

# Unsupervised Learning with Gaussian Processes

---

Carl Henrik Ek - [carlhenrik.ek@bristol.ac.uk](mailto:carlhenrik.ek@bristol.ac.uk)

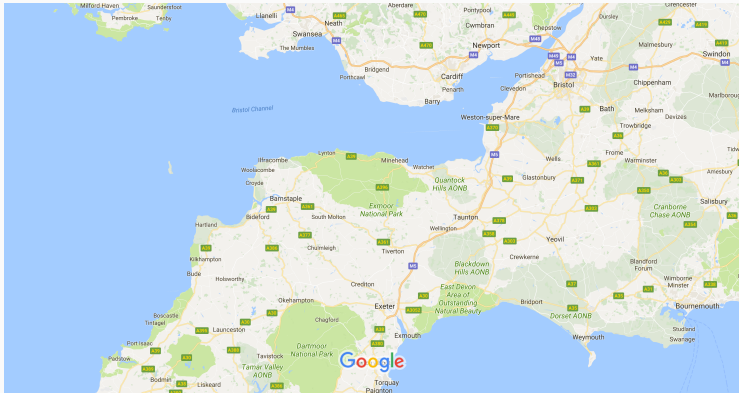
September 11, 2019

<http://www.carlhenrik.com>

# Introductions

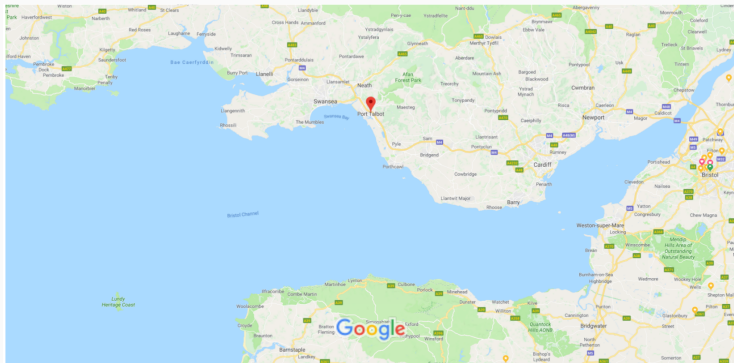
---

# This where I live



This is what I do



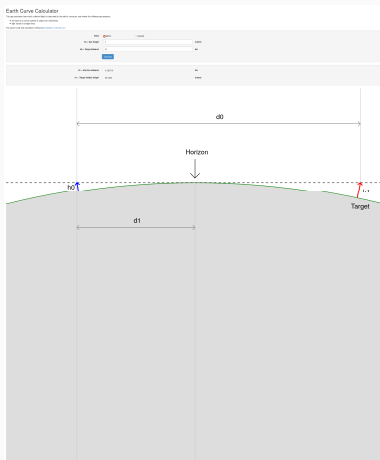










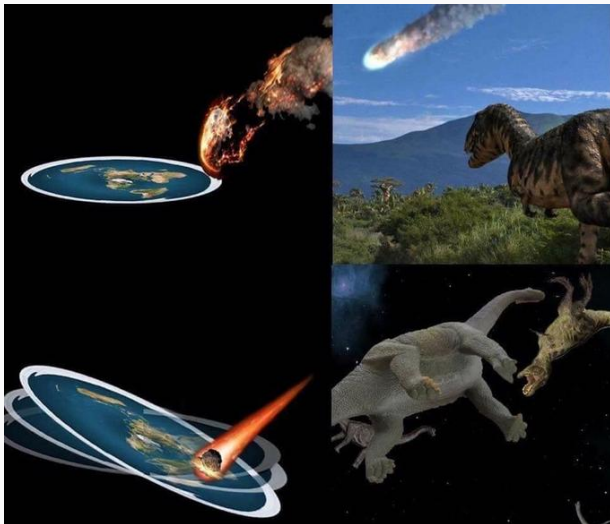


Distance to horizon  $6.2km$

Hidden height  $125.6m$









- $\mathcal{F}$  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

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

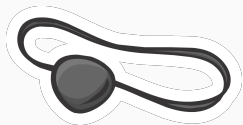
$$\begin{aligned}e(\mathcal{S}, \mathcal{A}, \mathcal{F}) &= \mathbb{E}_{P(\{\mathcal{X}, \mathcal{Y}\})} [\ell(\mathcal{A}_{\mathcal{F}}(\mathcal{S}), x, y)] \\ &\approx \frac{1}{M} \sum_{n=1}^M \ell(\mathcal{A}_{\mathcal{F}}(\mathcal{S}), x_n, y_n)\end{aligned}$$

## No Free Lunch

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



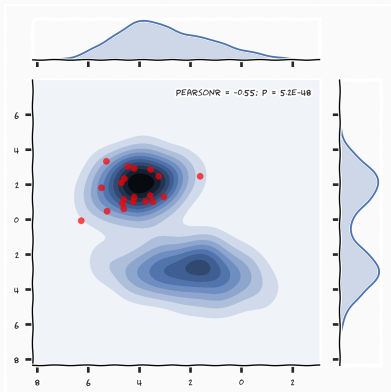
# Assumptions: Algorithms



Statistical Learning

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

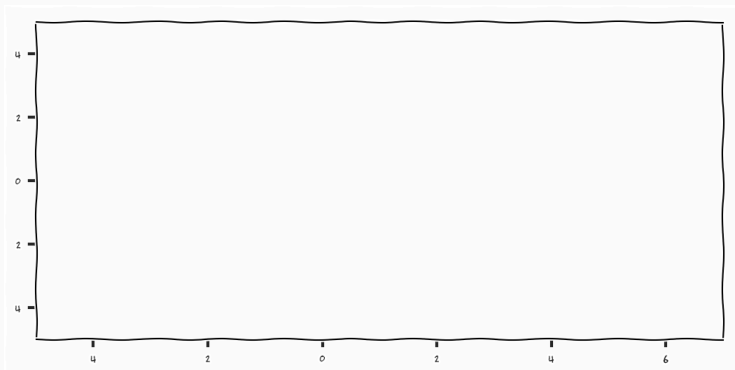
# Assumptions: Biased Sample



Statistical Learning

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

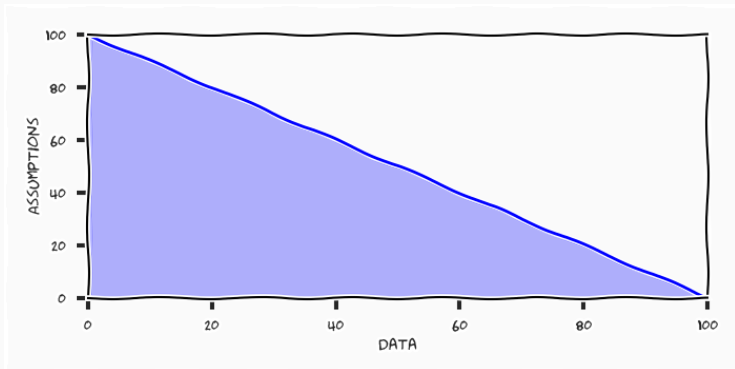
# Assumptions: Hypothesis space



Statistical Learning

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

# Data and Knowledge





*IUDICIUM POSTERIUM DISCIPULUS EST PRIORIS<sup>1</sup>*

---

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

September 11, 2019



UNIVERSITY OF  
**BATH**



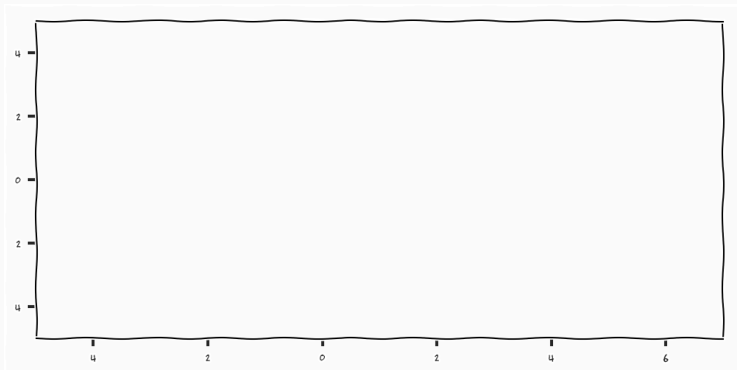
Neill Campbell, Carl Henrik Ek, David Fernandes, Ivan Ustyuzhaninov,  
Aidan Scannell, Emelie Barman, Erik Bodin, Andrew Lawrence, Markus  
Kaiser, Alessandro di Martino, Ieva Kazlauskaitė, Akshaya Thippur

