

Unsupervised Learning with Gaussian Processes

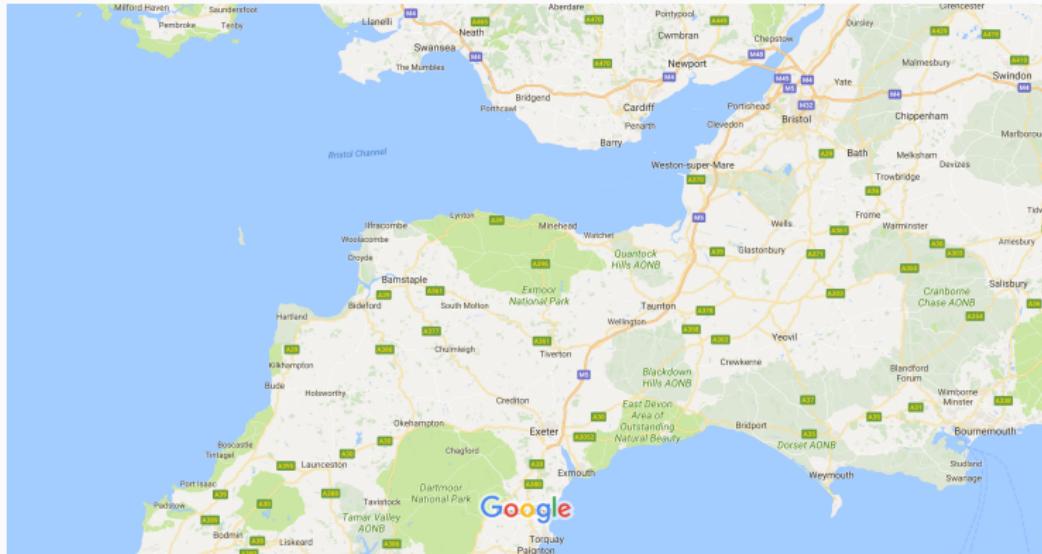
Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

