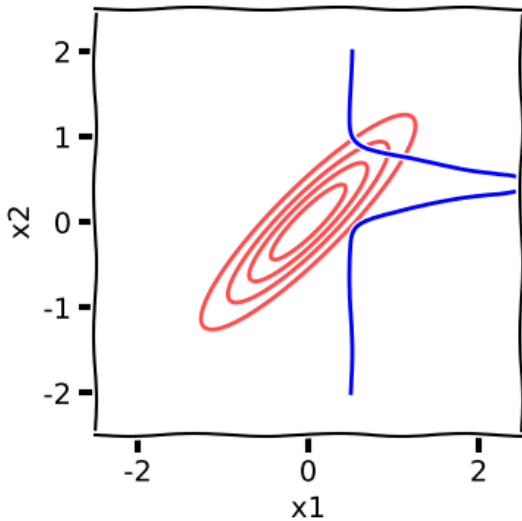


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}\right)$$

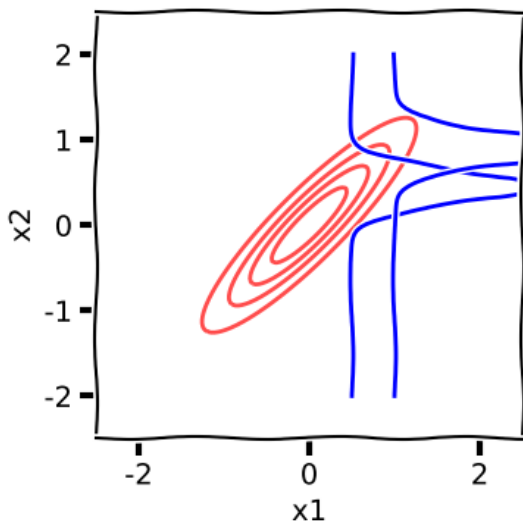
Conditional Gaussians



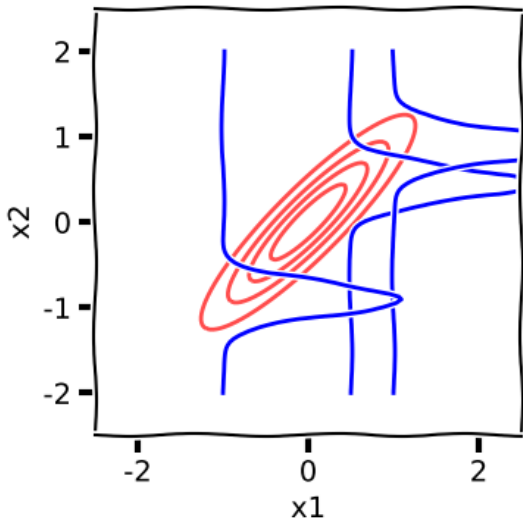
Conditional Gaussians



Conditional Gaussians

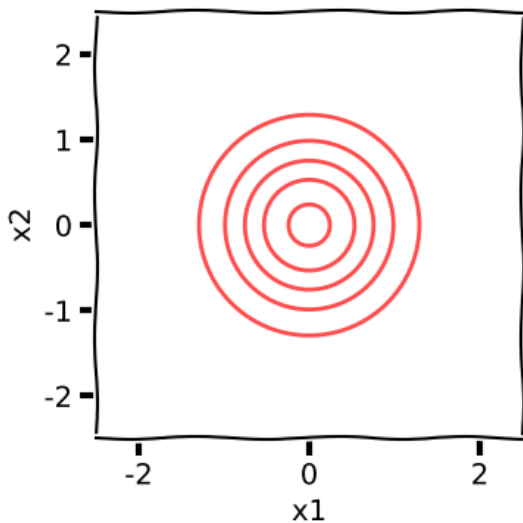


Conditional Gaussians

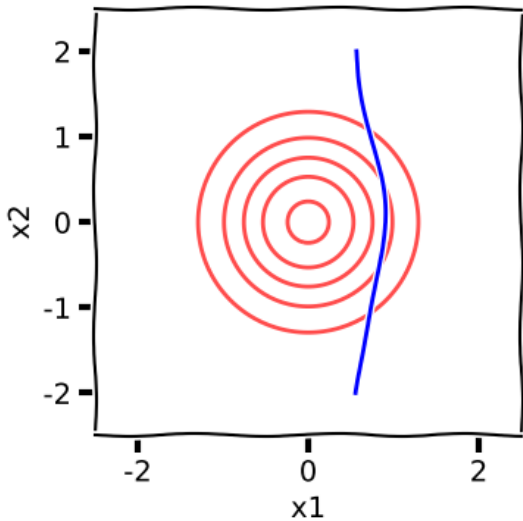


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

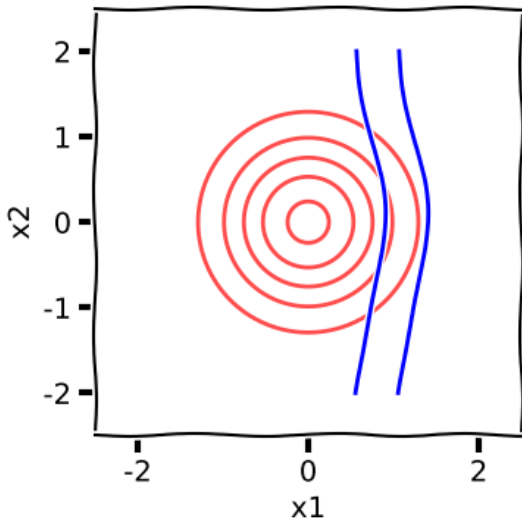
Conditional Gaussians



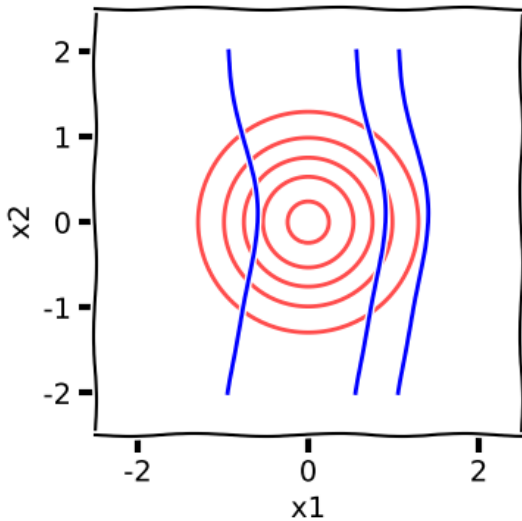
Conditional Gaussians



Conditional Gaussians

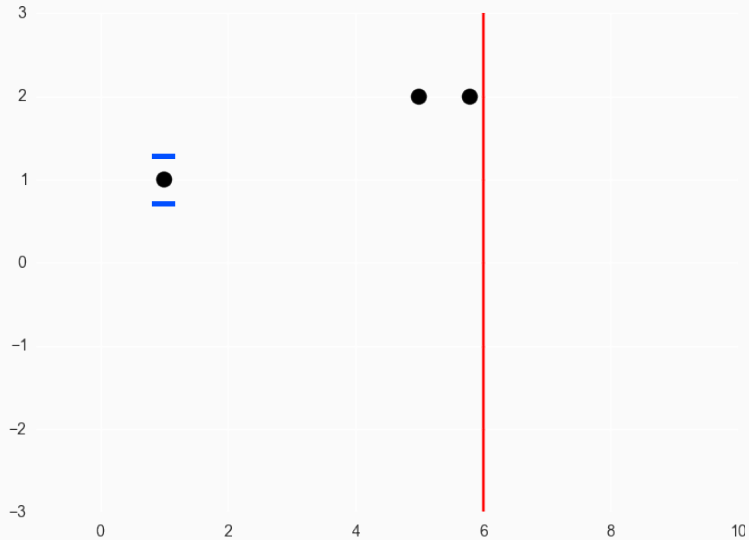


Conditional Gaussians

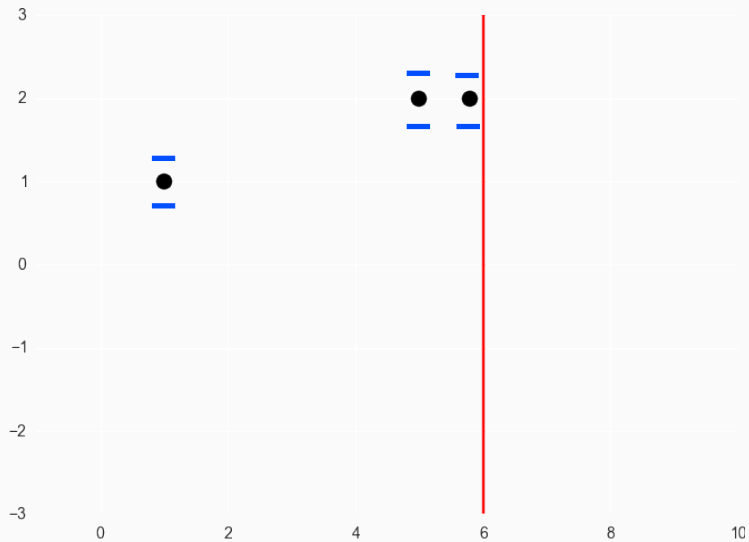




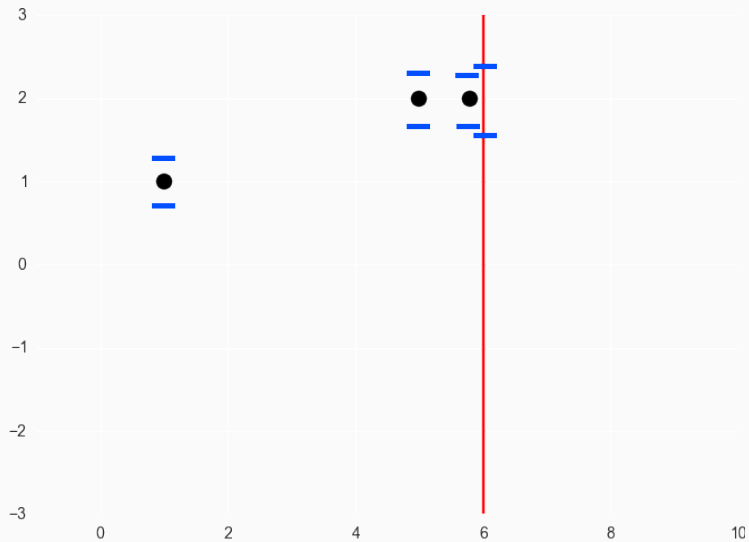
Functions



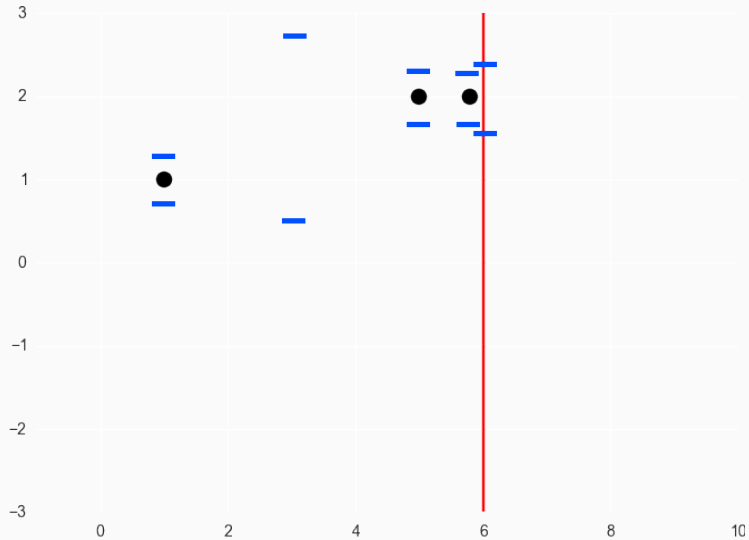
Functions



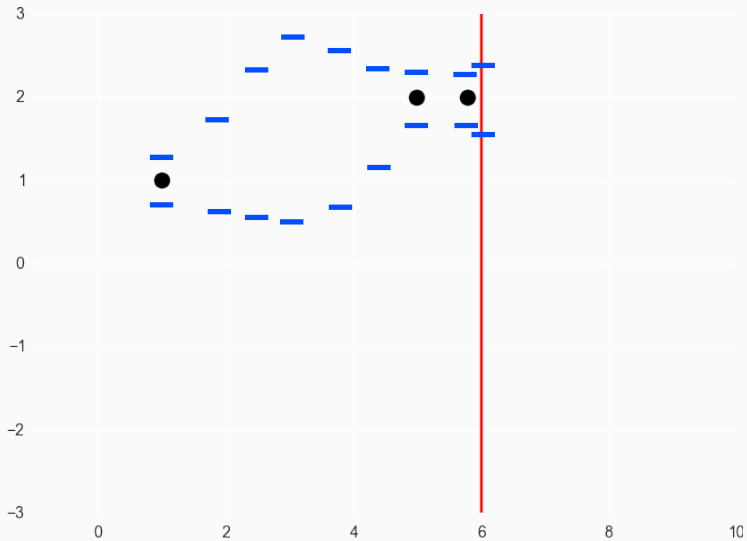
Functions



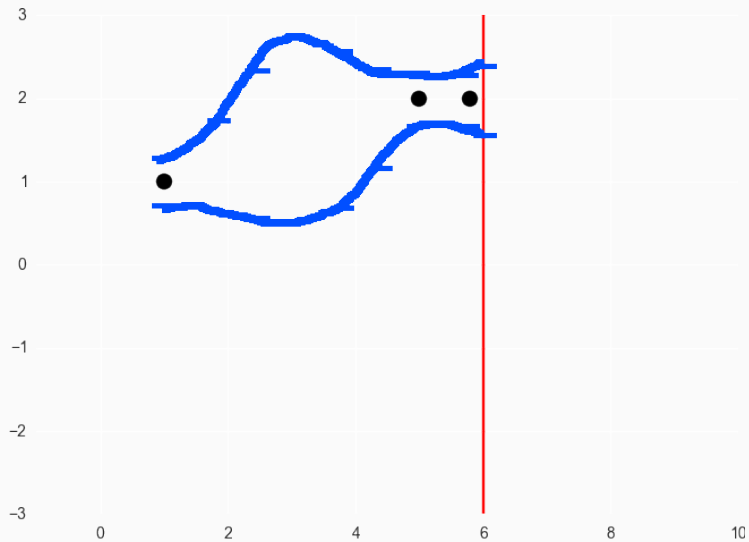
Functions



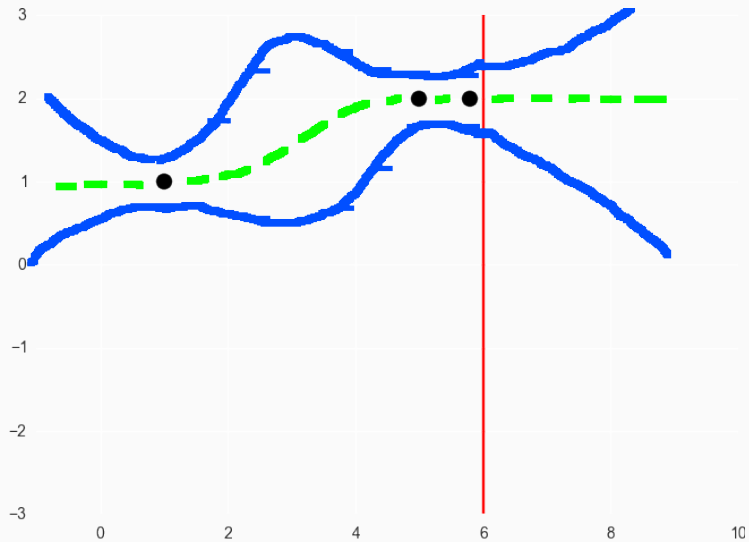
Functions



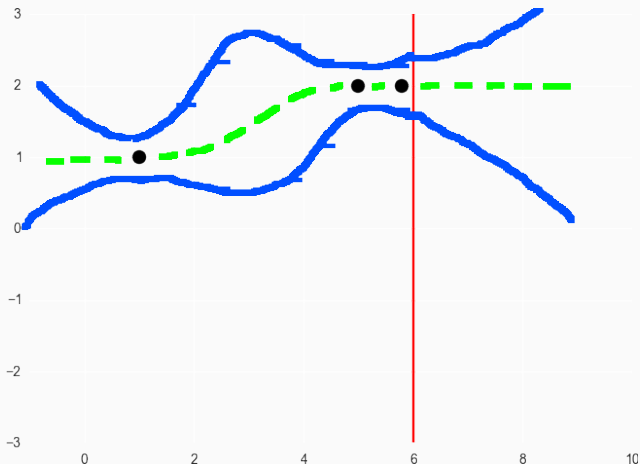
Functions



Functions



Functions



If all instantiations of the function are jointly Gaussian such that the co-variance structure depends on how much information an observation provides for the other, we will get the curve above.