# Gaussian Processes

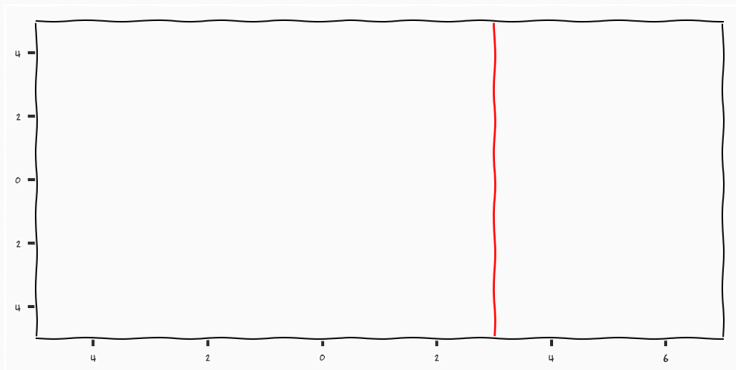
---



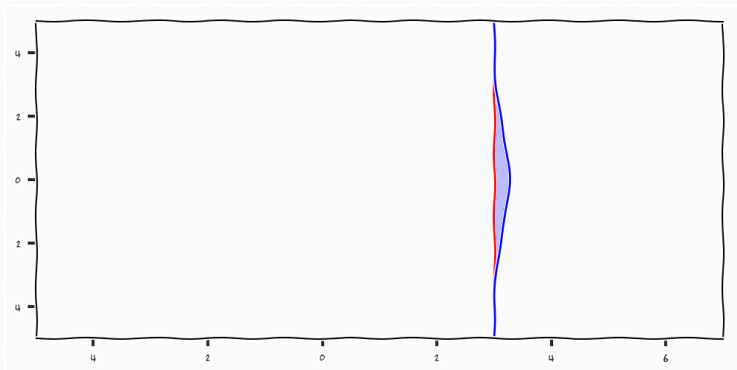
# Gaussian Processes



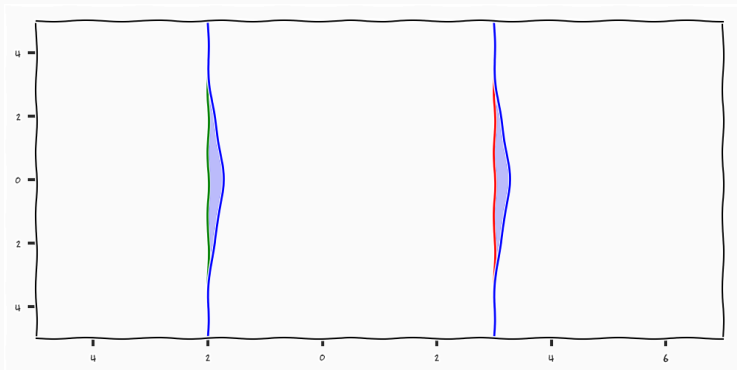
# Gaussian Processes



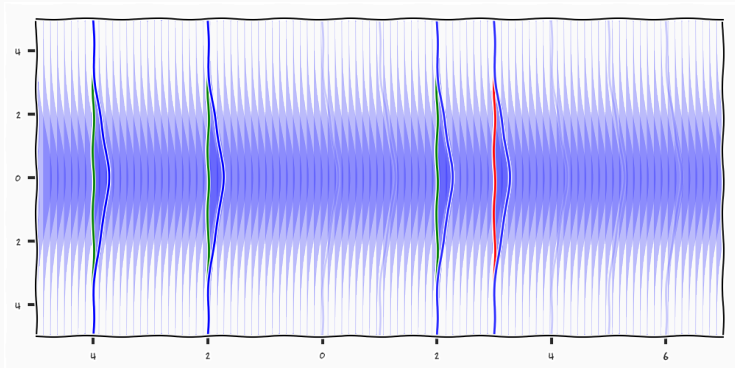
# Gaussian Processes



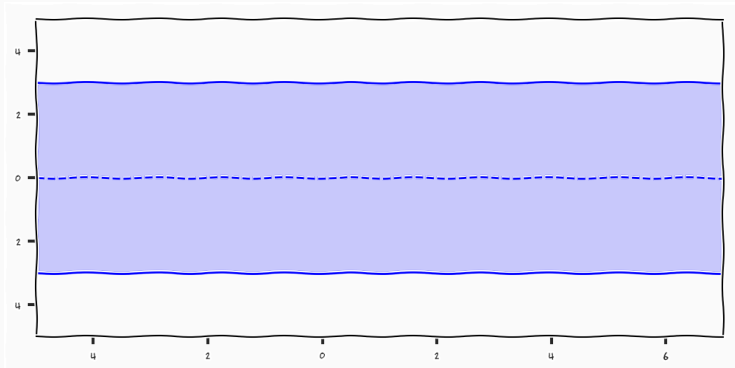
# Gaussian Processes



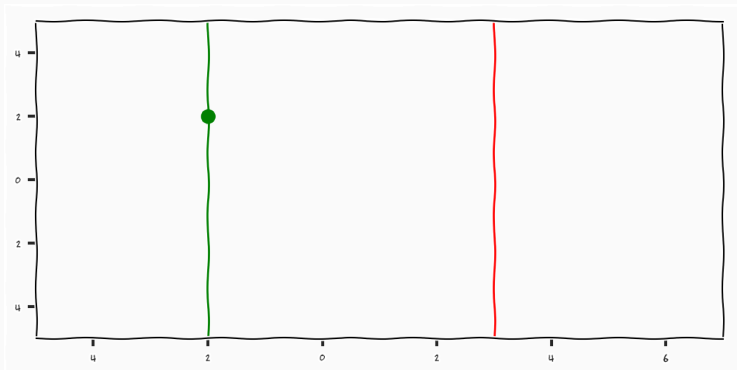
# Gaussian Processes



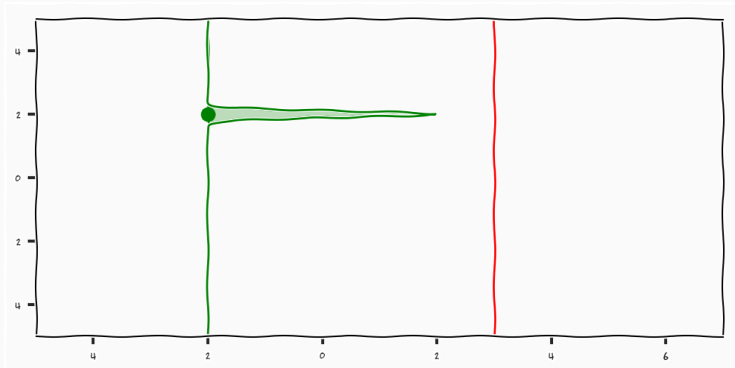
# Gaussian Processes



# Gaussian Processes

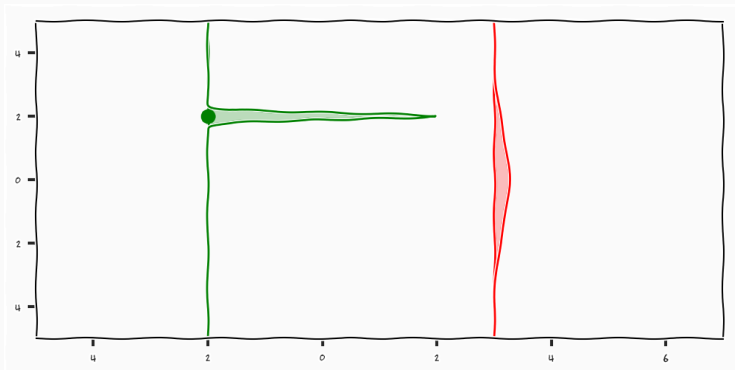


# Gaussian Processes

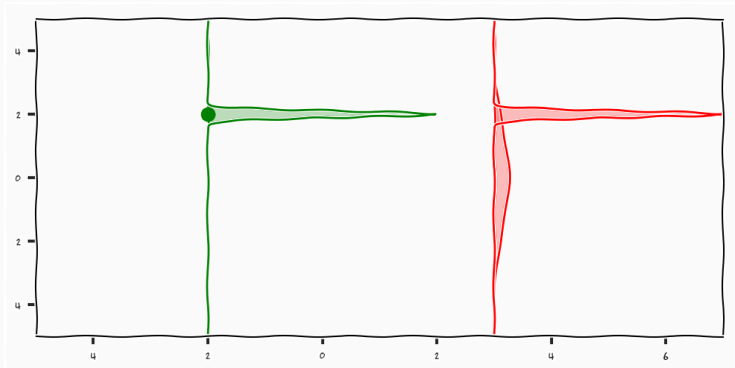




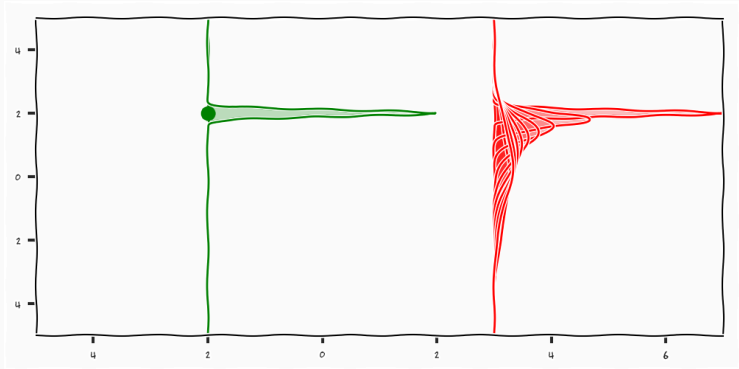
# Gaussian Processes



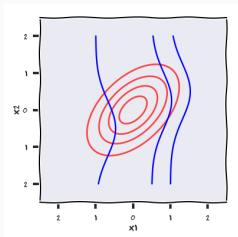
# Gaussian Processes



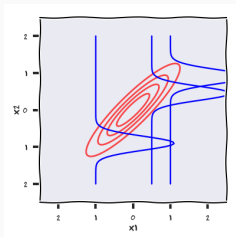
# Gaussian Processes



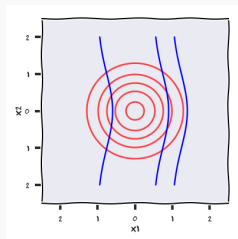
# Conditional Gaussians



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

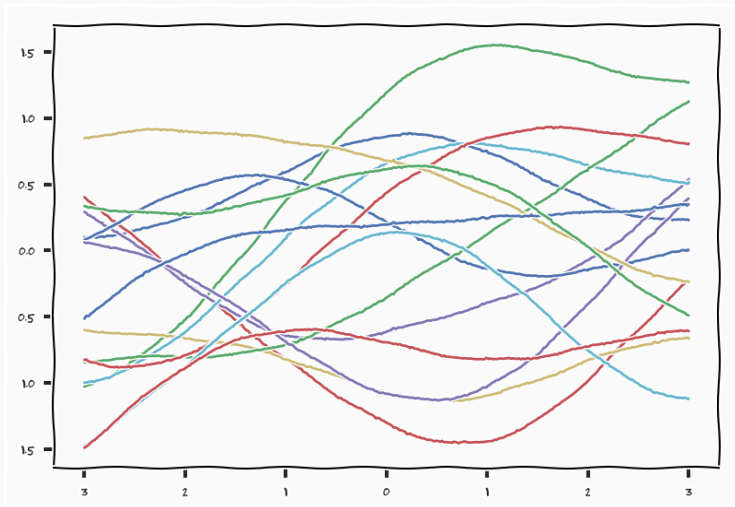


$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}\right)$$

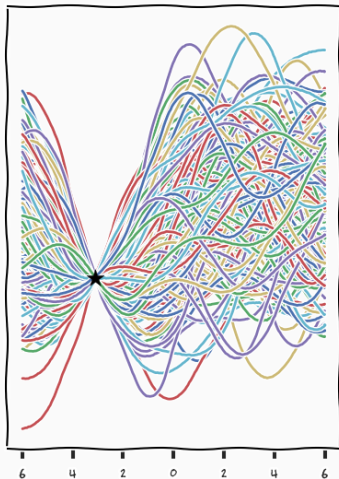
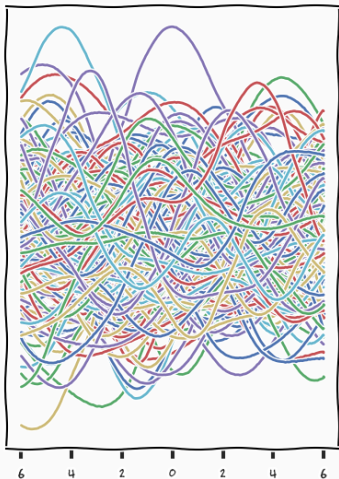