September 5, 2018

<http://www.carlhenrik.com>

Introductions

This where I live



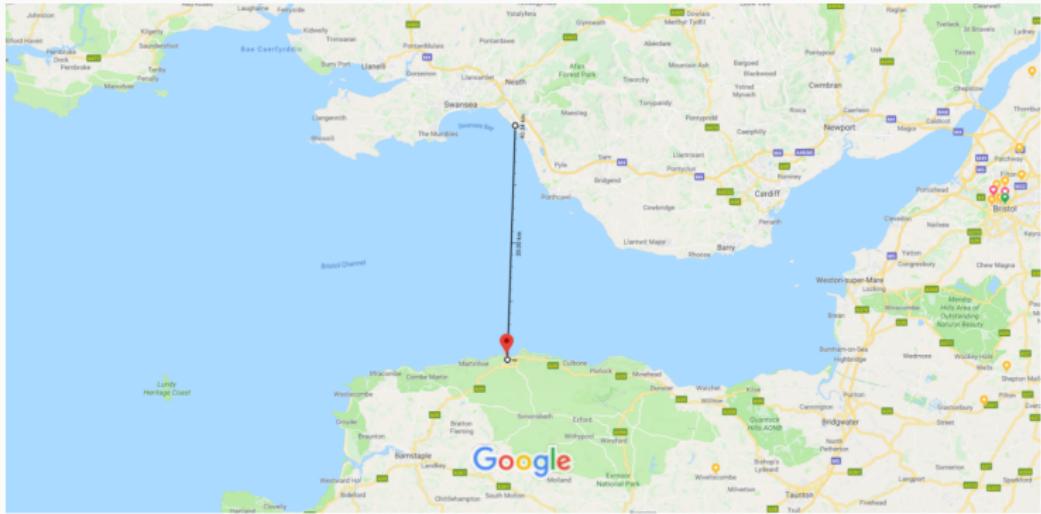
This is what I do

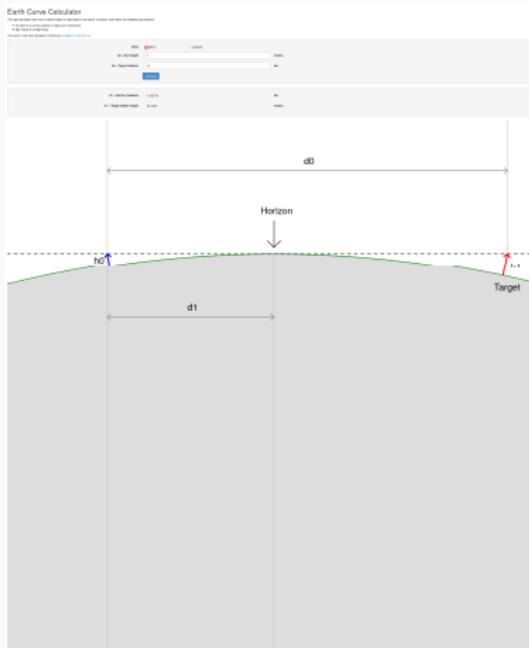










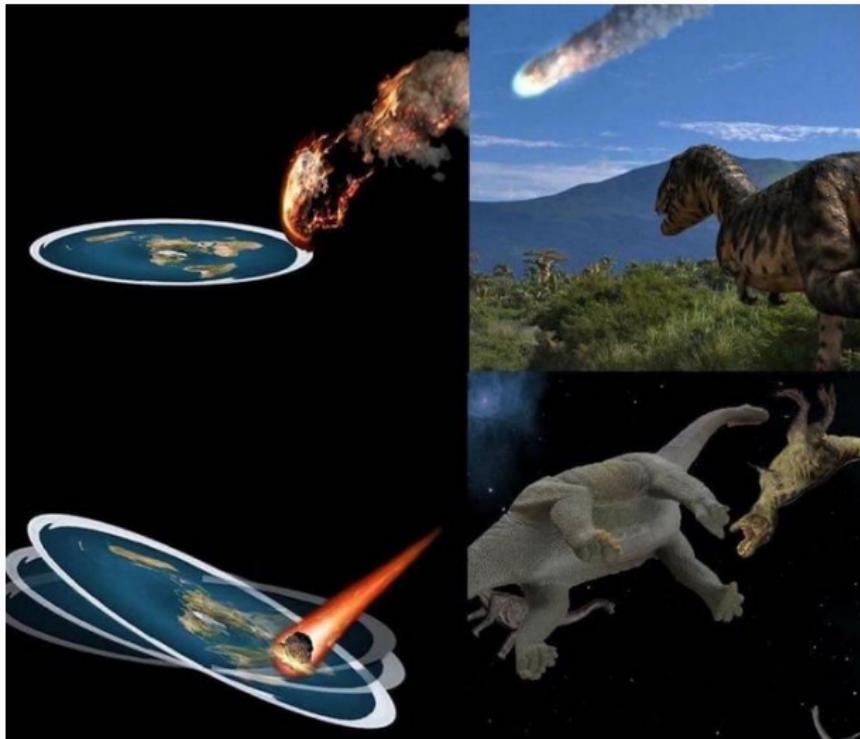


Distance to horizon 6.2km

Hidden height 125.6m









Machine Learning

"In inductive inference, we go from the specific to the general. We make many observations, discern a pattern, make a generalization, and infer an explanation or a theory"