Uncertainty over functions

- Regression model,

$$\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon$$
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

- Introduce f_i as *instantiation* of function,

$$f_i = f(\mathbf{x}_i),$$

- as a new random variable.

Uncertainty over functions

- Regression model,

$$\begin{aligned}y_i &= f(\mathbf{x}_i) + \epsilon \\ \epsilon &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})\end{aligned}$$

- Introduce f_i as *instantiation* of function,

$$f_i = f(\mathbf{x}_i),$$

- as a new random variable.
- now we have a "handle" to specify our assumptions over

Uncertainty over functions

Model,

$$p(\mathbf{y}, \mathbf{f}, \mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x})p(\boldsymbol{\theta})$$

Want to "push" \mathbf{x} through a mapping f of which we are uncertain,

$$p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta}),$$

prior over instantiations of function.

Uncertainty over functions

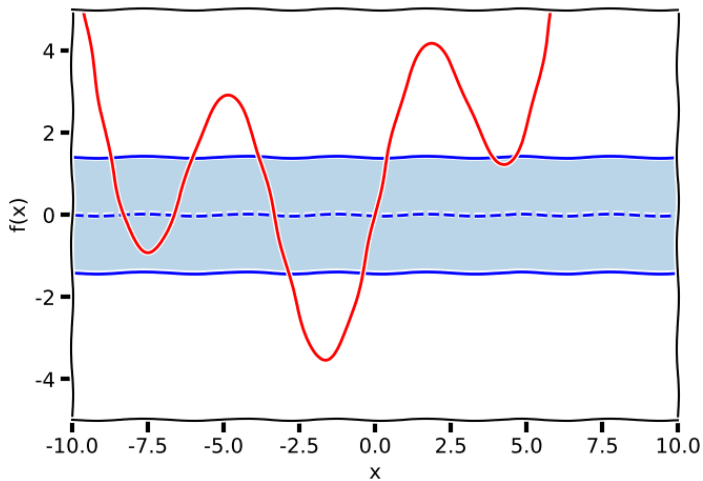
- As everything is gaussian both the marginal and predictive posterior are analytically tractable
- Marginal

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

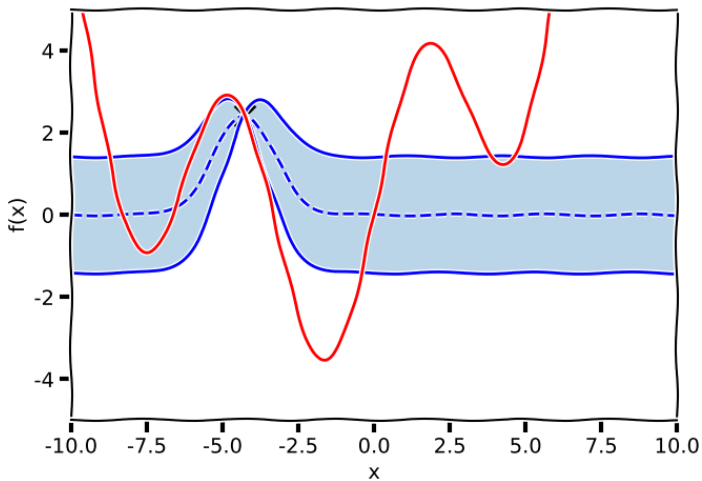
- Predictive posterior

$$p(\mathbf{f}_*|\mathbf{x}, \mathbf{x}_*, \mathbf{f}) = \frac{p(\mathbf{f}, \mathbf{f}_*|\mathbf{x}, \mathbf{x}_*)}{p(\mathbf{f}|\mathbf{x})}$$