$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

# Gaussian Processes



# Gaussian Processes



# The Gaussian Identities

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

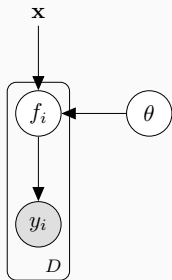
Gaussian Identities

# Unsupervised Learning with GPs

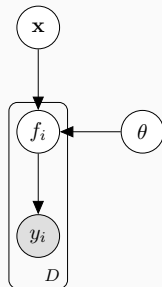
---



# Unsupervised Learning

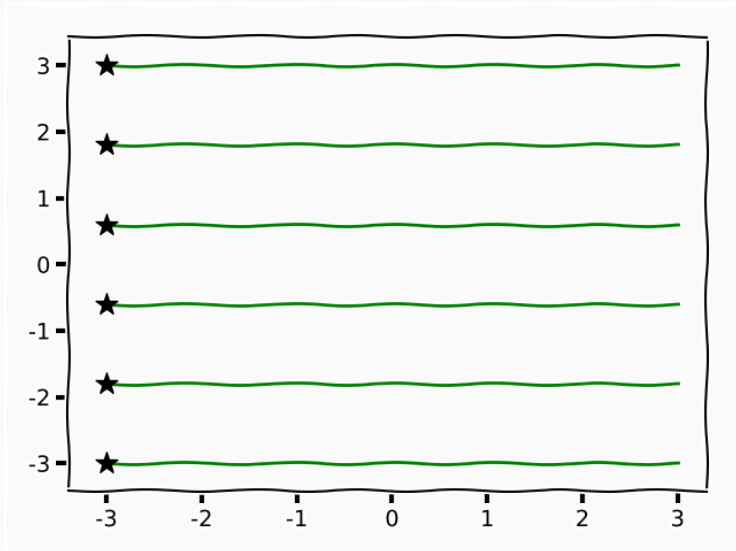


$$p(y|x)$$

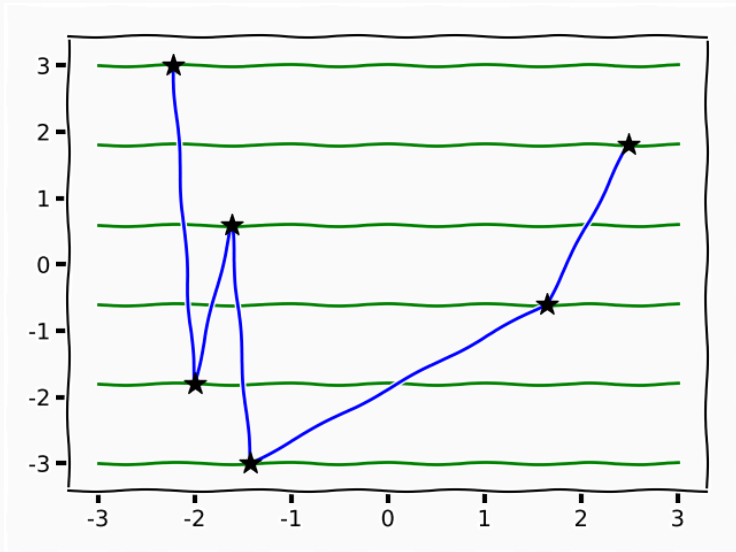


$$p(y)$$

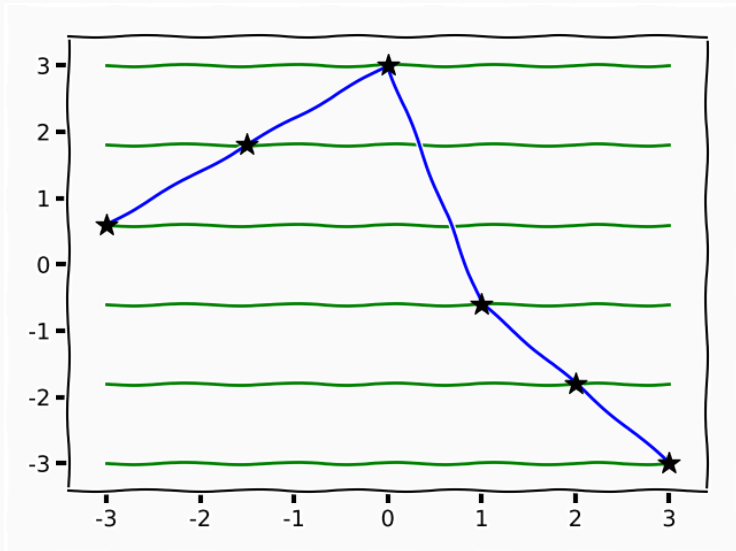
# Unsupervised Learning



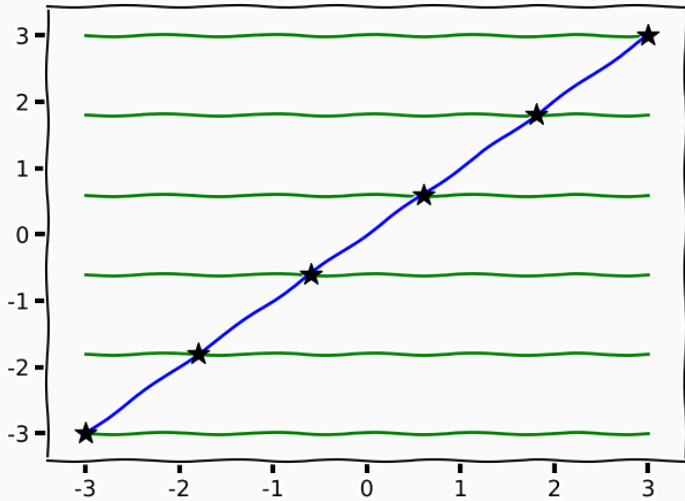
# Unsupervised Learning



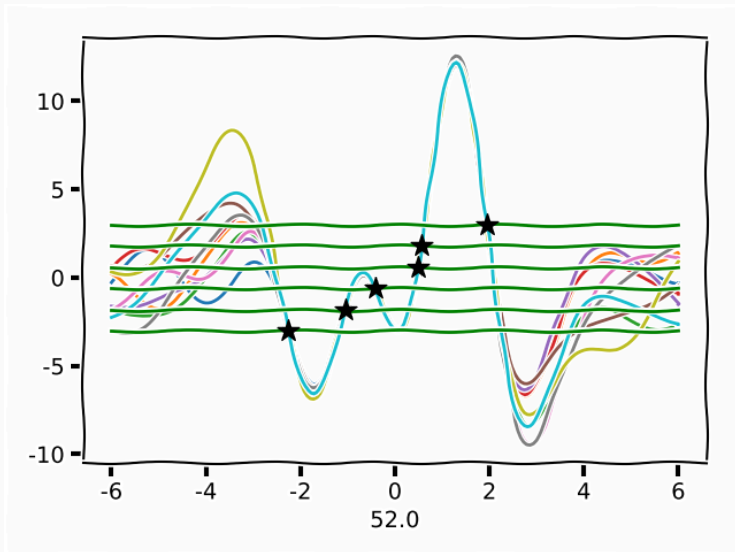
# Unsupervised Learning



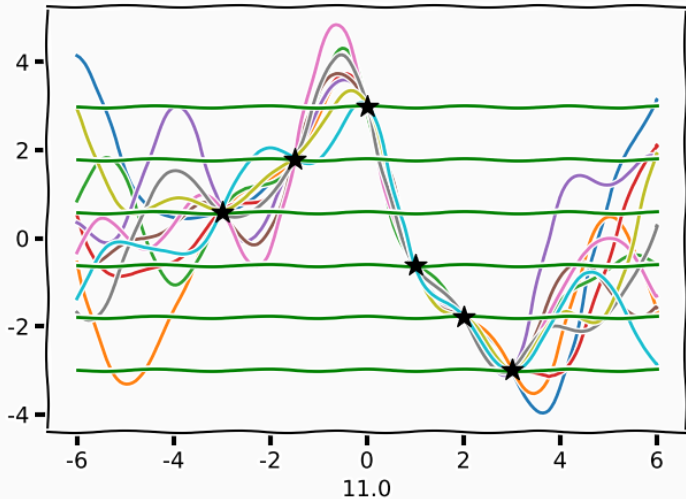
# Unsupervised Learning



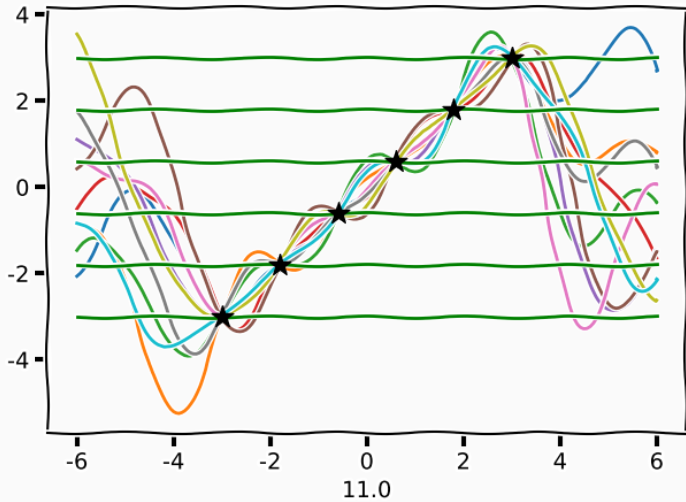
# Unsupervised Learning



# Unsupervised Learning

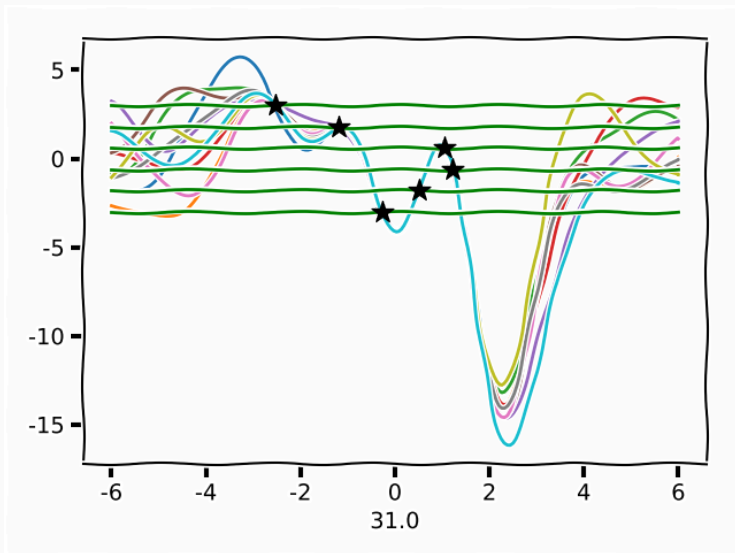


# Unsupervised Learning





# Unsupervised Learning





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

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

1. Priors that makes sense

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

$p(\mathbf{x})$  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

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

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

## 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



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

# Approximate Inference

---



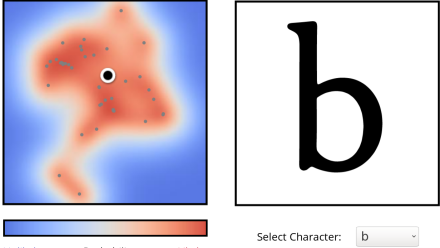
$$\begin{aligned}\hat{x} &= \operatorname{argmax}_x \int p(y|f)p(f|x)dfp(x) \\ &= \operatorname{argmin}_x \frac{1}{2}y^T\mathbf{K}^{-1}y + \frac{1}{2}|\mathbf{K}| - \log p(x)\end{aligned}$$

---

<sup>2</sup>Lawrence, N. D. (2005). Probabilistic non-linear principal component analysis with Gaussian process latent variable models.

- Li, W., Viola, F., Starck, J., Brostow, G. J., & Campbell, N. D. (2016). Roto++: accelerating professional rotoscoping using shape manifolds. (In proceeding of ACM SIGGRAPH'16)
- Grochow, K., Martin, S. L., Hertzmann, A., & Popović, Zoran (2004). Style-based inverse kinematics. SIGGRAPH '04: SIGGRAPH 2004
- Urtasun, R., Fleet, D. J., & Fua, P. (2006). 3D people tracking with Gaussian process dynamical models. Computer Vision and Pattern Recognition, 2006

Please drag the black and white circle around the heat map to explore the 2D font manifold.



Unlikely Probability Likely

Select Character:

URL

- Challenges with ML estimation
  - How to initialise  $x$ ?
  - What is the dimensionality  $q$ ?
- *Our assumption on the latent space does not reach the data*

---

<sup>3</sup>Titsias, M. (2009). Variational learning of inducing variables in sparse Gaussian processes.

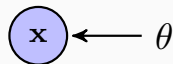
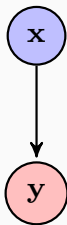
<sup>4</sup>Titsias, M., & Lawrence, N. D. (2010). Bayesian Gaussian Process Latent Variable Model

- Challenges with ML estimation
  - How to initialise  $x$ ?
  - What is the dimensionality  $q$ ?
- *Our assumption on the latent space does not reach the data*
- Approximate integration!<sup>3</sup>

---

<sup>3</sup>Titsias, M. (2009). Variational learning of inducing variables in sparse Gaussian processes.

<sup>4</sup>Titsias, M., & Lawrence, N. D. (2010). Bayesian Gaussian Process Latent Variable Model



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

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