– Wassertheil-Smoller

Learning Theory

- \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

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

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

Statistical Learning

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

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



IUDICIUM POSTERIUM DISCIPULUS EST PRIORIS

Today

September 5, 2018

Learning

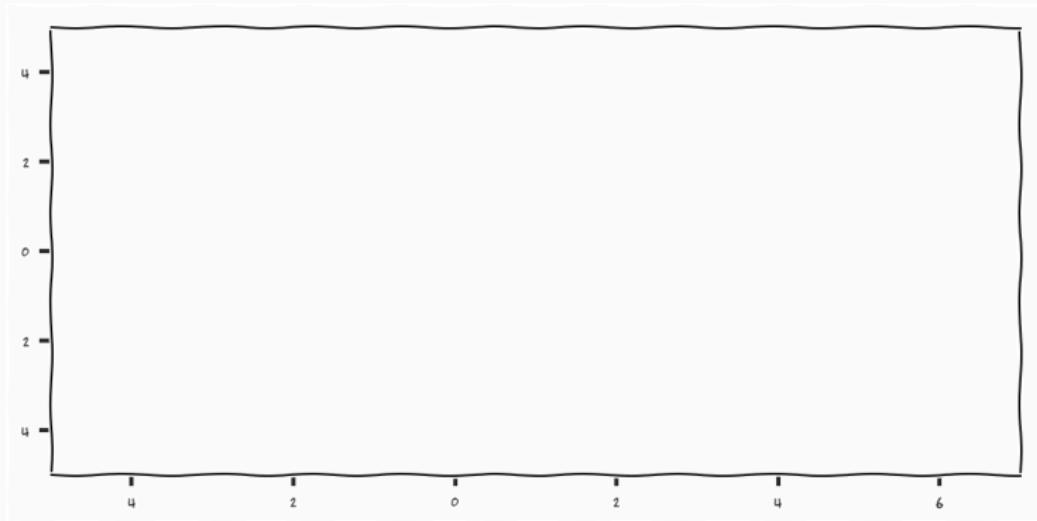
Machine Learning



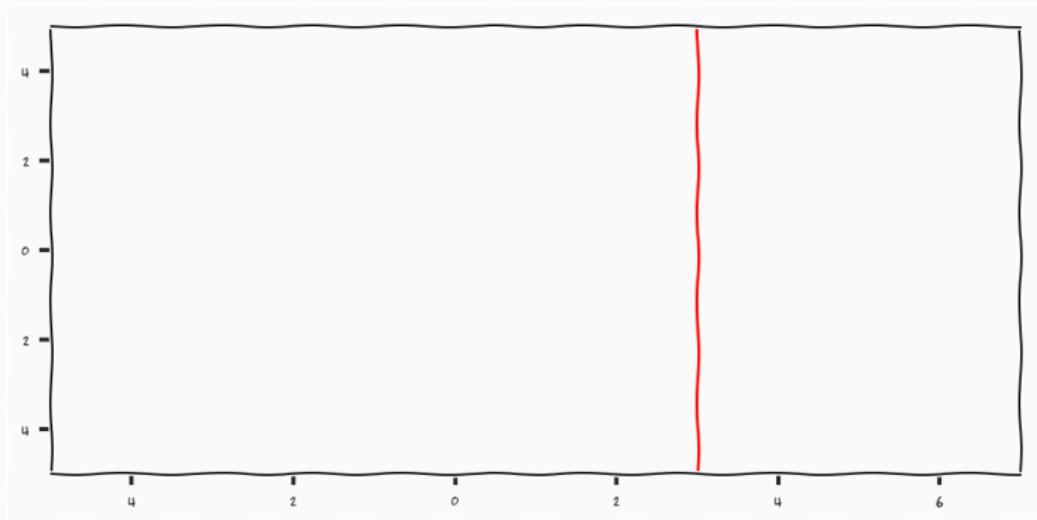
"Machine Learning is nothing but curve fitting, but its amazing what you can do by fitting curves."