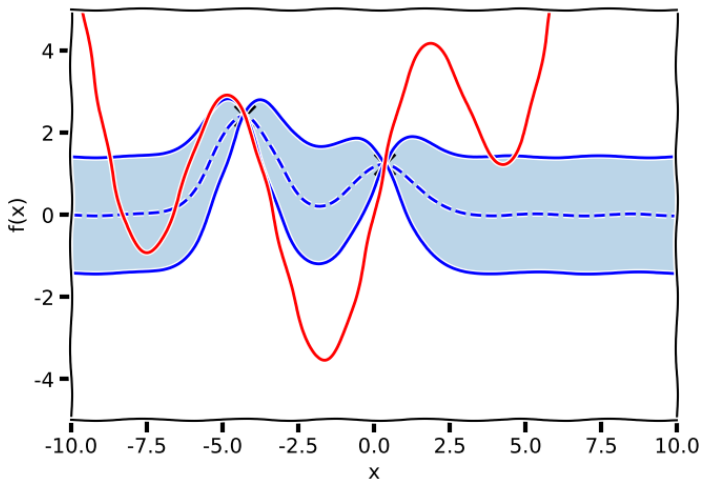
Gaussian Processes



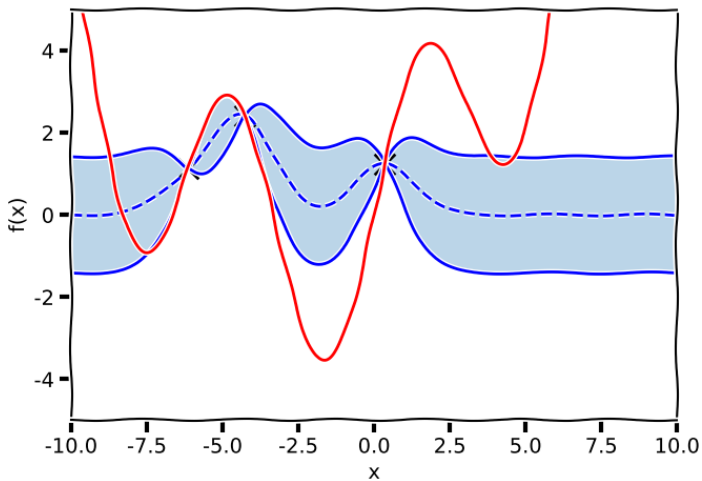
Gaussian Processes



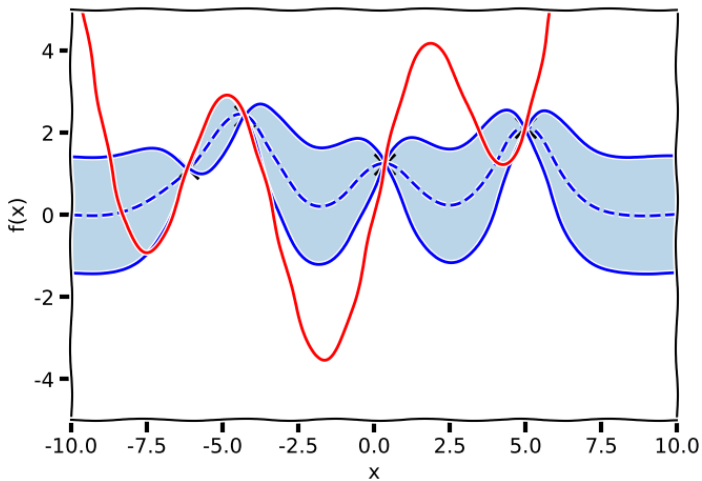
Gaussian Processes



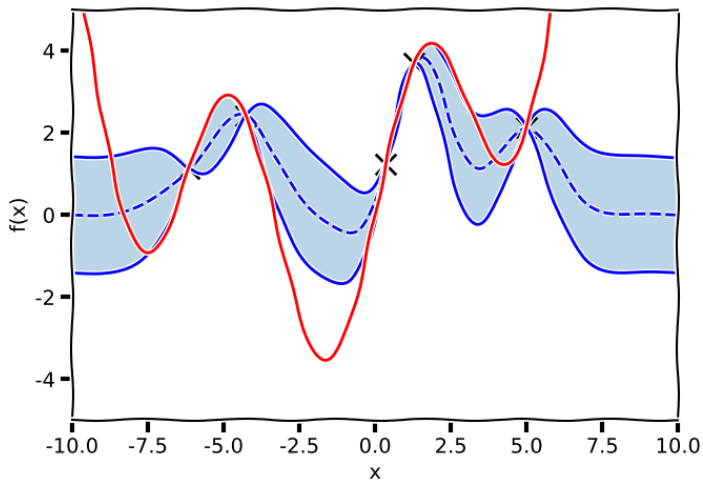
Gaussian Processes



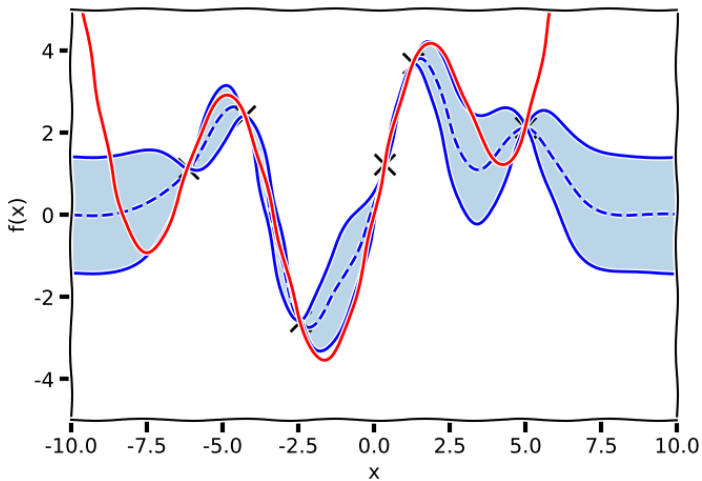
Gaussian Processes



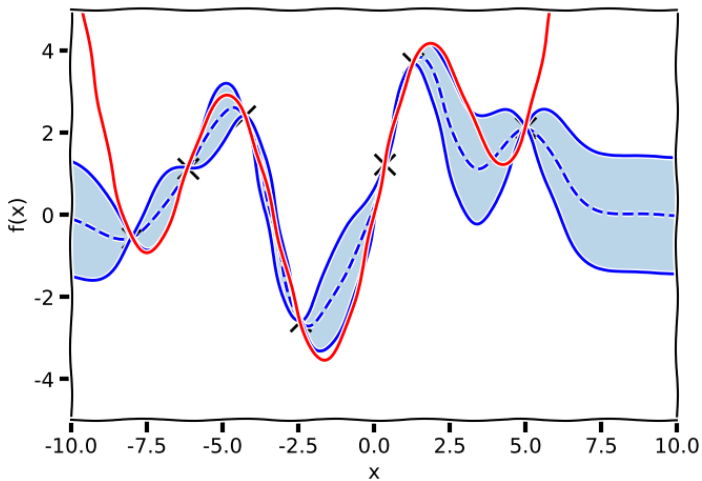
Gaussian Processes



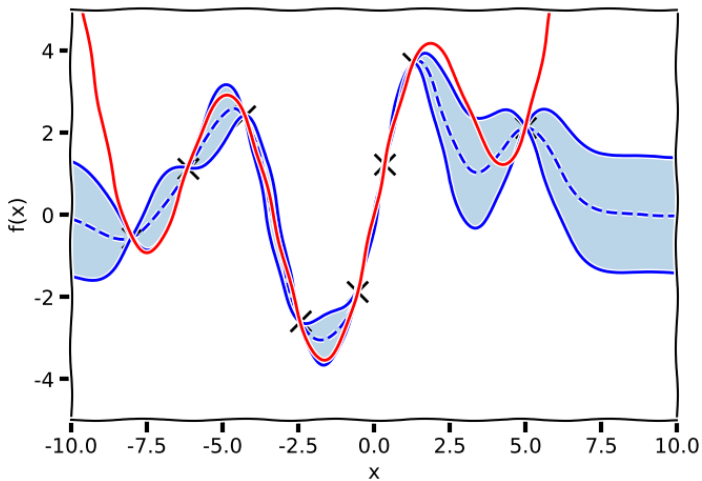
Gaussian Processes



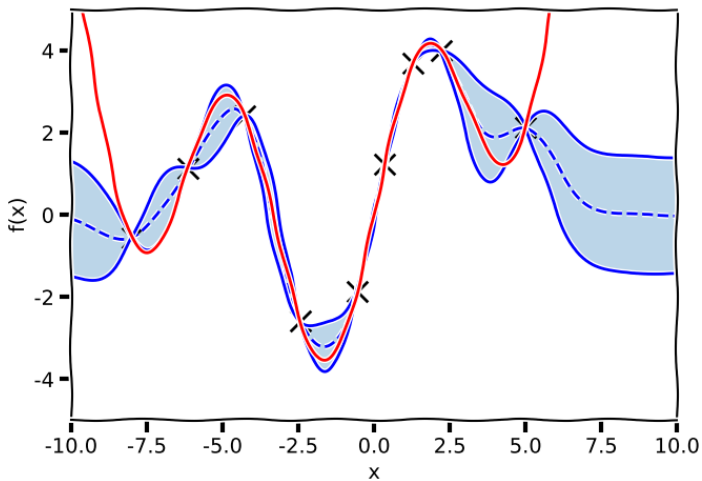
Gaussian Processes



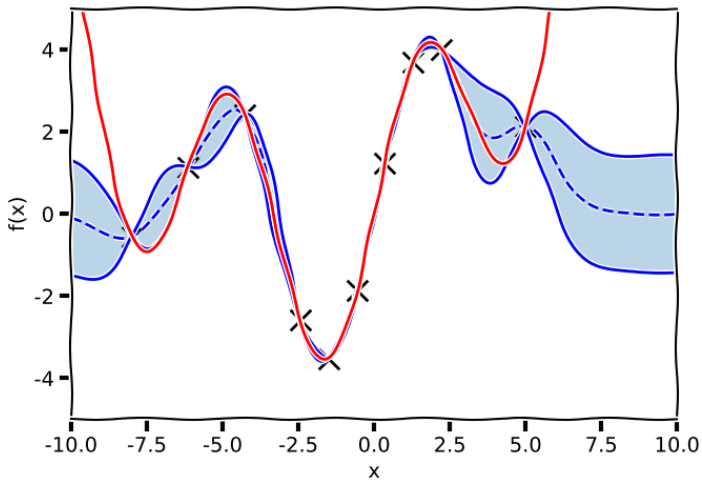
Gaussian Processes



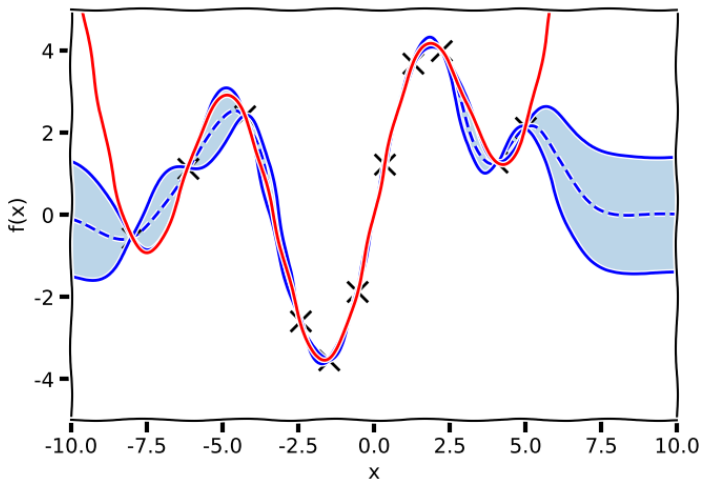
Gaussian Processes



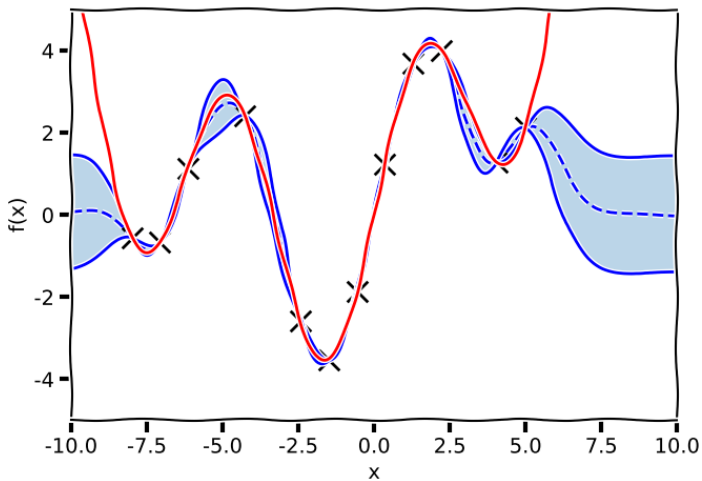
Gaussian Processes



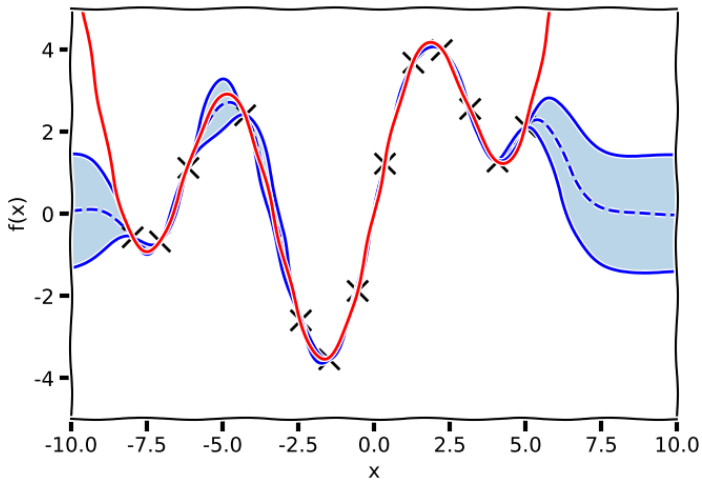
Gaussian Processes



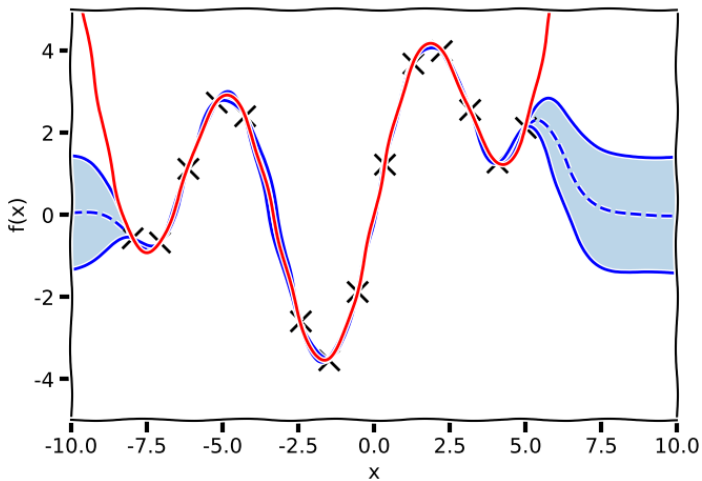
Gaussian Processes



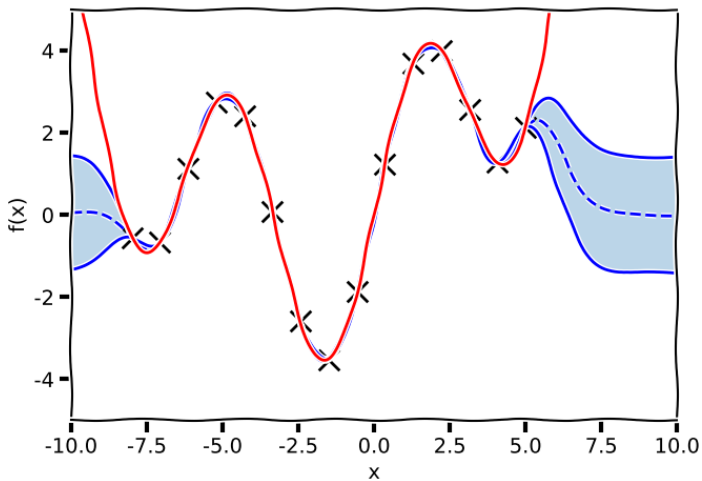
Gaussian Processes



Gaussian Processes



Gaussian Processes



Process \rightarrow Distribution \rightarrow value

- Each evaluation of a process is a distribution

$$\mathcal{N}(0, \Sigma) \sim \mathcal{N}(0, k(\mathbf{X}, \mathbf{X}))$$

Process \rightarrow Distribution \rightarrow value

- Each evaluation of a process is a distribution

$$\mathcal{N}(0, \Sigma) \sim \mathcal{N}(0, k(\mathbf{X}, \mathbf{X}))$$

- Each evaluation of a distribution is a value

$$y \sim \mathcal{N}(y|0, \Sigma)$$

Process \rightarrow Distribution \rightarrow value

- Each evaluation of a process is a distribution

$$\mathcal{N}(0, \Sigma) \sim \mathcal{N}(0, k(\mathbf{X}, \mathbf{X}))$$

- Each evaluation of a distribution is a value

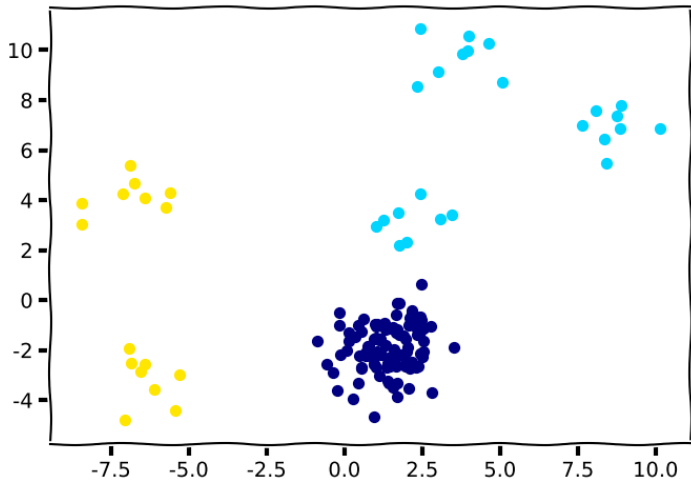
$$y \sim \mathcal{N}(y|0, \Sigma)$$

- Kolmogorov's Existence Theorem defines which distributions have an infinite generalisation

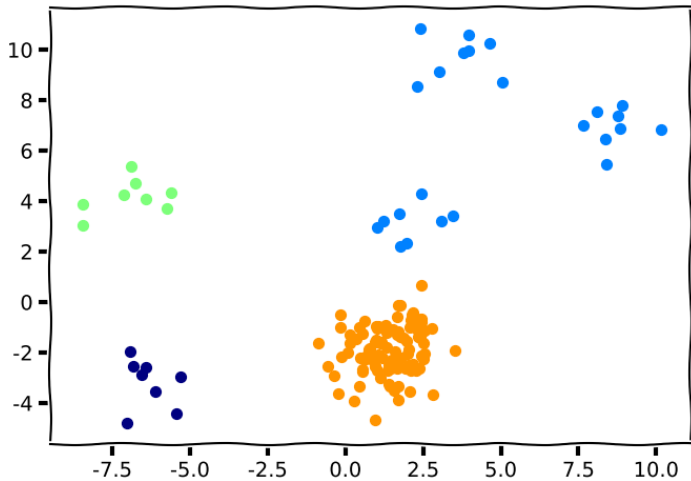
- Formulate process
- Evaluate process at specific location $x \rightarrow$ distribution
- Evaluate distribution at any location y
- GP is defined over uncountable infinite space

- Formulate process
- Evaluate process at specific location $x \rightarrow$ distribution
- Evaluate distribution at any location y
- GP is defined over uncountable infinite space
- *What about countable objects?*

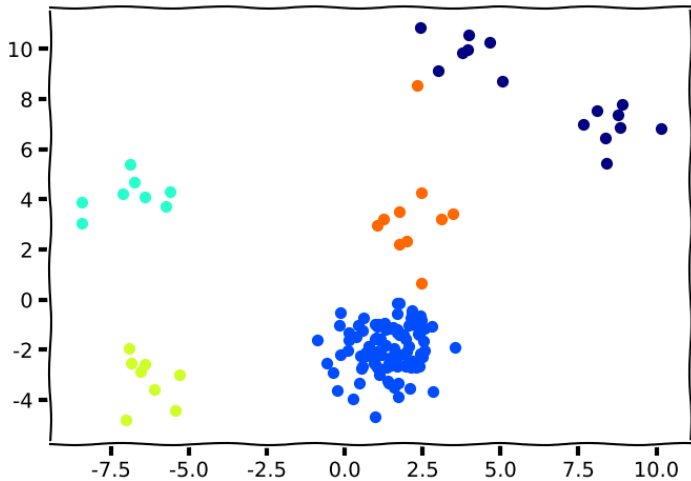
Gaussian Mixture Model



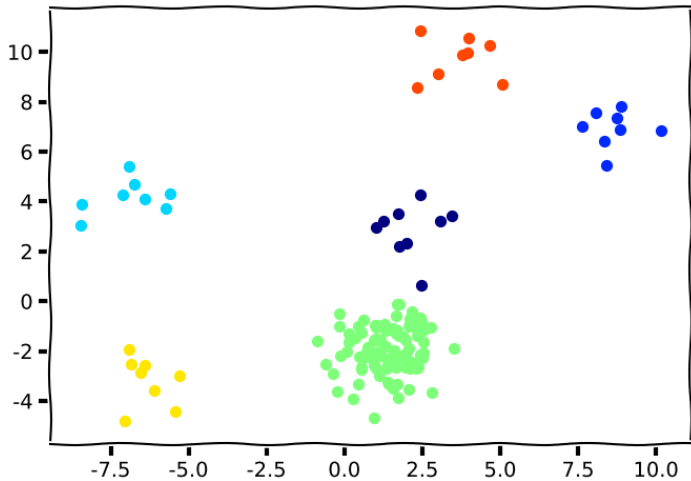
Gaussian Mixture Model



Gaussian Mixture Model



Gaussian Mixture Model

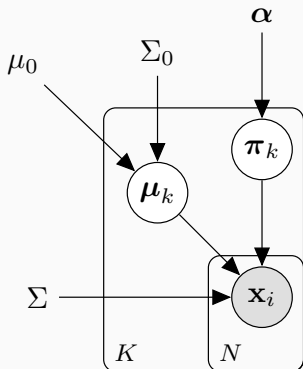


Gaussian Mixture Models

$$p(\mathbf{X}) = \sum_{k=1}^K p(\mathbf{X}|k)p(k) = \sum_{k=1}^K \mathcal{N}(\mathbf{X}|\boldsymbol{\mu}_k, \Sigma_k)p(k)$$

- Represent the probability of \mathbf{X} as a combination or *mixture* of distributions
- What should K be?
- Can we make K infinite?

Gaussian Mixture Model



1. Sample proportions
2. Sample cluster id given proportions
3. Sample cluster mean
4. Sample data

$$p(\mathbf{X}) = \sum_{k=1}^{\infty} p(\mathbf{X}|k)p(k) = \sum_{k=1}^{\infty} \mathcal{N}(\mathbf{X}|\mu_k, \Sigma_k)p(k)$$

- Multinomial

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

Distributions over partitionings

- Multinomial

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

- Conjugate prior

$$p(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k-1}$$

Distributions over partitionings

- Multinomial

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

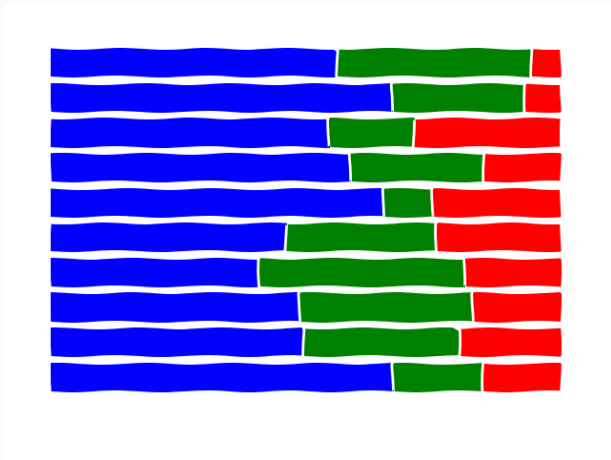
- Conjugate prior

$$p(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k-1}$$

- Dirichlet Distribution

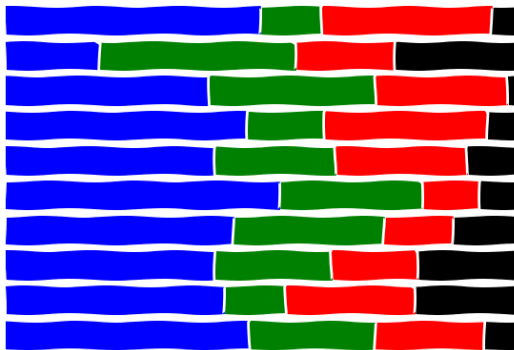
$$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdot \dots \cdot \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k-1}$$