$$p(y)$$

$$\log p(y)$$



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

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

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

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

## Variational Bayes

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

## The log term

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

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

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



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

- if we maximise the ELBO we,
  - find an approximate posterior
  - 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 Adams<sup>5</sup>

---

<sup>5</sup>Talking 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 Adams<sup>5</sup>

---

<sup>5</sup>Talking 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 Adams<sup>5</sup>

---

<sup>5</sup>[Talking Machines Season 2, Episode 5](#)

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

---

<sup>6</sup>Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

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

---

<sup>6</sup>Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

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

---

<sup>6</sup>Damianou, A. C. (2015). Deep Gaussian Processes and Variational Propagation of Uncertainty (Doctoral dissertation)

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

---

<sup>6</sup>Damianou, 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 approximations<sup>7</sup>

---

<sup>7</sup>Quinonero-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,

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

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

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

- 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  $\mathbf{F}$

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

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

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

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

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

## Latent space priors

---

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

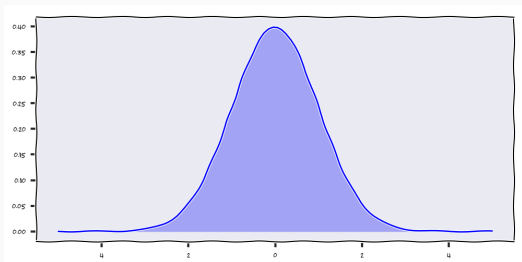
- Importantly  $p(\mathbf{X})$  appears only in KL term
- Allows us to express stronger assumptions about the model

---

<sup>8</sup>Damianou, A. C., Titsias, M., & Lawrence, Neil D, Variational Inference for Uncertainty on the Inputs of Gaussian Process Models (2014)



# The Gaussian blob



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

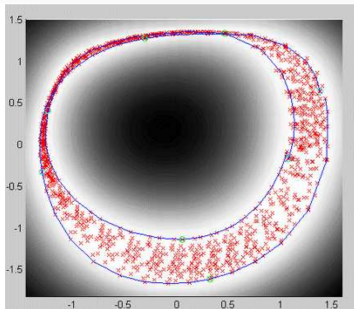
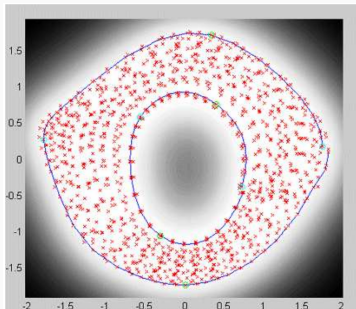
$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma e^{-\sum_d^D \alpha_d \cdot (x_{i,d} - x_{j,d})^2}$$

## GPy

### Code

```
RBF(..., ARD=True)  
Matern32(..., ARD=True)
```

# Dynamic Gaussian Processes<sup>9, 10</sup>



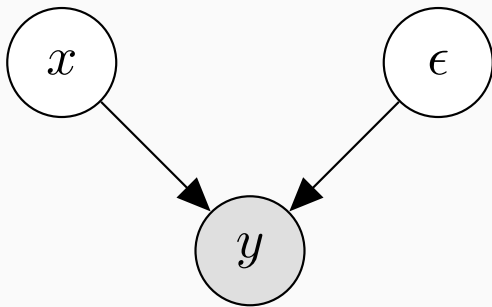
$$p(y, f, x|t) = p(y|f)p(f|x) \underbrace{p(x|t)}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{I})}$$

<sup>9</sup>Urtasun, R., Fleet, D. J., & Fua, P., 3d people tracking with gaussian process dynamical models, CVPR(2006)

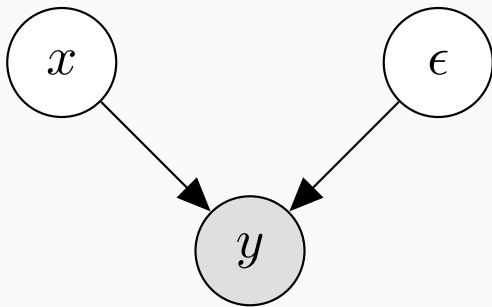
<sup>10</sup>Damianou, A. C., Titsias, M., & Lawrence, N. D., Variational Gaussian Process Dynamical Systems, 2011

## Latent space structures

---

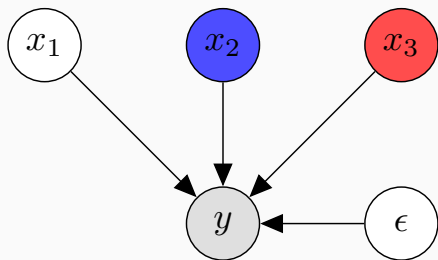


$$y = f(x) + \epsilon$$



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

# Factor Analysis



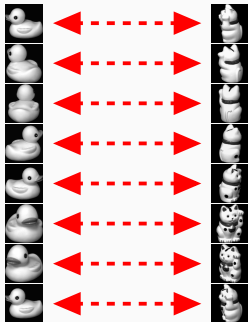
$$y = f(x_1, x_2, x_3) + \epsilon$$

# Alignments

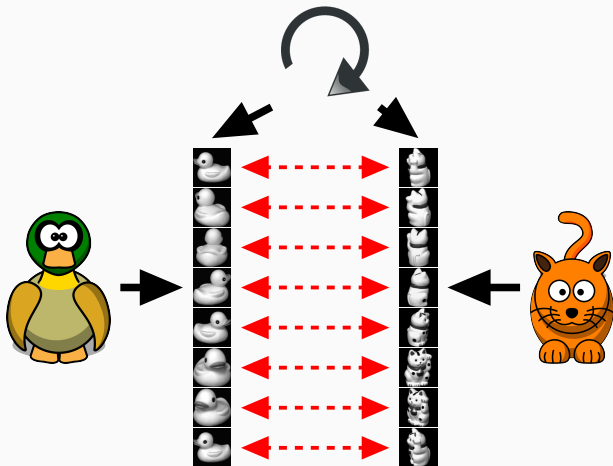




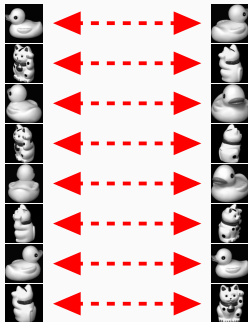
# Alignments



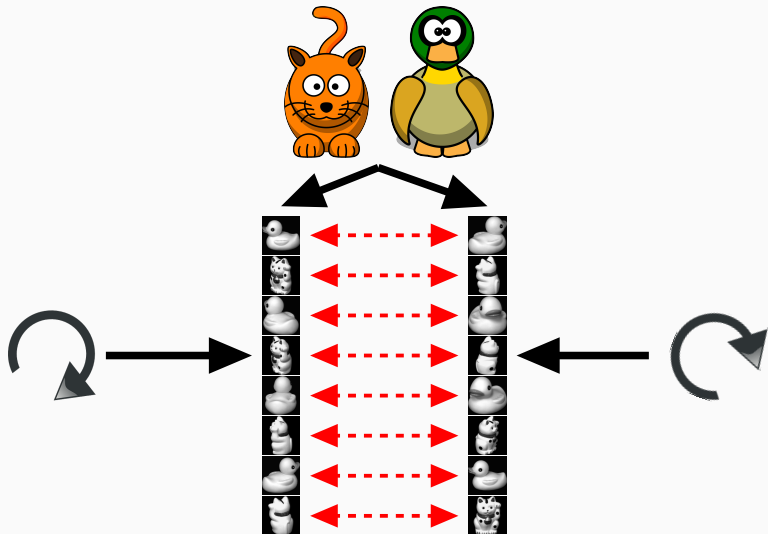
# Alignments



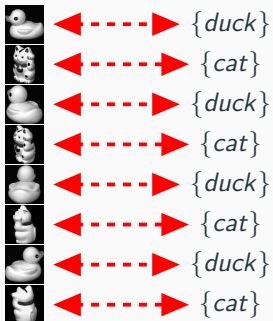
# Alignments



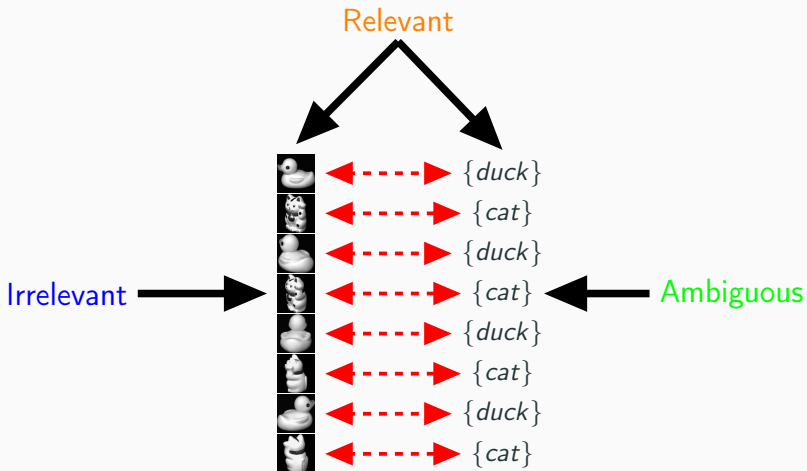
# Alignments

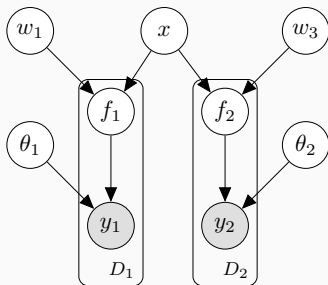


# Alignments



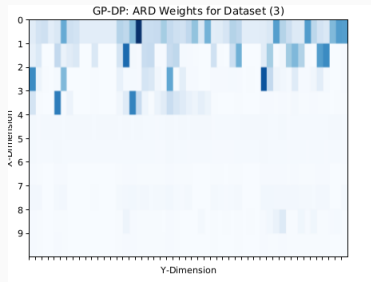
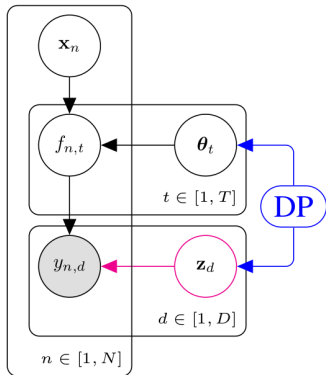
# Alignments





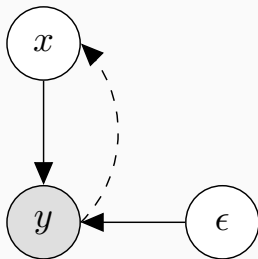
$$y_1 = f(w_1^T x) \quad y_2 = f(w_2^T x)$$

<sup>11</sup>Damianou, A., Lawrence, N. D., & Ek, C. H. (2016). Multi-view learning as a nonparametric nonlinear inter-battery factor analysis



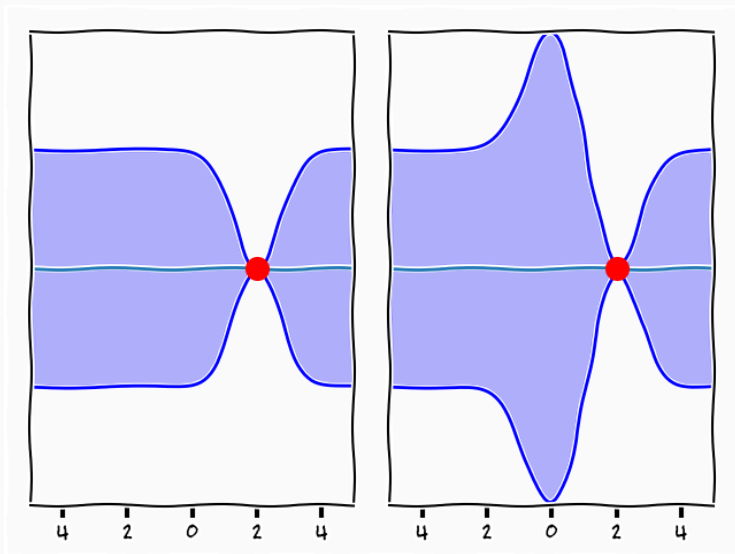
<sup>12</sup>Lawrence, A. R., Ek, C. H., & Campbell, N. W., DP-GP-LVM: A bayesian non-parametric model for learning multivariate dependency structures, ICML (2019)





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

<sup>13</sup>Lawrence, N. D., & Quinero-Candela, Joaquin, Local distance preservation in the gp-lvm through back constraints, ICML, 2006

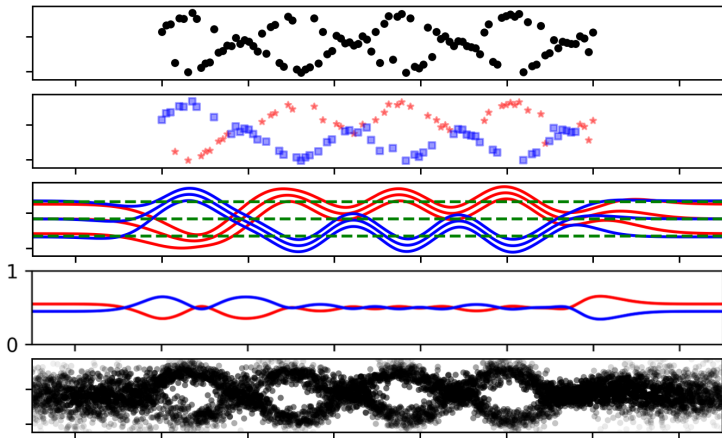


$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{F}) p(\mathbf{F}|\mathbf{X}, \mathbf{X}^{(C)}) p(\mathbf{X}^{(C)}) d\mathbf{F} d\mathbf{X}^{(C)}.$$

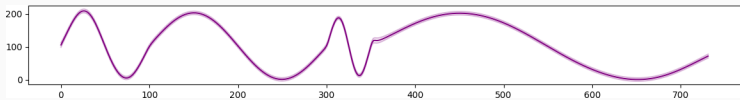
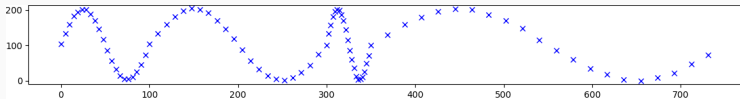
---

<sup>14</sup>Bodin, E., Campbell, N. D. F., & Ek, C. H., Latent Gaussian Process Regression (2017).