Dirichlet Distribution



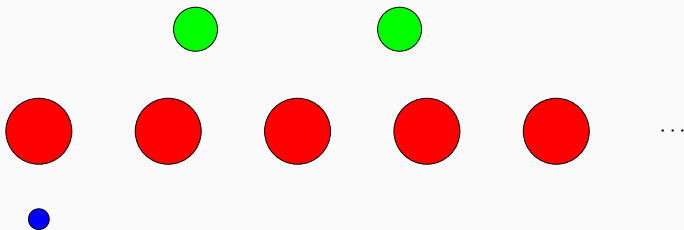
$\text{Dir}(10, 5, 3)$

Dirichlet Distribution

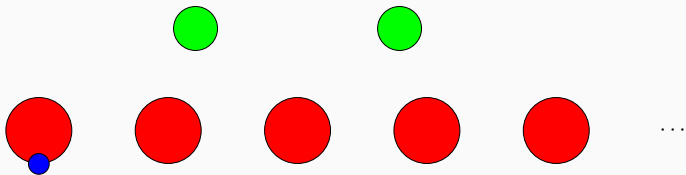


$\text{Dir}(7, 5, 3, 2)$

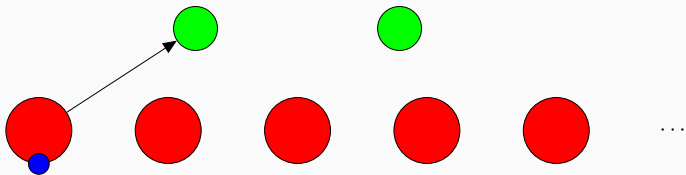
Chinese Restaurant Process



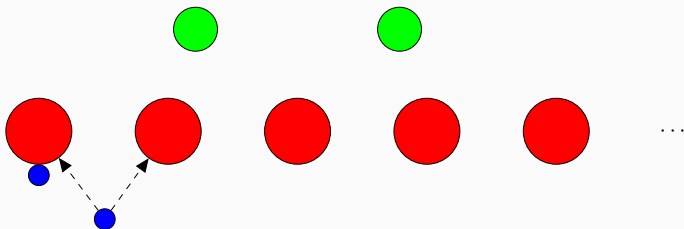
Chinese Restaurant Process



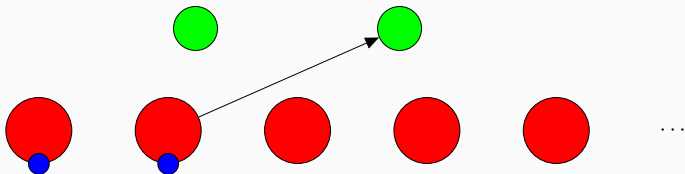
Chinese Restaurant Process



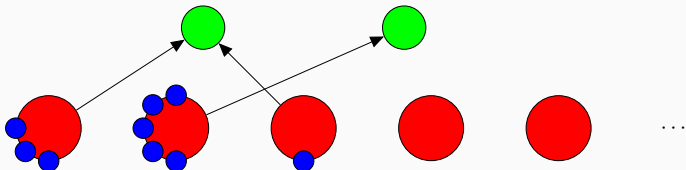
Chinese Restaurant Process



Chinese Restaurant Process

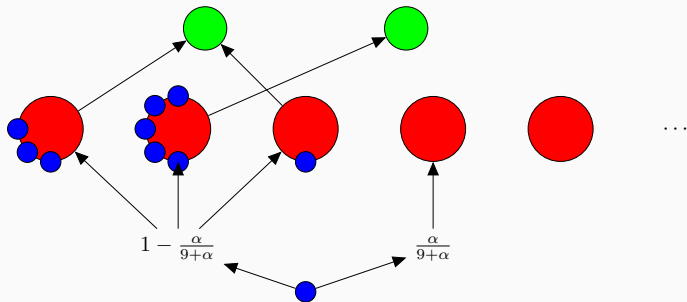


Chinese Restaurant Process

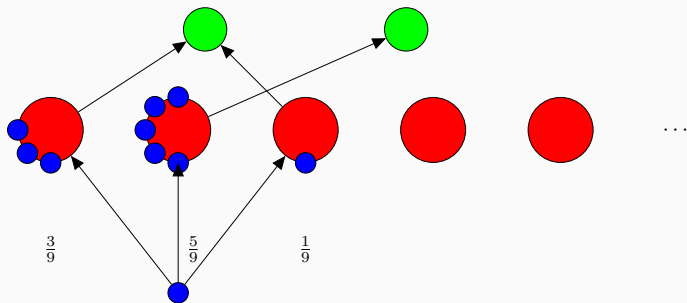


- Go to new table $\frac{\alpha}{N-1+\alpha}$
- If not choose table as $\frac{n_i}{N}$ where n_i number of diners at table i

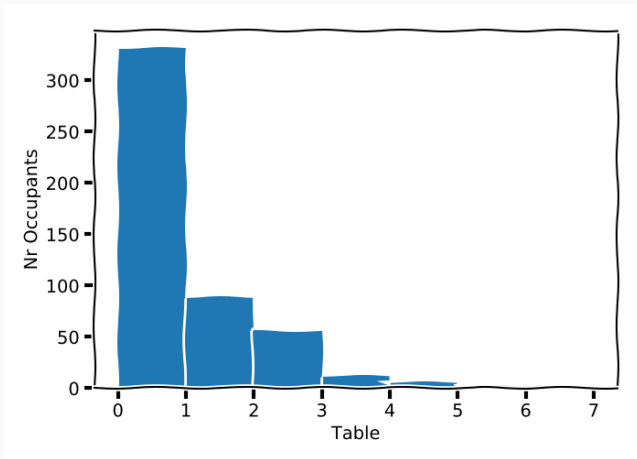
Chinese Restaurant Process



Chinese Restaurant Process

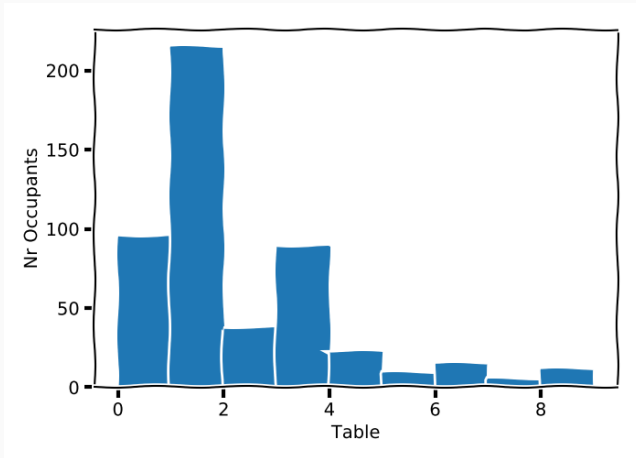


Chinese Restaurant Process



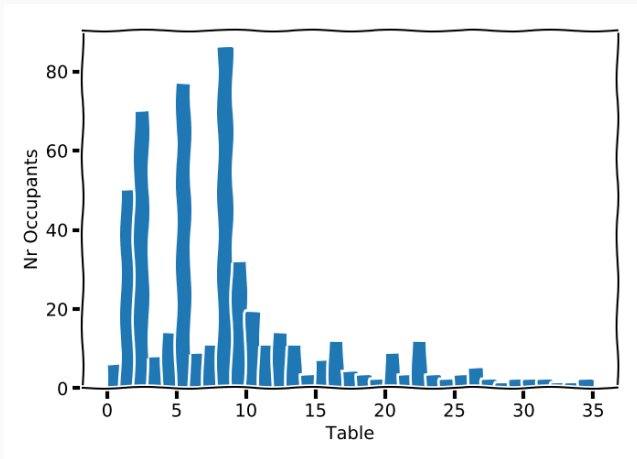
$$N = 500 \quad \alpha = 1.0$$

Chinese Restaurant Process



$$N = 500 \quad \alpha = 2.0$$

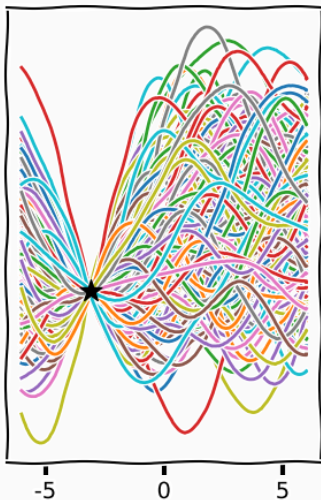
Chinese Restaurant Process



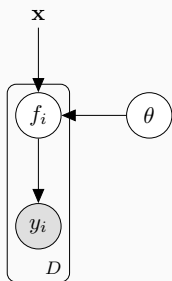
$N = 500$ $\alpha = 10.0$

Unsupervised Learning

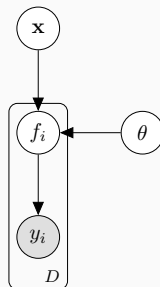
Gaussian Processes



Unsupervised Learning

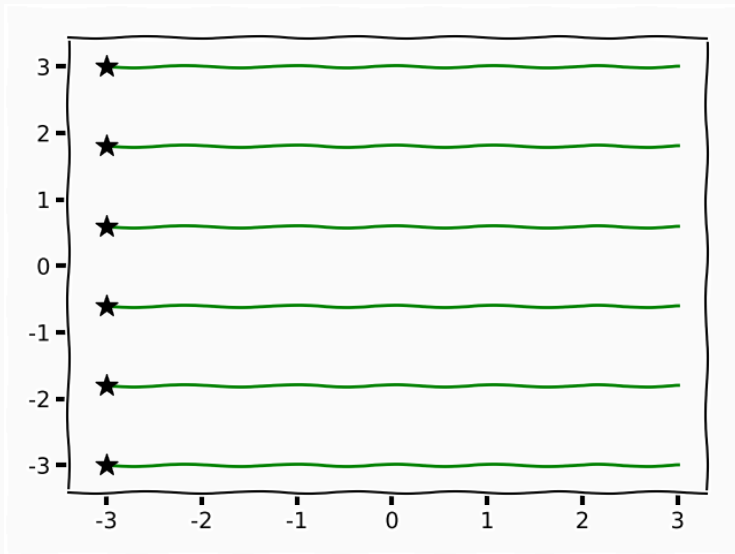


$$p(y|x)$$

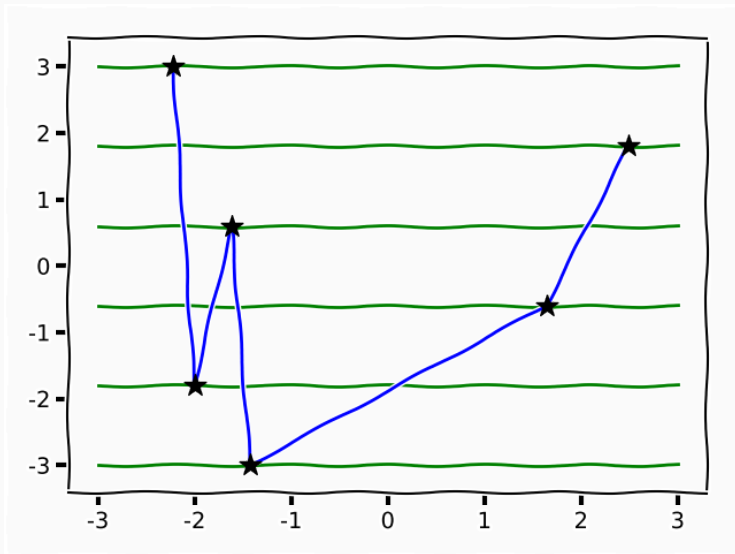


$$p(y)$$

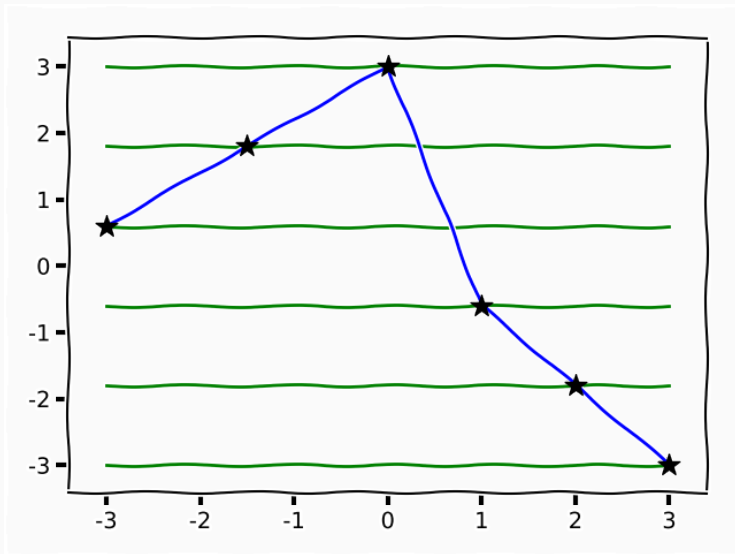
Unsupervised Learning



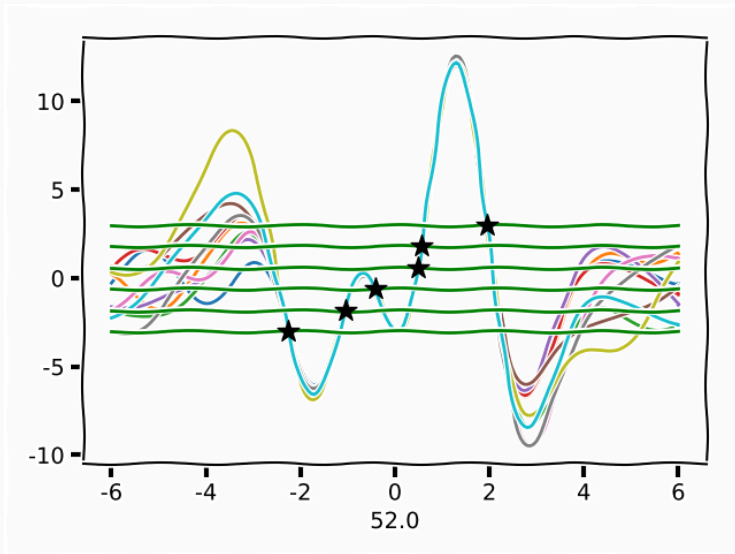
Unsupervised Learning



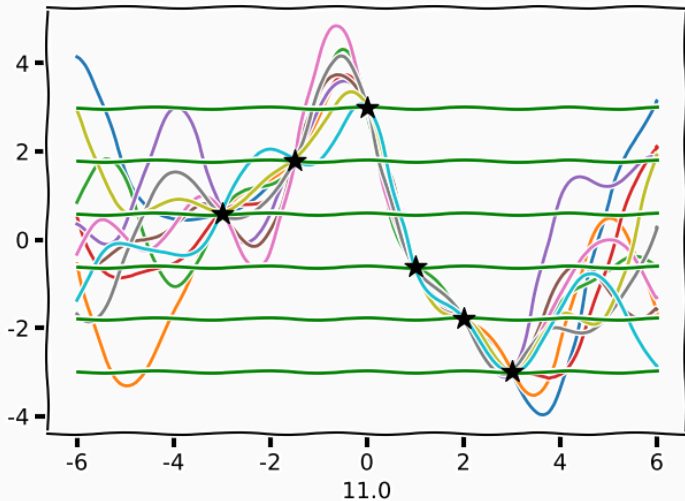
Unsupervised Learning



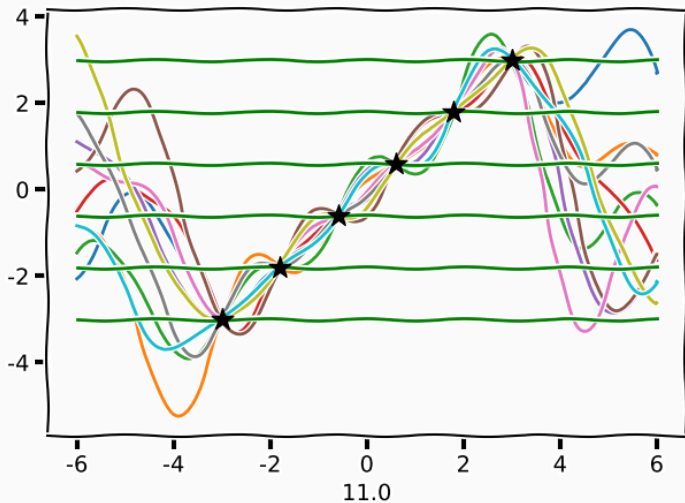
Unsupervised Learning



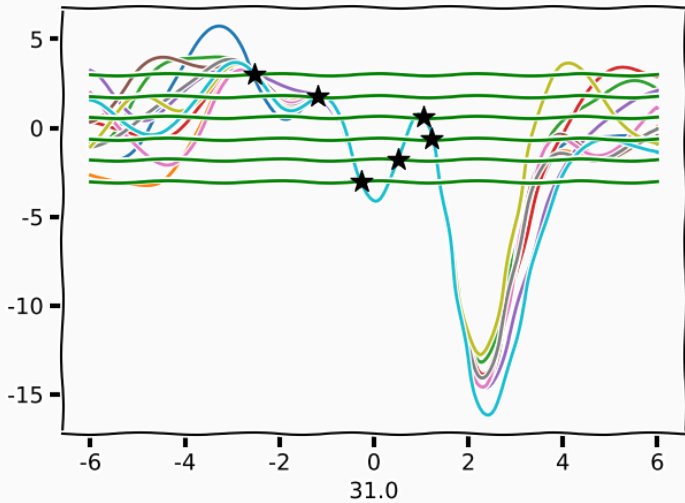
Unsupervised Learning



Unsupervised Learning



Unsupervised Learning



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

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

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. The priors are *"balanced"*