# Discrete



# Continuous



# Composite Functions

---

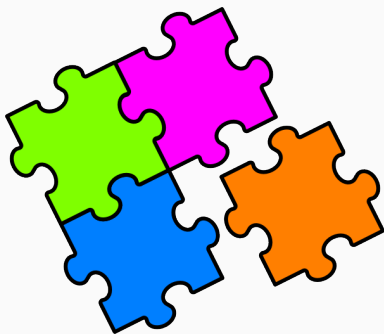


- Place a GP as a warping function, that is warped, ...

---

<sup>15</sup>Damianou, A. C., & Lawrence, N. D. (2013). Deep Gaussian Processes

# Composite Functions



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



## 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"

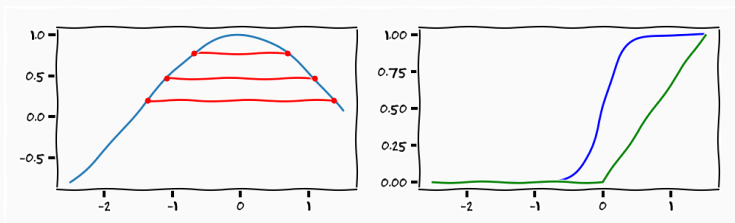


1st layer  
"Edges"



Pixels

# Composite functions



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

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

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

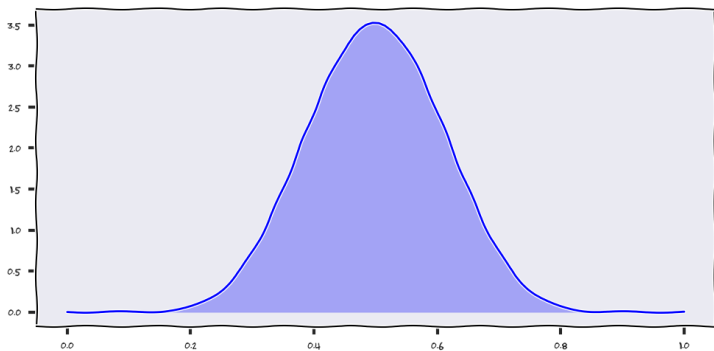
## Theorem (Change of Variable)

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

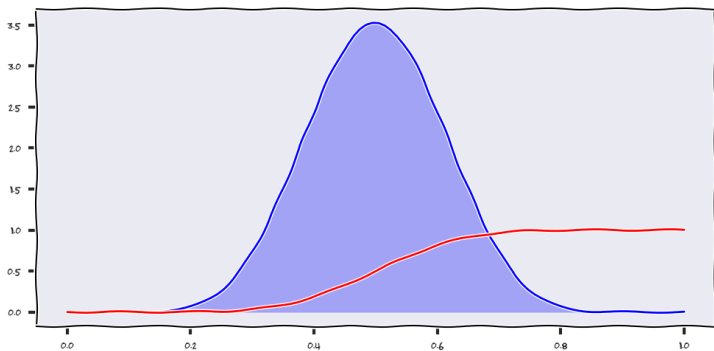
$$p_y(y) = p_x(\psi(y)) |\nabla \psi(y)|$$

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

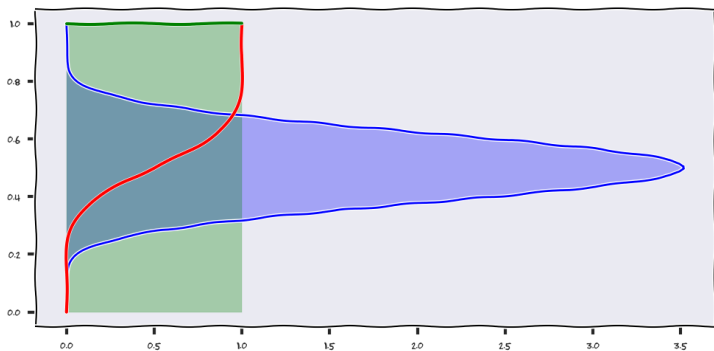
# Sampling



# Sampling



# Sampling



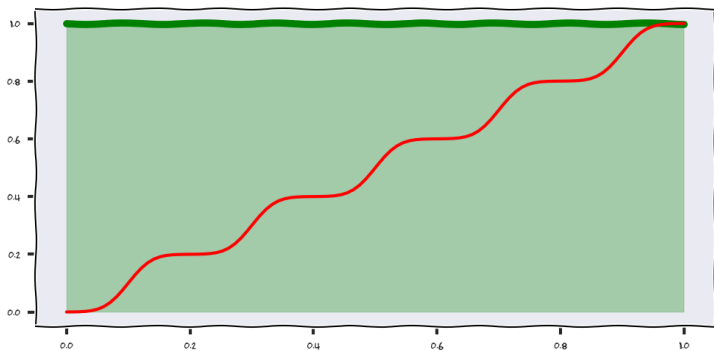
## Theorem (Change of Variable)

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

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

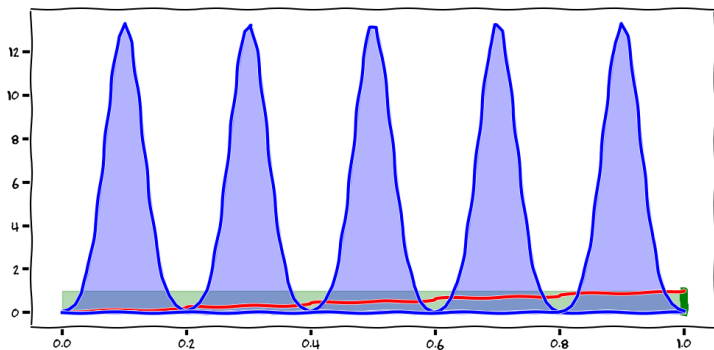
where  $\nabla \psi(\cdot)$  denotes the Jacobian of  $\psi(\cdot)$ , and  $|\cdot|$  denotes the determinant operator.

# Change of Variables

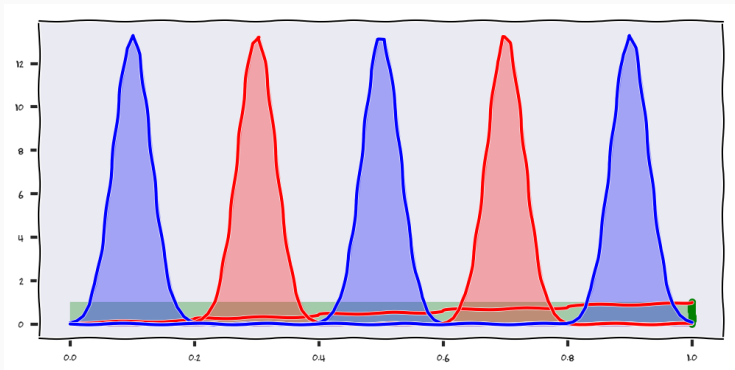




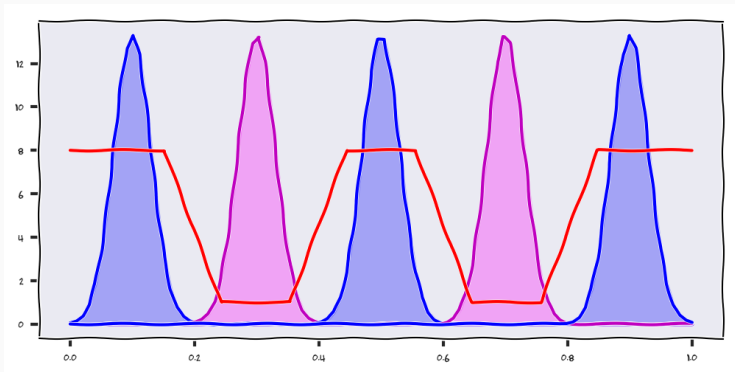
# Change of Variables



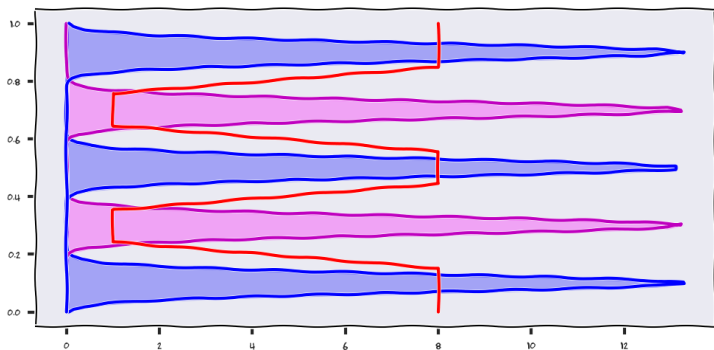
# Change of Variables



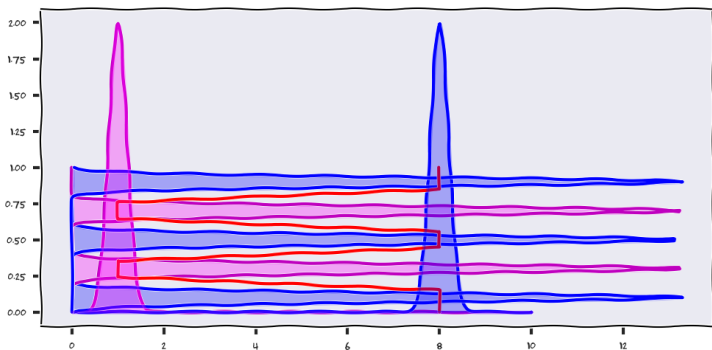
# Change of Variables



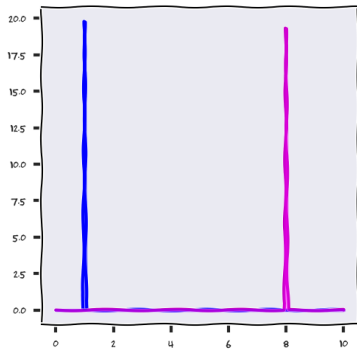
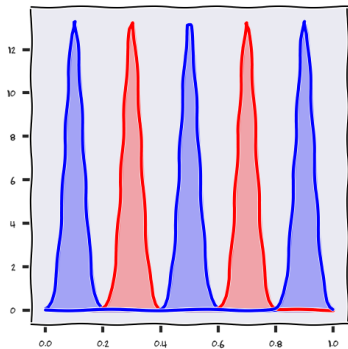
# Change of Variables