3. 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

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

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

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

- Analytically intractable (Non Elementary Integral) and infinitely differentiable

Laplace Integration



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

Unsupervised Learning with GPs

$$p(\mathbf{Y})$$

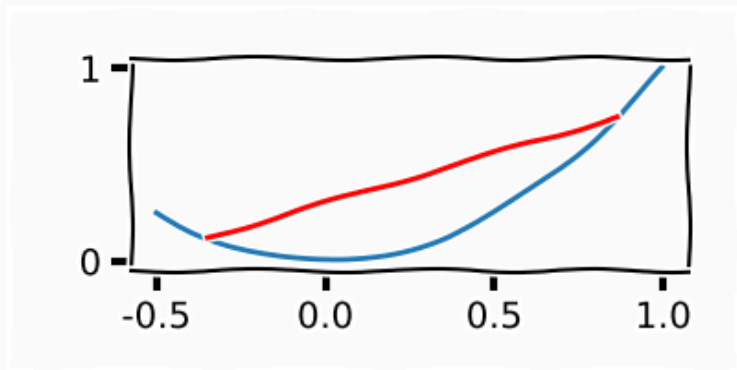
$$\log p(\mathbf{Y})$$

$$\log p(\mathbf{Y}) = \log \int p(\mathbf{Y}, \mathbf{X}) d\mathbf{X}$$

$$\log p(\mathbf{Y}) = \log \int p(\mathbf{Y}, \mathbf{X}) d\mathbf{X} = \log \int p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X}$$

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int p(\mathbf{Y}, \mathbf{X}) d\mathbf{X} = \log \int p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X} \\ &= \log \int \frac{q(\mathbf{X})}{q(\mathbf{X})} p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X}\end{aligned}$$

Jensen Inequality



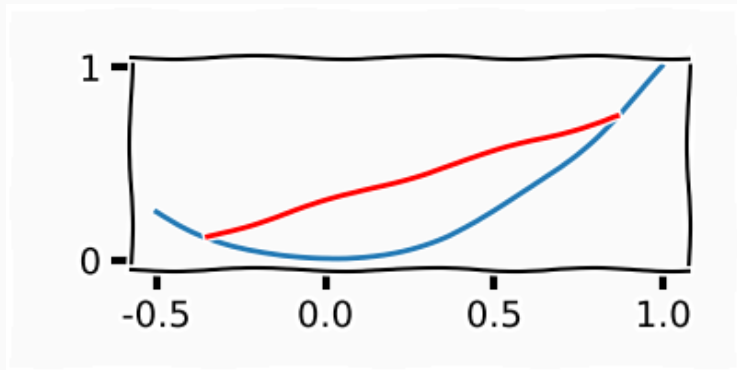
Convex Function

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

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

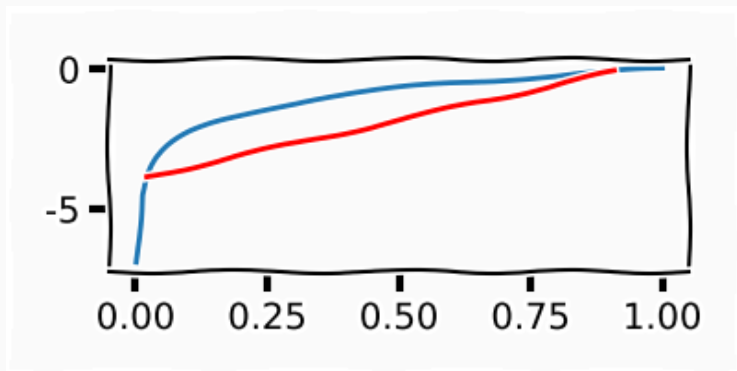
$$\lambda \in [0, 1]$$

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

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

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int \frac{q(\mathbf{X})}{q(\mathbf{X})} p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X} = \\ &\geq \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X}\end{aligned}$$

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int \frac{q(\mathbf{X})}{q(\mathbf{X})} p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X} = \\ &\geq \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} \\ &= \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} + \int q(\mathbf{X}) d\mathbf{X} \log p(\mathbf{Y})\end{aligned}$$

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int \frac{q(\mathbf{X})}{q(\mathbf{X})} p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X} = \\ &\geq \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} \\ &= \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} + \int q(\mathbf{X}) d\mathbf{X} \log p(\mathbf{Y}) \\ &= -\text{KL}(q(\mathbf{X}) || p(\mathbf{X}|\mathbf{Y})) + \log p(\mathbf{Y})\end{aligned}$$

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int \frac{q(\mathbf{X})}{q(\mathbf{X})} p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X} = \\ &\geq \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} \\ &= \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} + \int q(\mathbf{X}) d\mathbf{X} \log p(\mathbf{Y}) \\ &= -\text{KL}(q(\mathbf{X}) || p(\mathbf{X}|\mathbf{Y})) + \log p(\mathbf{Y})\end{aligned}$$

- if $q(\mathbf{X})$ is the true posterior we have an equality, therefore match the distributions

Variational Bayes cont.

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int \frac{q(\mathbf{X})}{q(\mathbf{X})} p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y}) d\mathbf{X} = \\ &\geq \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y}) p(\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} \\ &= \int q(\mathbf{X}) \log \frac{p(\mathbf{X}|\mathbf{Y})}{q(\mathbf{X})} d\mathbf{X} + \int q(\mathbf{X}) d\mathbf{X} \log p(\mathbf{Y}) \\ &= -\text{KL}(q(\mathbf{X}) || p(\mathbf{X}|\mathbf{Y})) + \log p(\mathbf{Y})\end{aligned}$$

- if $q(\mathbf{X})$ is the true posterior we have an equality, therefore match the distributions
- i.e. $\text{argmin}_q \text{KL}(q(\mathbf{X}) || p(\mathbf{X}|\mathbf{Y}))$
 \Rightarrow variational distributions are approximations to intractable posteriors

$$\text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y}))$$

$$\text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})) = \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X}|\mathbf{Y})} d\mathbf{X}$$

$$\begin{aligned}\text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})) &= \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X}|\mathbf{Y})} d\mathbf{X} \\ &= \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X}, \mathbf{Y})} d\mathbf{X} + \log p(\mathbf{Y})\end{aligned}$$

$$\begin{aligned}\text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})) &= \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X}|\mathbf{Y})} d\mathbf{X} \\ &= \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X}, \mathbf{Y})} d\mathbf{X} + \log p(\mathbf{Y}) \\ &= H(q(\mathbf{X})) - \mathbb{E}_{q(\mathbf{X})} [\log p(\mathbf{X}, \mathbf{Y})] + \log p(\mathbf{Y})\end{aligned}$$

$$\log p(\mathbf{Y}) = \text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})) + \underbrace{\mathbb{E}_{q(\mathbf{x})} [\log p(\mathbf{X}, \mathbf{Y})] - H(q(\mathbf{X}))}_{\text{ELBO}}$$

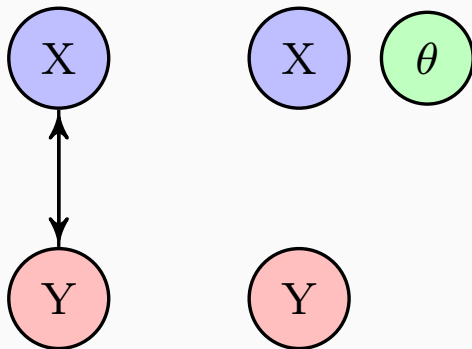
$$\log p(\mathbf{Y}) = \text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})) + \underbrace{\mathbb{E}_{q(\mathbf{X})} [\log p(\mathbf{X}, \mathbf{Y})] - H(q(\mathbf{X}))}_{\text{ELBO}}$$

$$\geq \mathbb{E}_{q(\mathbf{X})} [\log p(\mathbf{X}, \mathbf{Y})] - H(q(\mathbf{X})) = \mathcal{L}(q(\mathbf{X}))$$

$$\log p(\mathbf{Y}) = \text{KL}(q(\mathbf{X})||p(\mathbf{X}|\mathbf{Y})) + \underbrace{\mathbb{E}_{q(\mathbf{x})} [\log p(\mathbf{X}, \mathbf{Y})] - H(q(\mathbf{X}))}_{\text{ELBO}}$$

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

- if we maximise the ELBO we,
 - find an approximate posterior
 - get an approximation to the marginal likelihood
- *maximising* $p(\mathbf{Y})$ is learning
- finding $p(\mathbf{X}|\mathbf{Y}) \approx q(\mathbf{X})$ is prediction



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 Adams³

³Talking Machines Season 2, Episode 5

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 Adams³

³Talking Machines Season 2, Episode 5

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 Adams³

³Talking Machines Season 2, Episode 5

$$\mathcal{L} = \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{X})}{q(\mathbf{X})} \right)$$

⁴Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

$$\mathcal{L} = \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{X})}{q(\mathbf{X})} \right) \\ \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{X})p(\mathbf{X})}{q(\mathbf{X})} \right)$$

⁴Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

$$\begin{aligned}\mathcal{L} &= \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{X})}{q(\mathbf{X})} \right) \\ &\quad \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{X})p(\mathbf{X})}{q(\mathbf{X})} \right) \\ &= \int_{\mathbf{F}, \mathbf{X}} q(\mathbf{X}) \log p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{X}) - \int_{\mathbf{X}} q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X})}\end{aligned}$$

⁴Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

$$\begin{aligned}\mathcal{L} &= \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{X})}{q(\mathbf{X})} \right) \\ &= \int_{\mathbf{X}, \mathbf{F}} q(\mathbf{X}) \log \left(\frac{p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{X})p(\mathbf{X})}{q(\mathbf{X})} \right) \\ &= \int_{\mathbf{F}, \mathbf{X}} q(\mathbf{X}) \log p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{X}) - \int_{\mathbf{X}} q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X})} \\ &= \tilde{\mathcal{L}} - \text{KL} (q(\mathbf{X}) \parallel p(\mathbf{X}))\end{aligned}$$

⁴Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

$$\tilde{\mathcal{L}} = \int_{\mathbf{F}, \mathbf{X}} q(\mathbf{X}) \log p(\mathbf{Y}|\mathbf{F}) p(\mathbf{F}|\mathbf{X})$$

- Has not eliviate the problem at all, X still needs to go through F to reach the data
- Idea of sparse approximations⁵

⁵Quinonero-Candela, Joaquin, & Rasmussen, C. E. (2005). A unifying view of sparse approximate Gaussian process regression & Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs

- Add another set of samples from the same prior

$$p(\mathbf{U}|\mathbf{Z}) = \prod_{j=1}^d \mathcal{N}(\mathbf{u}_{:,j}|\mathbf{0}, \mathbf{K})$$

- Add another set of samples from the same prior

$$p(\mathbf{U}|\mathbf{Z}) = \prod_{j=1}^d \mathcal{N}(\mathbf{u}_{:,j}|\mathbf{0}, \mathbf{K})$$

- Conditional distribution

$$\begin{aligned} p(\mathbf{f}_{:,j}, \mathbf{u}_{:,j}|\mathbf{X}, \mathbf{Z}) &= p(\mathbf{f}_{:,j}|\mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z})p(\mathbf{u}_{:,j}|\mathbf{Z}) \\ &= \mathcal{N}(\mathbf{f}_{:,j}|\mathbf{K}_{fu}(\mathbf{K}_{uu})^{-1}\mathbf{u}_{:,j}, \mathbf{K}_{ff} - \mathbf{K}_{fu}(\mathbf{K}_{uu})^{-1}\mathbf{K}_{uf}) \mathcal{N}(\mathbf{u}_{:,j}|\mathbf{0}, \mathbf{K}_{uu}), \end{aligned}$$

$$p(\mathbf{Y}, \mathbf{F}, \mathbf{U}, \mathbf{X} | \mathbf{Z}) = p(\mathbf{X}) \prod_{j=1}^d p(\mathbf{y}_{:,j} | \mathbf{f}_{:,j}) p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}) p(\mathbf{u}_{:,j} | \mathbf{Z})$$

- we have done nothing to the model, just added *halucinated* observations

$$p(\mathbf{Y}, \mathbf{F}, \mathbf{U}, \mathbf{X} | \mathbf{Z}) = p(\mathbf{X}) \prod_{j=1}^d p(\mathbf{y}_{:,j} | \mathbf{f}_{:,j}) p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}) p(\mathbf{u}_{:,j} | \mathbf{Z})$$

- we have done nothing to the model, just added *halucinated* observations
- however, we will now interpret \mathbf{U} and \mathbf{X}_u **not** as random variables but **variational** parameters

$$p(\mathbf{Y}, \mathbf{F}, \mathbf{U}, \mathbf{X} | \mathbf{Z}) = p(\mathbf{X}) \prod_{j=1}^d p(\mathbf{y}_{:,j} | \mathbf{f}_{:,j}) p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}) p(\mathbf{u}_{:,j} | \mathbf{Z})$$

- we have done nothing to the model, just added *halucinated* observations
- however, we will now interpret \mathbf{U} and \mathbf{X}_u **not** as random variables but **variational** parameters
- i.e. parametrise approximate posterior using these parameters (remember sparse motivation)

- Variational distributions are approximations to intractable posteriors,

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

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

$$q(\mathbf{X}) \approx p(\mathbf{X}|\mathbf{Y})$$

Lower Bound

- Variational distributions are approximations to intractable posteriors,

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

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

$$q(\mathbf{X}) \approx p(\mathbf{X}|\mathbf{Y})$$

- Assume that we can *find* \mathbf{U} that completely represents \mathbf{F} , i.e. \mathbf{U} is sufficient statistics of \mathbf{F} ,

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

$$\tilde{\mathcal{L}} = \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) \log \frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{U}|\mathbf{X}, \mathbf{Z})}{q(\mathbf{F})q(\mathbf{U})}$$

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) \log \frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{U}|\mathbf{X}, \mathbf{Z})}{q(\mathbf{F})q(\mathbf{U})} \\ &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) \log \frac{\prod_{j=1}^d p(\mathbf{y}_{:,j}|\mathbf{f}_{:,j})p(\mathbf{f}_{:,j}|\mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z})p(\mathbf{u}_{:,j}|\mathbf{Z})}{q(\mathbf{F})q(\mathbf{U})}\end{aligned}$$

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) \log \frac{p(\mathbf{Y}, \mathbf{F}, \mathbf{U}|\mathbf{X}, \mathbf{Z})}{q(\mathbf{F})q(\mathbf{U})} \\ &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) \log \frac{\prod_{j=1}^d p(\mathbf{y}_{:,j}|\mathbf{f}_{:,j})p(\mathbf{f}_{:,j}|\mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z})p(\mathbf{u}_{:,j}|\mathbf{Z})}{q(\mathbf{F})q(\mathbf{U})}\end{aligned}$$

- Assume that \mathbf{U} is sufficient statistics for \mathbf{F}

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

$$\begin{aligned}\tilde{\mathcal{L}} &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} \prod_{j=1}^d p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_{:,j}) q(\mathbf{X}) \\ \log \frac{\prod_{j=1}^d p(\mathbf{y}_{:,j} | \mathbf{f}_{:,j}) p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_{:,j} | \mathbf{Z})}{\prod_{j=1}^d p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_{:,j})} &= \end{aligned}$$

$$\begin{aligned}
 \tilde{\mathcal{L}} &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} \prod_{j=1}^d p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_{:,j}) q(\mathbf{X}) \\
 &\quad \log \frac{\prod_{j=1}^d p(\mathbf{y}_{:,j} | \mathbf{f}_{:,j}) p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_{:,j} | \mathbf{Z})}{\prod_{j=1}^d p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_{:,j})} = \\
 &= \int_{\mathbf{X}, \mathbf{F}, \mathbf{U}} \prod_{j=1}^p p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_{:,j}) q(\mathbf{X}) \log \frac{\prod_{j=1}^p p(\mathbf{y}_{:,j} | \mathbf{f}_{:,j}) p(\mathbf{u}_{:,j} | \mathbf{Z})}{\prod_{j=1}^p q(\mathbf{u}_{:,j})} \\
 &= \mathbb{E}_{q(\mathbf{F}), q(\mathbf{X}), q(\mathbf{U})} [p(\mathbf{Y} | \mathbf{F})] - \text{KL}(q(\mathbf{U}) || p(\mathbf{U} | \mathbf{Z}))
 \end{aligned}$$

$$\mathbb{E}_{q(\mathbf{F}), q(\mathbf{X}), q(\mathbf{U})} [p(\mathbf{Y}|\mathbf{F})] - \text{KL}(q(\mathbf{U})||p(\mathbf{U}|\mathbf{Z})) - \text{KL}(q(\mathbf{X})||p(\mathbf{X}))$$

- Expectation tractable (for some co-variances)
- Reduces to expectations over co-variance functions known as Ψ statistics
- Allows us to place priors and not "regularisers" over the latent representation

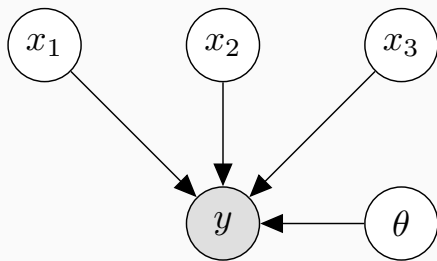
Latent space priors

$$\mathbb{E}_{q(\mathbf{F}), q(\mathbf{X}), q(\mathbf{U})} [p(\mathbf{Y}|\mathbf{F})] - \text{KL}(q(\mathbf{U})||p(\mathbf{U}|\mathbf{Z})) - \text{KL}(q(\mathbf{X})||p(\mathbf{X}))$$

- Importantly $p(\mathbf{X})$ appears only in KL term
- Allows us to express stronger assumptions about the model

⁶Damianou, A. C., Titsias, M., & Lawrence, Neil D, Variational Inference for Uncertainty on the Inputs of Gaussian Process Models (2014)

Factor Analysis

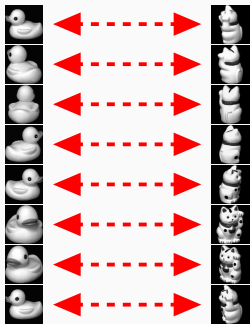


$$y = f(x_1, x_2, x_3) + \epsilon$$

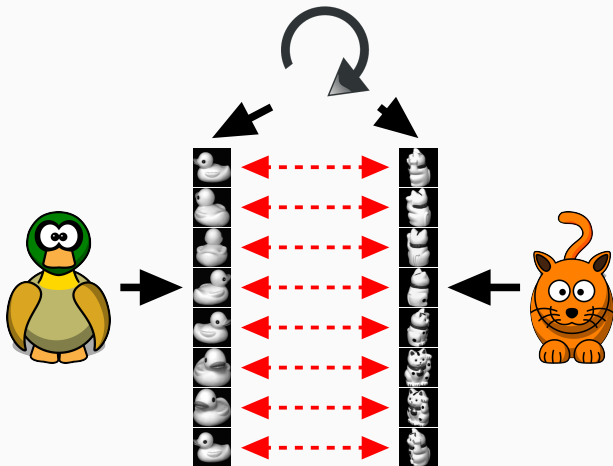
Alignments



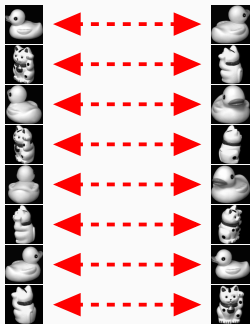
Alignments



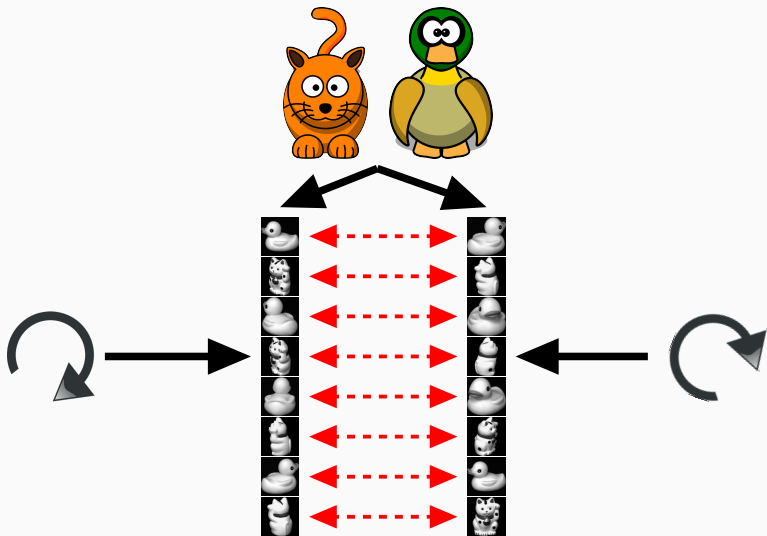
Alignments



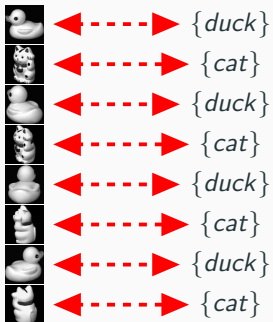
Alignments



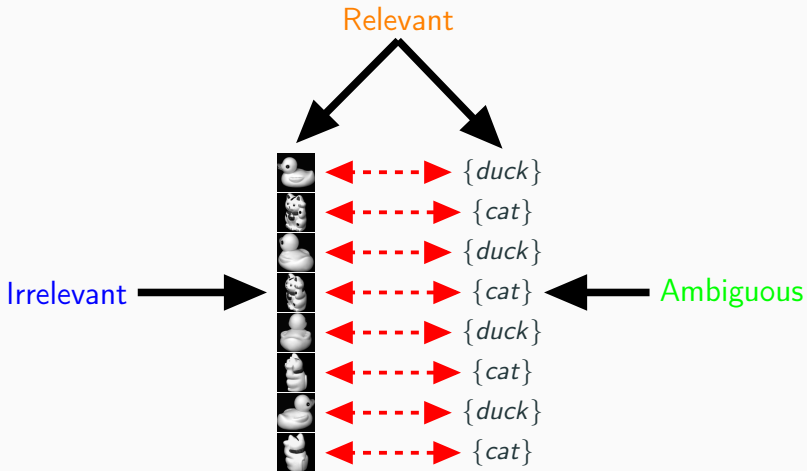
Alignments



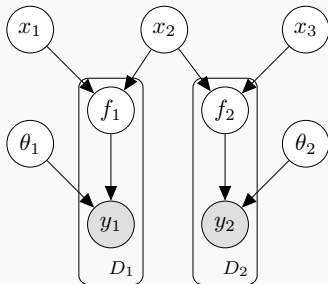
Alignments

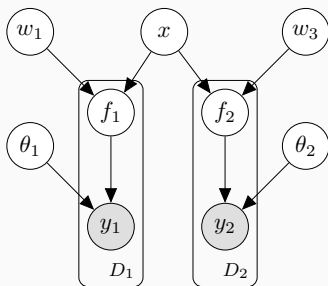


Alignments



Explaining Away cont.

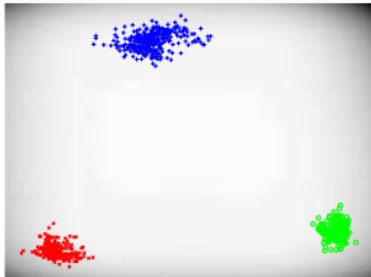
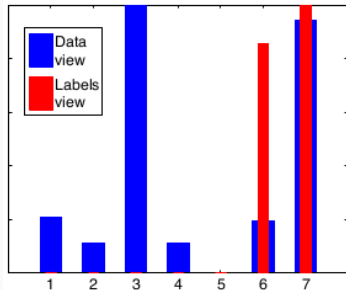


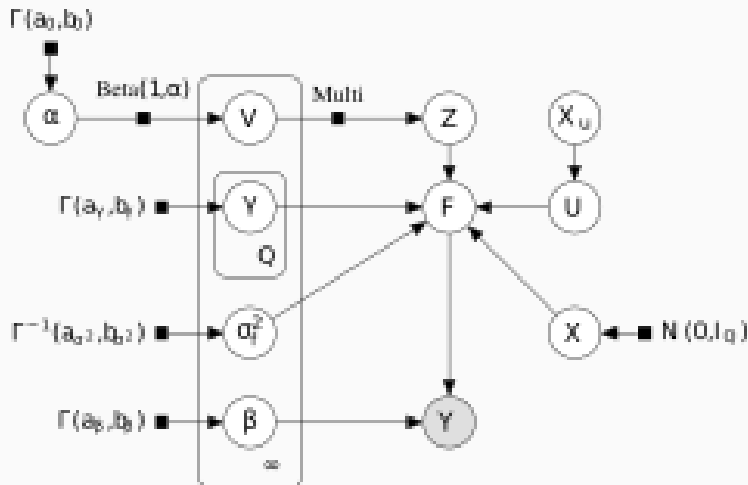


$$y_1 = f(w_1^T x) \quad y_2 = f(w_2^T x)$$

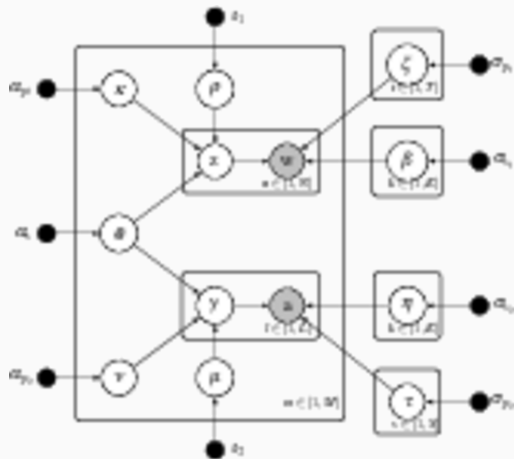
⁷Damianou, A., Lawrence, N. D., & Ek, C. H. (2016). Multi-view learning as a nonparametric nonlinear inter-battery factor analysis

IBFA with GP-LVM

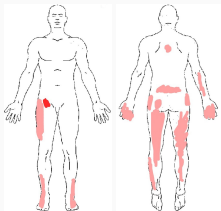




⁸Joint work with Andrew Lawrence and Neill Campbell at University of Bath, Will be presented at *Advances in Modeling and Learning Interactions from Complex Data* NIPS 2017



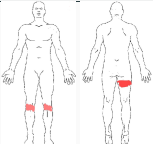
⁹Zhang, C., Kjellström, Hedvig, & Ek, C. H., Inter-battery topic representation learning, In ECCV 2016



Symptom diagnoses: Interscapular discomfort; R arm discomfort; B hands discomfort; Lumbago; B crest of the ilium discomfort; L side thigh discomfort; B back thigh discomfort; B calf discomfort; B achilles tendinitis; B shin discomfort; R inguinal discomfort;

Pattern diagnoses B L5 Radiculopathy; B S1 Radiculopathy; B C7 Radiculopathy;

Pathophysiological diagnoses DLI L4-L5; DLI S1-S2; DLI C6-C7

	<p>6 Prd: R back thigh dcf; L PFS (Patellofemoral pain syndrome); R PFS; L L5 Rdc; R L5 Rdc; DLI L4-L5;</p> <p>6 GT: R back thigh dcf; L PFS; R PFS; L L5 Rdc; R L5 Rdc; DLI L4-L5;</p>
-----------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------



50 Prd: Headache; L neck dcf; R neck dcf; Neck dcf; L upper trapezius dcf; R upper trapezius dcf; L shoulder dcf; L hand dcf; R hand dcf; Interscapular dcf; Lumbago; Lateral abdominal dcf; L groin dcf; L side thigh dcf; R side thigh dcf; L calf dcf; L back thigh dcf; L crest of the ilium dcf; R crest of the ilium dcf; R foot arch dcf; L toe joint dcf; R toe joint dcf; L medial elbow dcf; L ankle dcf; R ankle dcf; L foot arch dcf; L PFS; L dorsal knee dcf; R medial knee dcf;

L C2 Rdc; R C2 Rdc; L C3 Rdc; R C3 Rdc; L C4 Rdc; R C4 Rdc; R C6 Rdc; L C6 Rdc; L C7 Rdc; R C7 Rdc; L L5 Rdc; R L5 Rdc; L S1 Rdc; R S1 Rdc; DLI C2-C3; DLI C3-C4; DLI C5-C6; DLI C6-C7; DLI L4-L5; DLI L5-S1; OB;