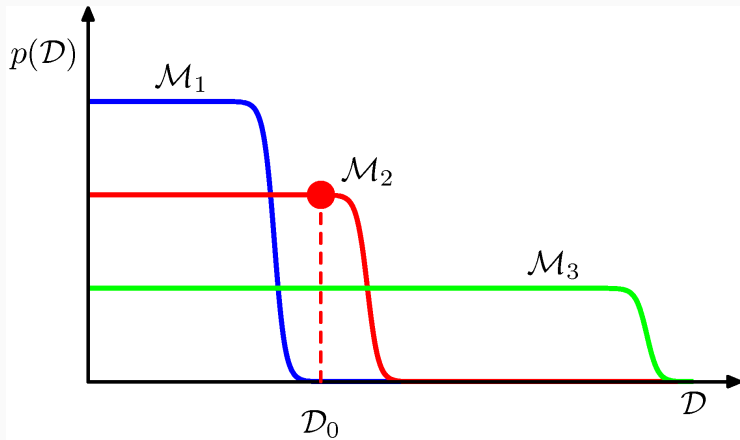
# Change of Variables



# Change of Variables



# MacKay plot



# 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

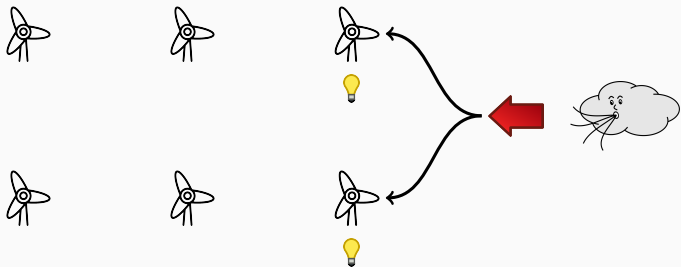


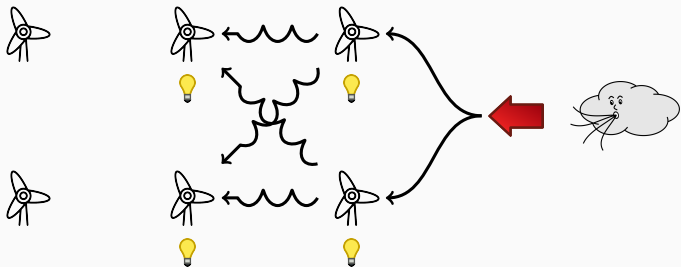
# Windfarms



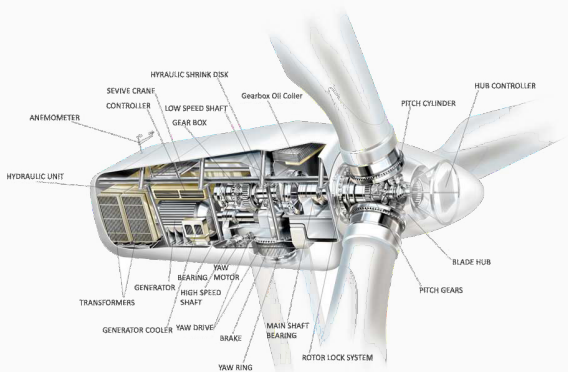
- Effectiveness of modern windfarm
  - 25-60% (of Betz Limit)
- Turbine has several parameters
  - angle and direction of blades
  - gear
  - etc.

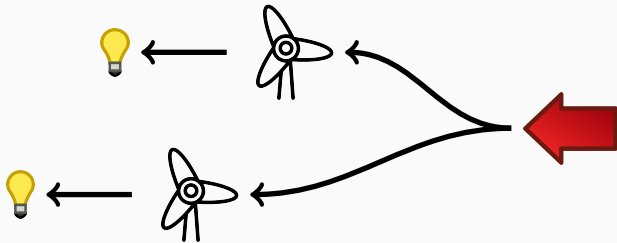
- Effectiveness of modern windfarm
  - 25-60% (of Betz Limit)
- Turbine has several parameters
  - angle and direction of blades
  - gear
  - etc.
- *How can we maximise the efficiency of a windfarm?*



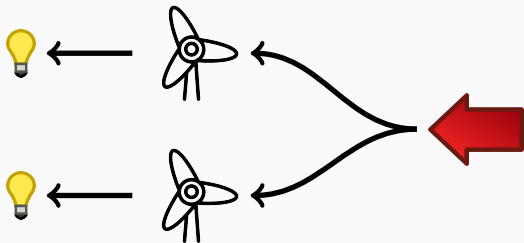


# The Wind Turbine





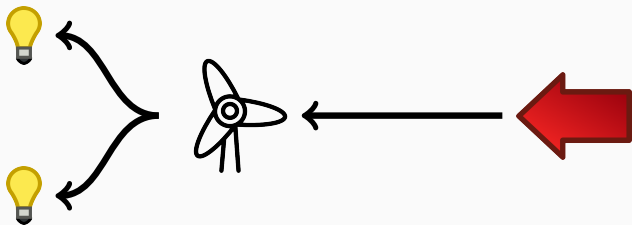
## Model: Alignment



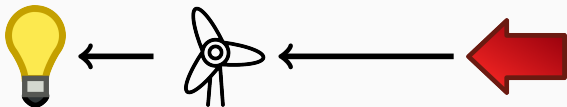
$$w_1(t) = w_2(a(t))$$



## Model: Windfront

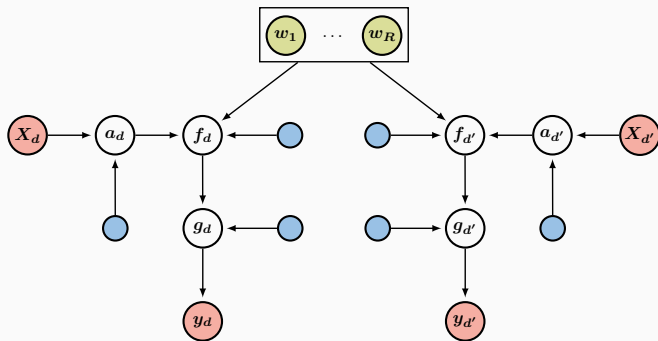


$$f_d(x) = \sum_{r=1}^R \int T_{d,r}(x-z) \cdot w_r(z) \frac{d}{dz}$$



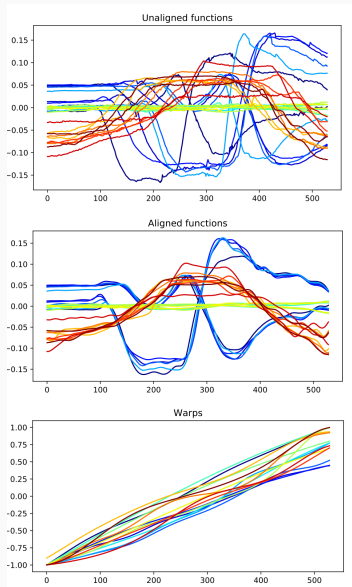
$$y_d = g_d(\mathbf{f}_d)$$

# Model: Graphical Model <sup>16</sup>

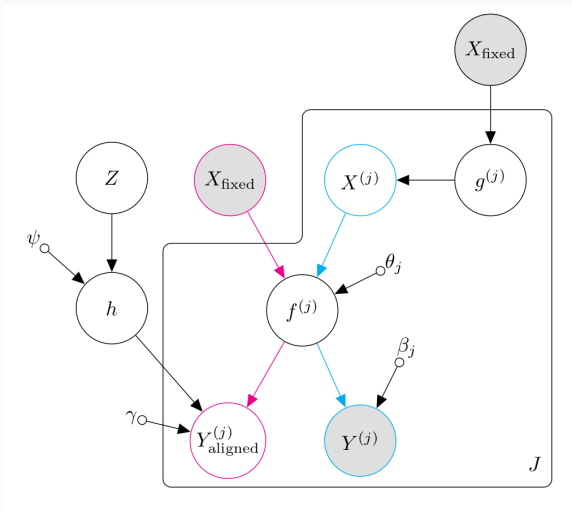


<sup>16</sup>Kaiser, M., Otte, C., Runkler, T., & Ek, C.~H., Bayesian alignments of warped multi-output gaussian processes, NIPS, 2018

# Alignment Learning

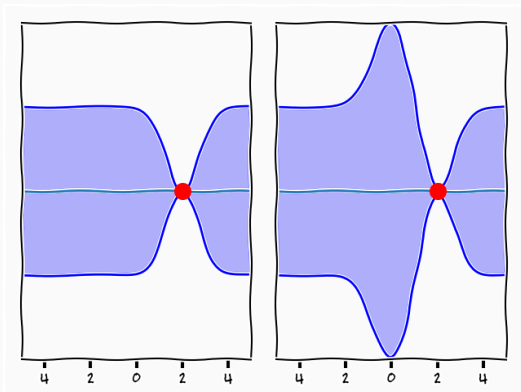


# Alignment Learning<sup>17</sup>



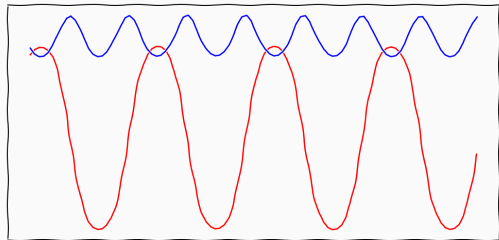
<sup>17</sup>Kazlauskaitė, I., Ek, C. H., & Campbell, N. D. F., Gaussian Process Latent Variable Alignment Learning, AISTATS 2019