50 GT: L back headache; R back headache; Neck dcf; L jaw dcf; L upper trapezius dcf; R upper trapezius dcf; L arm dcf; R arm dcf; L lateral elbow dcf; R lateral elbow dcf; L hand joint dcf; R hand joint dcf; L hand dcf; R hand dcf; L thumb dcf; R thumb dcf; L finger dcf; R finger dcf; Lumbago; L groin dcf; L back thigh dcf; L calf dcf; L medial knee dcf; L ankle dcf; R ankle dcf; R medial knee dcf; R big toe dcf; L big toe dcf;

L C2 Rdc; R C2 Rdc; L C3 Rdc; R C3 Rdc; L C4 Rdc; R C4 Rdc; L C5 Rdc; R C5 Rdc; L C6 Rdc; R C6 Rdc; L C7 Rdc; R C7 Rdc; L L4 Rdc; L L5 Rdc; R L5 Rdc; L S1 Rdc; R S1 Rdc; Craniocervical joint injury; DLI C4-C5; DLI L3-L4; DLI L4-L5; DLI L5-S1;

Convolutional Deep GPs

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

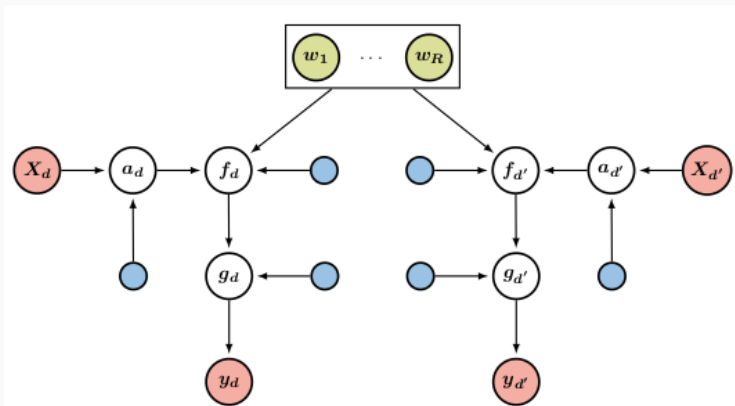
*"Scientific modelling is a scientific activity, the aim of which is to make a particular part or feature of the world easier to **understand**, define, quantify, visualize, or simulate by **referencing** it to existing and usually commonly **accepted** knowledge."* ¹⁰

¹⁰Wikipedia

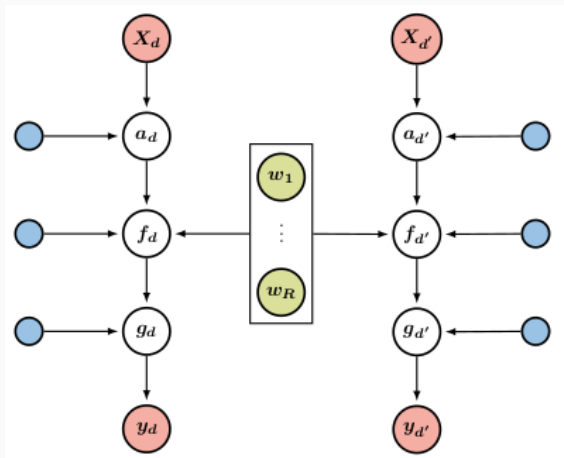
$$\mathbf{y}_d = g_d(f_d(a_d(\mathbf{x}))) + \epsilon$$

$$f_d(\mathbf{x}) = \sum_{r=1}^R \int T_{d,r}(\mathbf{x} - \mathbf{z}) w_r(\mathbf{z}) d\mathbf{z}$$

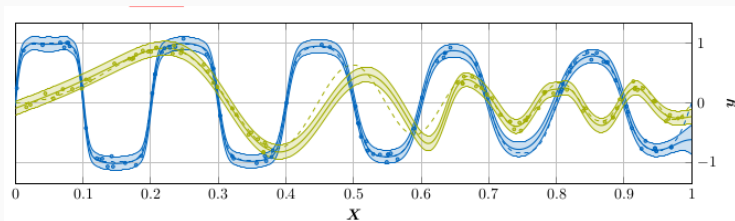
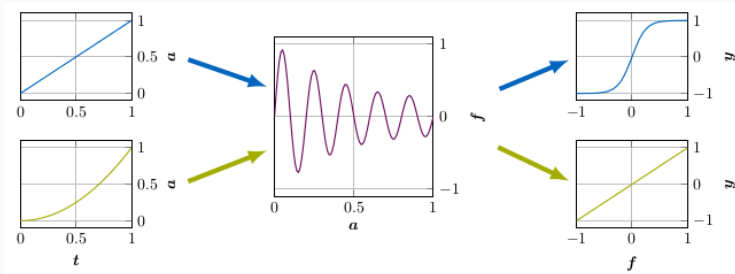
- Hierarchical set of function
- Convolution process with shared kernel



¹¹Kaiser, M., Otte, C., Runkler, T., & Ek, C. H. , Bayesian alignments of warped multi-output gaussian processes, ArXiv e-prints, (), (2017).



Hierarchies¹¹



Composite Functions

- Learning Composite functions have become very popular
- Composite functions are not as intuitive as one might think

Why are composite functions attractive?

$$y = g(\mathbf{x}) = f_K(f_{K-1}(f_{K-2}(\dots f_1(\mathbf{x}) \dots)))$$

- Kernel of a function

$$\text{Kern}(f_k) = \{(\mathbf{x}, \mathbf{x}') | f_k(\mathbf{x}) = f_k(\mathbf{x}')\}$$

- Kernel of a function

$$\text{Kern}(f_k) = \{(\mathbf{x}, \mathbf{x}') | f_k(\mathbf{x}) = f_k(\mathbf{x}')\}$$

- Image of a function

$$\text{Im}(f_k(\mathbf{x})) = \{\mathbf{y} \in Y | \mathbf{y} = f_k(\mathbf{x}), \mathbf{x} \in X\}$$

- Kernel of function

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

Composite Functions

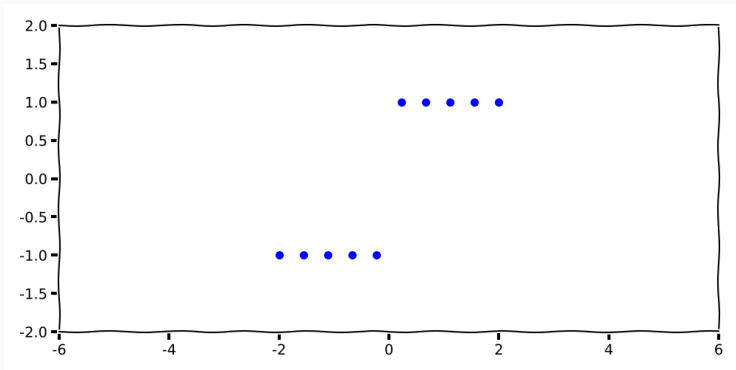
- Kernel of function

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

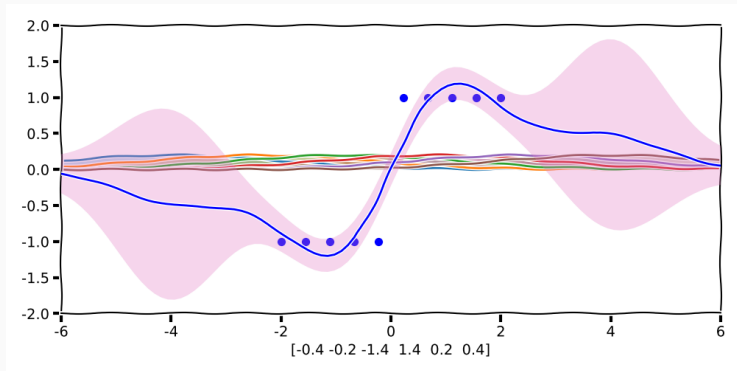
- Image of a function

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

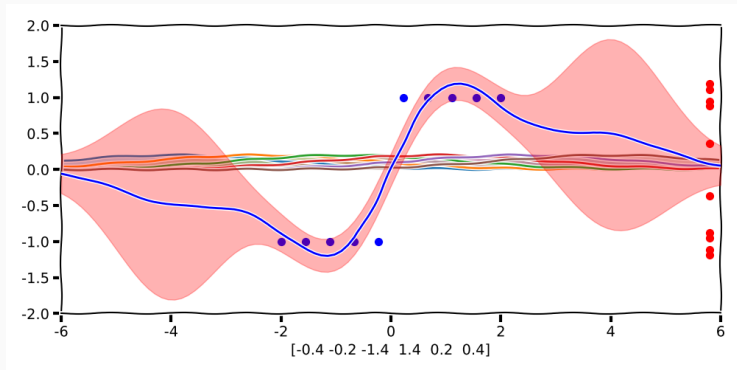
Composite Functions



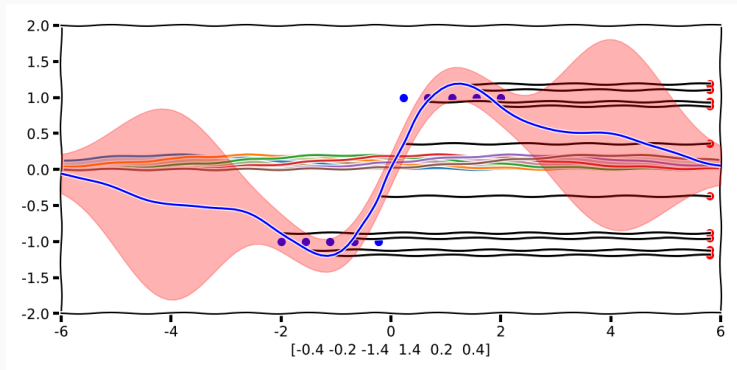
Composite Functions



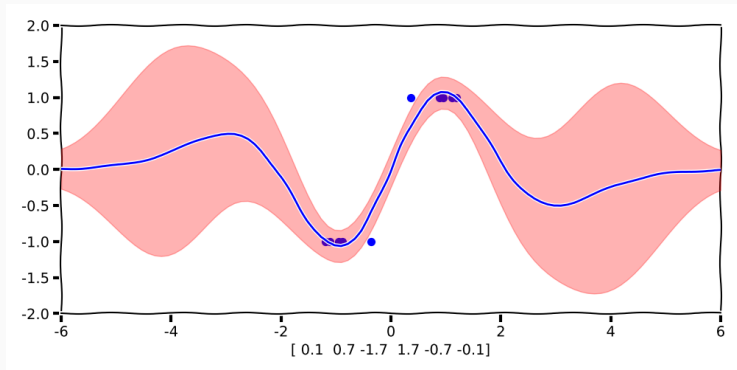
Composite Functions



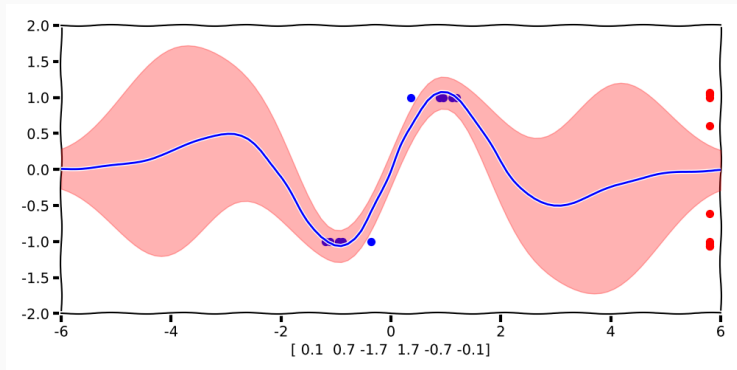
Composite Functions



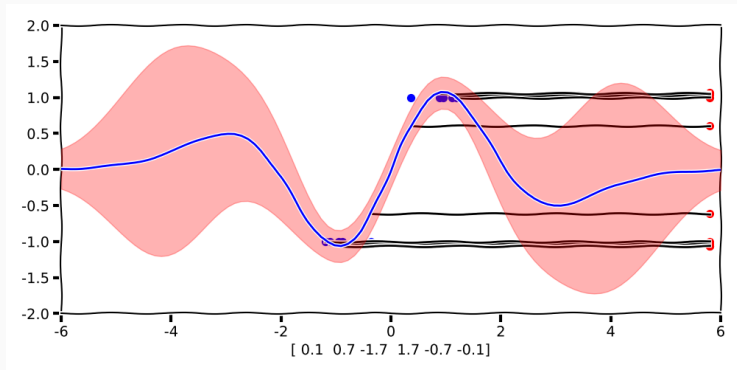
Composite Functions



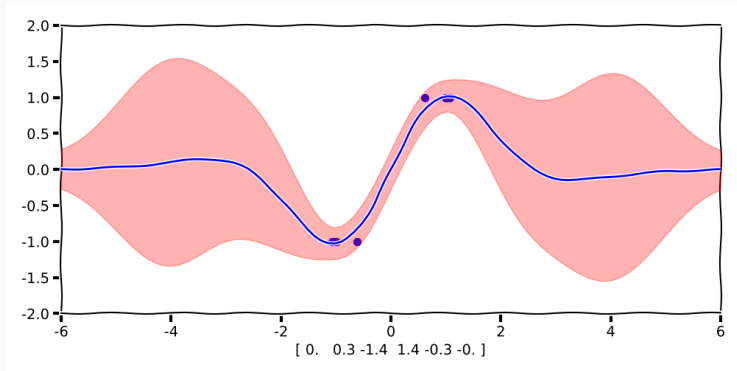
Composite Functions



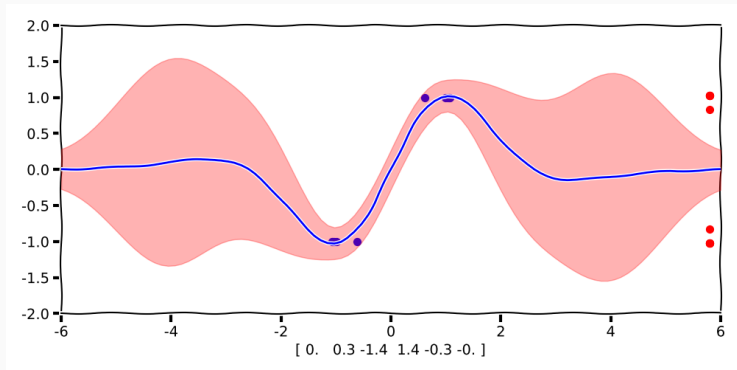
Composite Functions



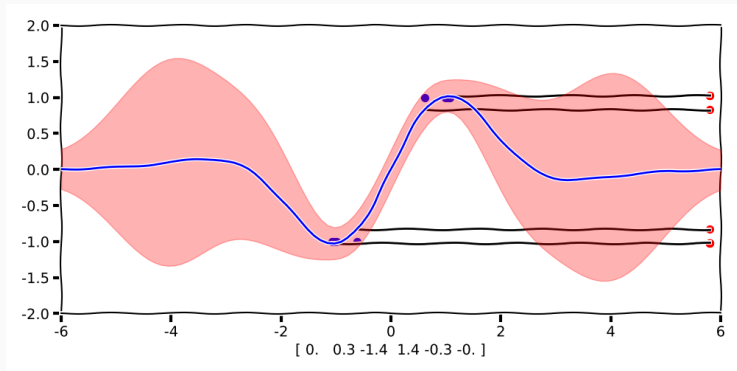
Composite Functions



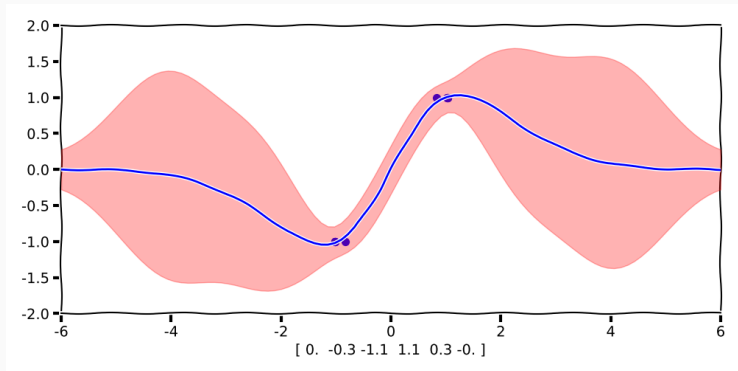
Composite Functions



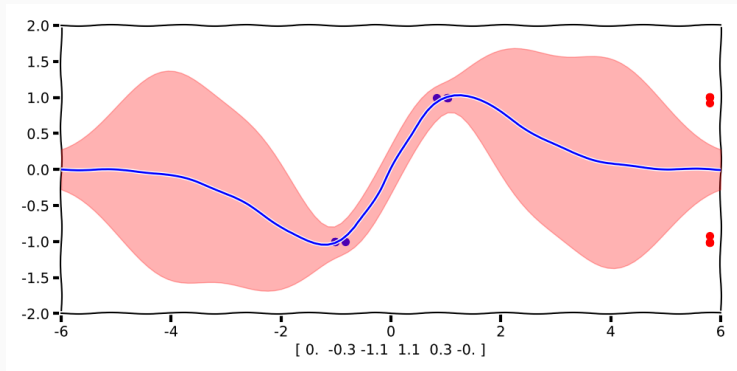
Composite Functions



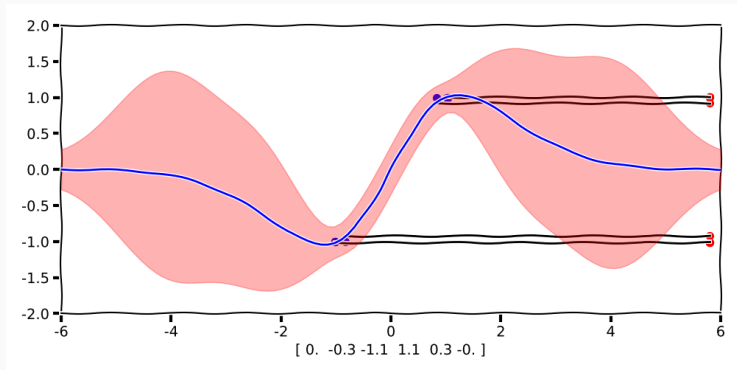
Composite Functions



Composite Functions



Composite Functions



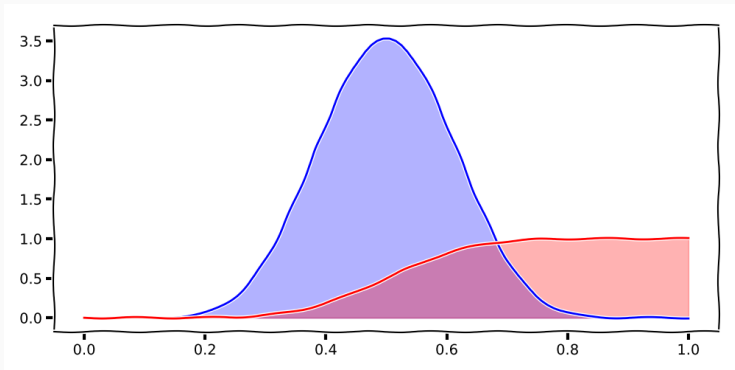
Theorem (Change of Variable)

Let $x \in \mathcal{X} \subseteq \mathbb{R}^n$ be a random vector with a probability density function given by $p_x(x)$, and let $y \in \mathcal{Y} \subseteq \mathbb{R}^n$ be a random vector such that $\psi(y) = x$, where the function $\psi : \mathcal{Y} \rightarrow \mathcal{X}$ is bijective of class of \mathcal{C}^1 and $|\nabla \psi(y)| > 0, \forall y \in \mathcal{Y}$. Then, the probability density function $p_y(\cdot)$ induced in \mathcal{Y} is given by

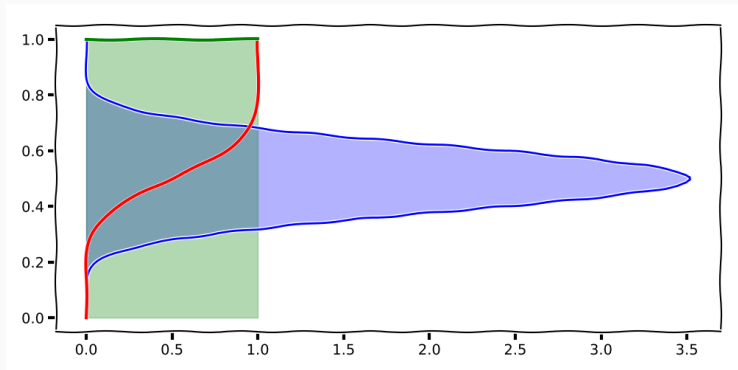
$$p_y(y) = p_x(\psi(y)) |\nabla \psi(y)|$$

where $\nabla \psi(\cdot)$ denotes the Jacobian of $\psi(\cdot)$, and $|\cdot|$ denotes the determinant operator.

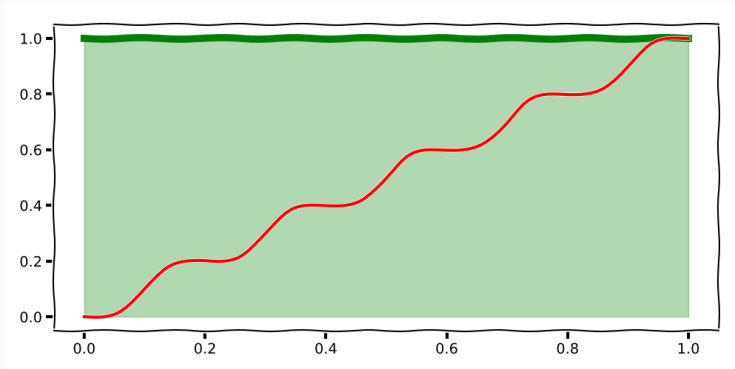
Sampling



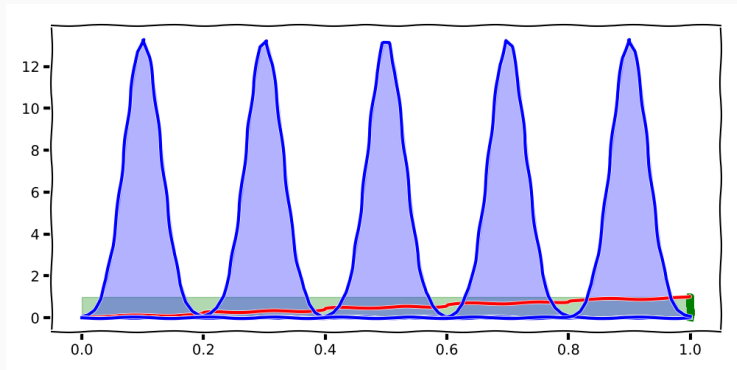
Sampling



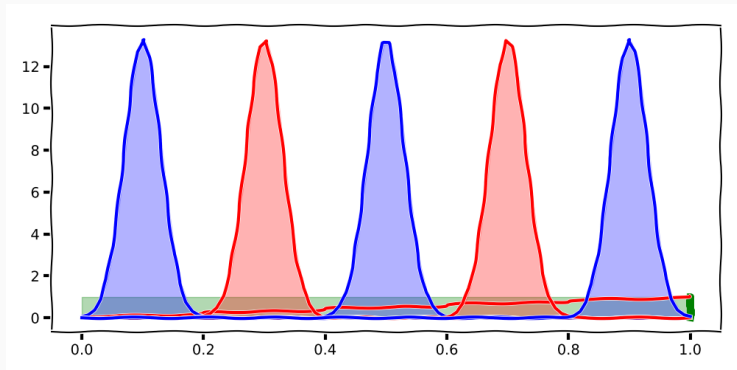
Change of Variables



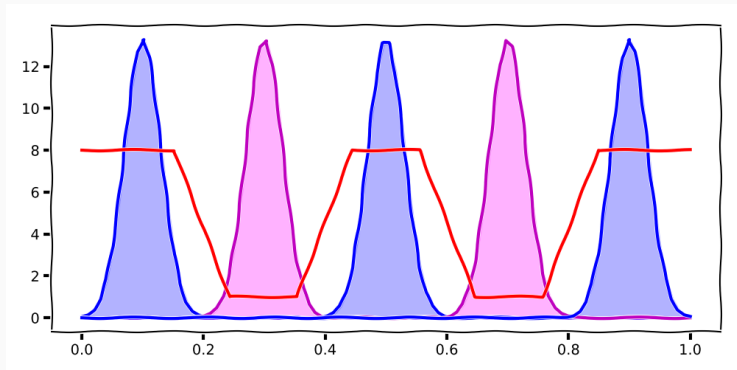
Change of Variables



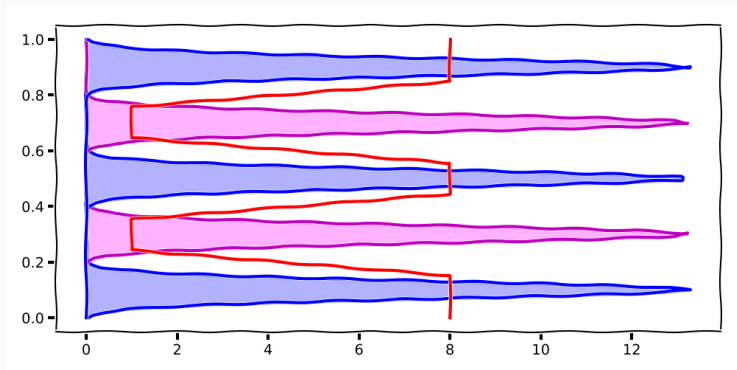
Change of Variables



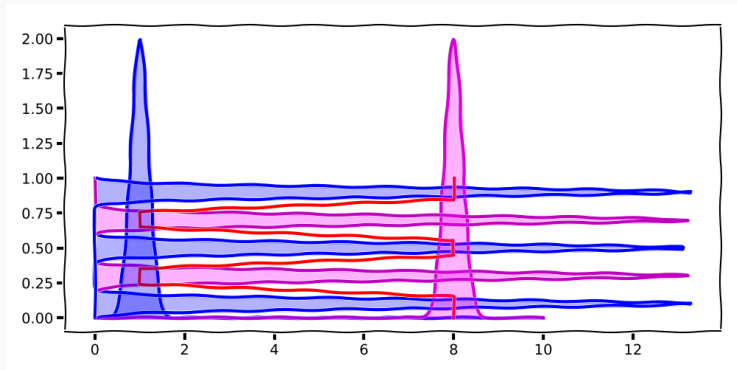
Change of Variables



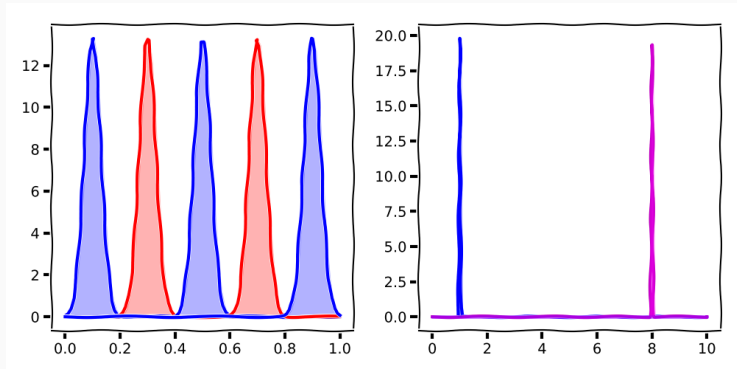
Change of Variables



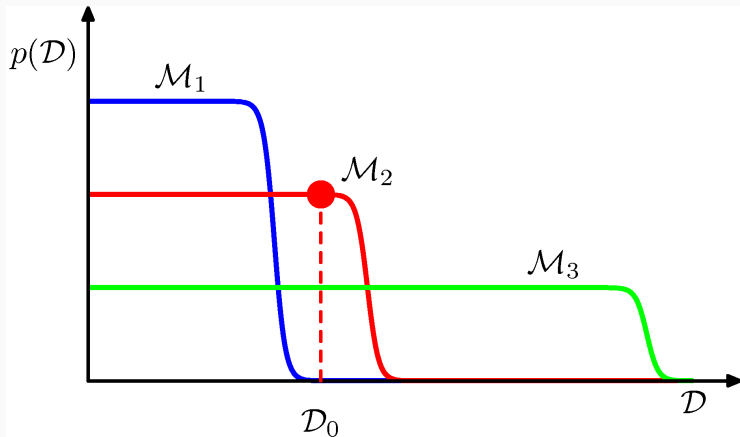
Change of Variables



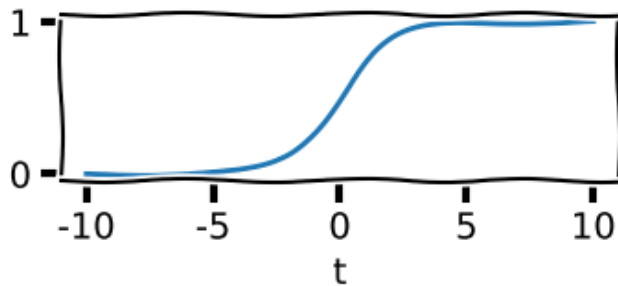
Change of Variables



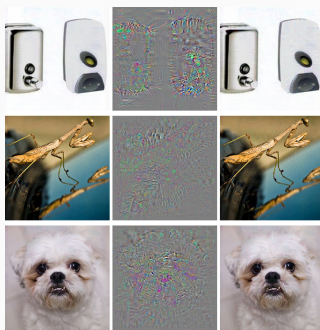
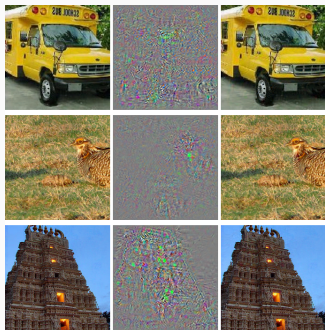
MacKay plot



Activation functions



Data inefficiency¹²



¹²Nguyen, A. M., Yosinski, J., & Clune, J., Deep neural networks are easily fooled: high confidence predictions for unrecognizable images, CoRR, abs/1412.1897(), (2014).

Summary

Summary

- Unsupervised learning is **very** hard

Summary

- Unsupervised learning is **very** hard
 - *Its actually not, its really really easy.*

Summary

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

Summary

- Unsupervised learning 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

Summary

- Unsupervised learning 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 (DPs,GPs) provide strong, interpretative assumptions that aligns well to our intuitions allowing us to make **relevant** assumptions

Summary II

- Composite functions **cannot** model more things

Summary II

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

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

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
- Intuitions needs to change, we need to think of priors over hierarchies

eof

References



Pierre Simon Laplace.

A philosophical essay on probabilities, 1814.