# Kernel Re-Parametrisation

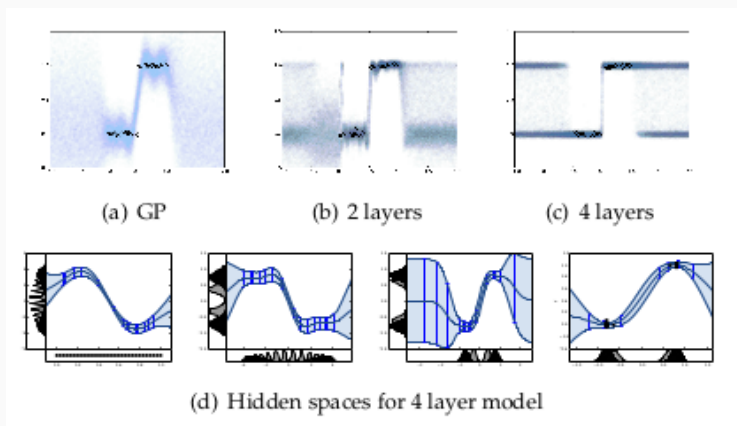


$$k(x'_1, x'_2) = k(f(x_1), f(x_2)) = k([x_1, z_1], [x_2, z_2])$$

## Composition: priors



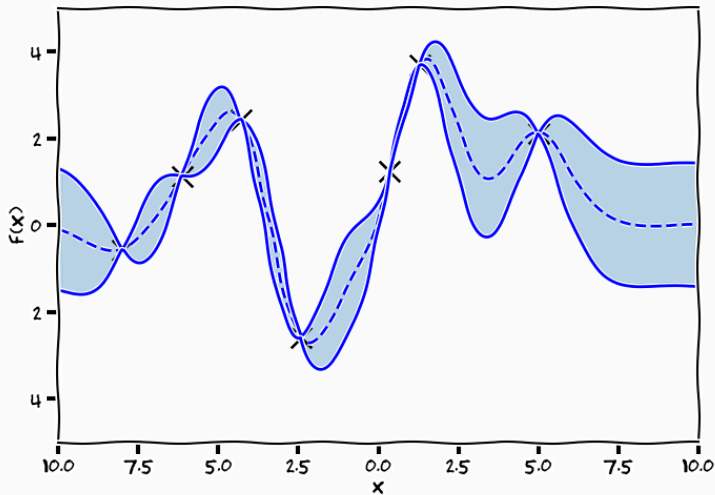
# Composition: priors<sup>18</sup>



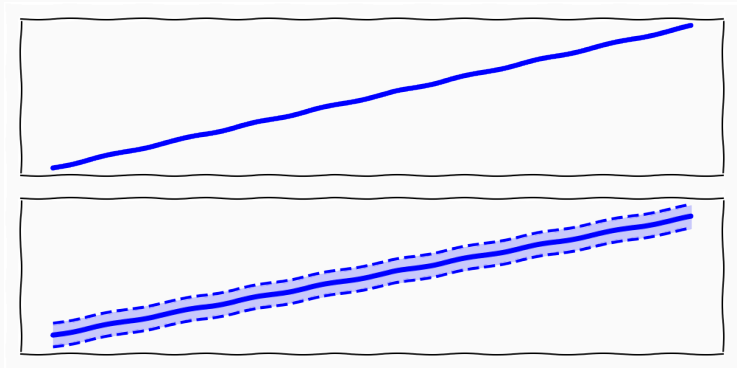
<sup>18</sup>Slides by James Hensman



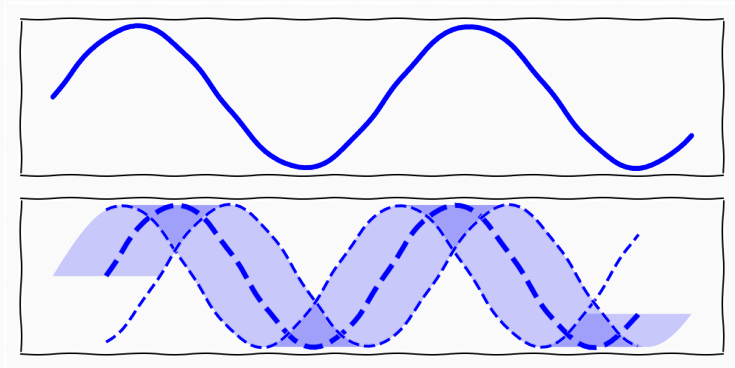
# Propagation of Uncertainty



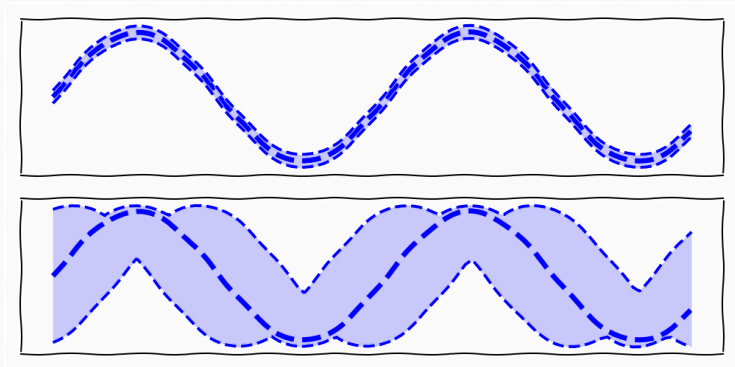
## Composition: uncertainty



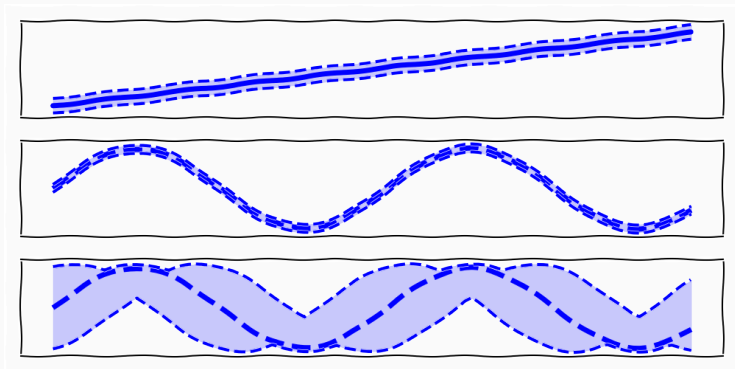
## Composition: uncertainty



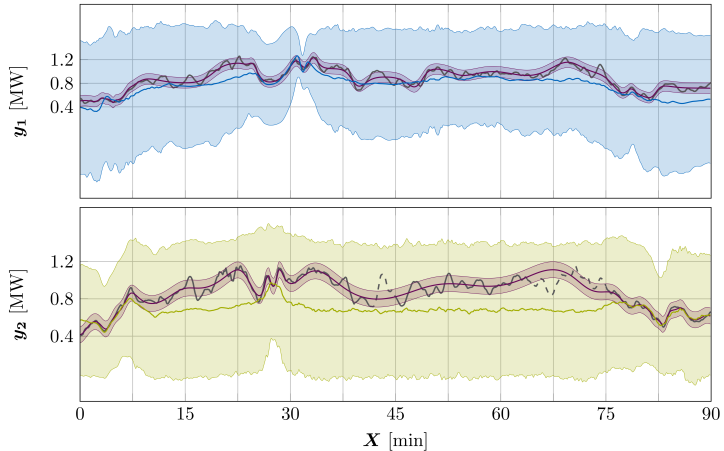
## Composition: uncertainty



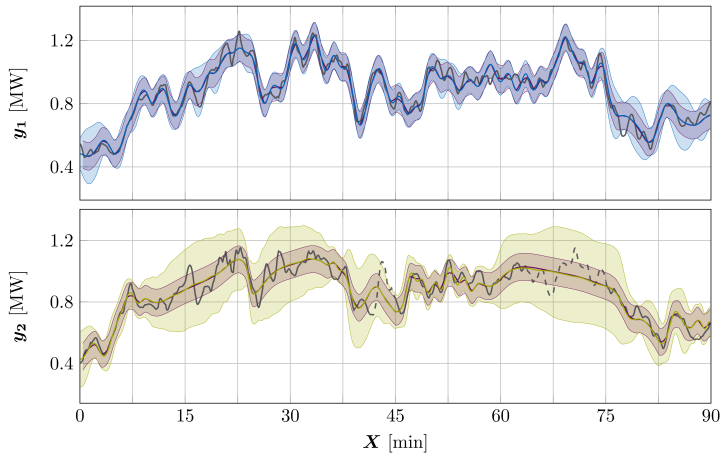
## Composition: uncertainty



# Composition: uncertainty



# Composition: uncertainty



## Summary

---



- Unsupervised learning is **very** hard<sup>19</sup>

---

<sup>19</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning is **very** hard<sup>19</sup>
  - *Its actually not, its really really easy.*

---

<sup>19</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning is **very** hard<sup>19</sup>
  - *Its actually not, its really really easy.*
- Relevant assumptions needed to learn anything useful

---

<sup>19</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning is **very** hard<sup>19</sup>
  - *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>19</sup>I would argue that there is no such thing

# Summary

- Unsupervised learning is **very** hard<sup>19</sup>
  - *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>19</sup>I would argue that there is no such thing

- 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

## 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
- We need to think about correlated uncertainty, not marginals

eof