– Judea Pearl

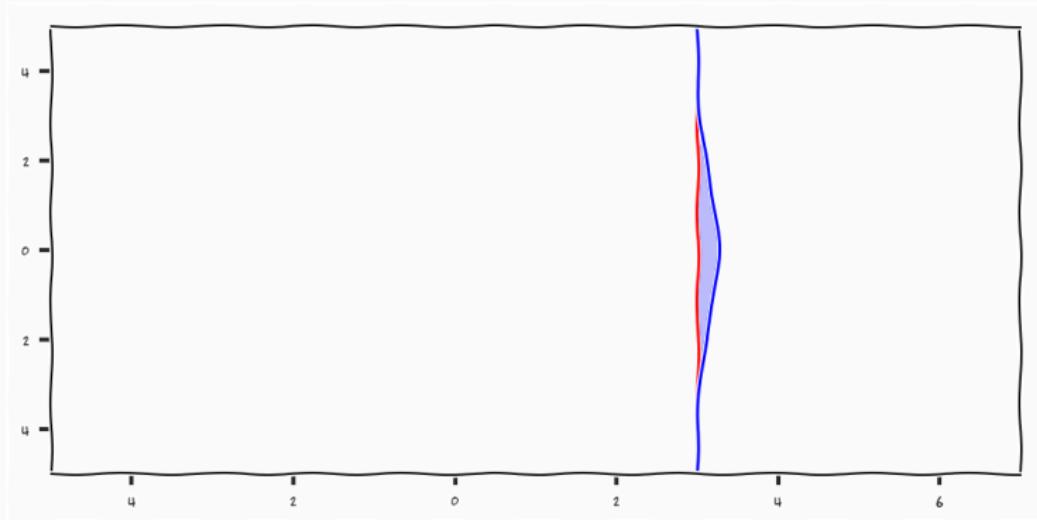
Gaussian Processes



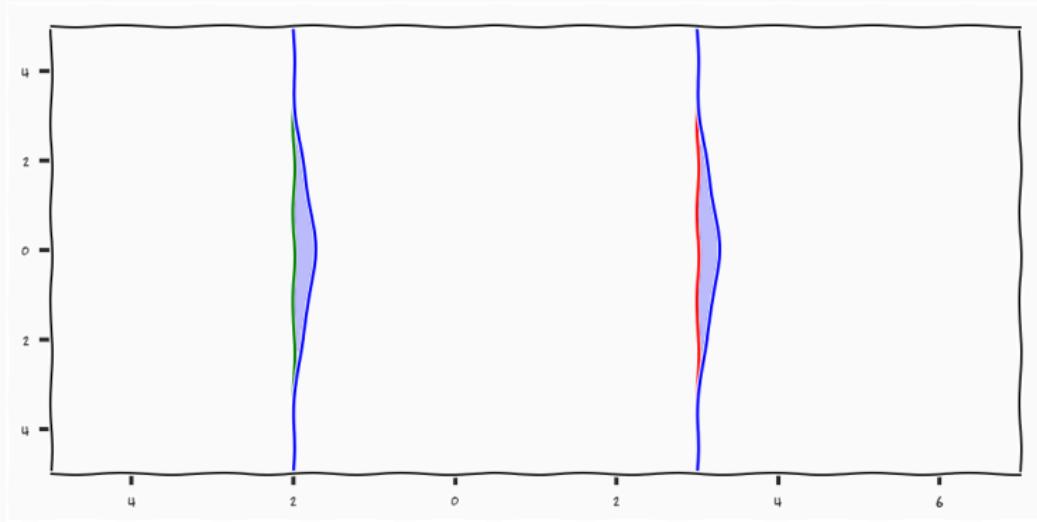
Gaussian Processes



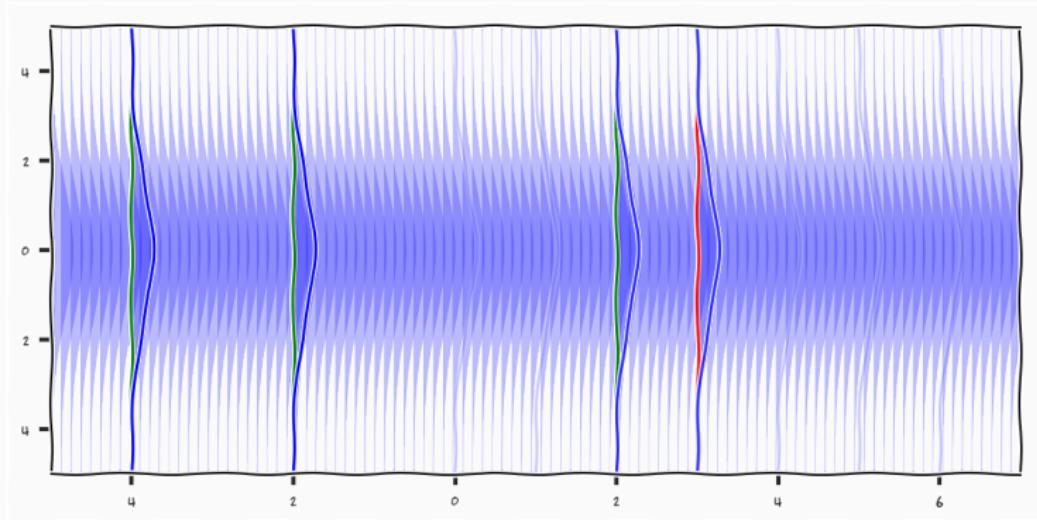
Gaussian Processes



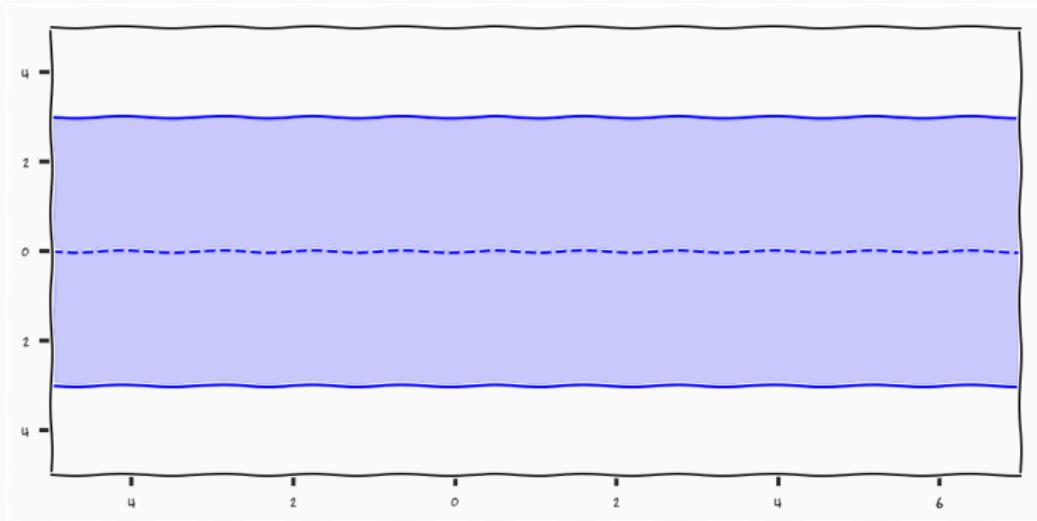
Gaussian Processes



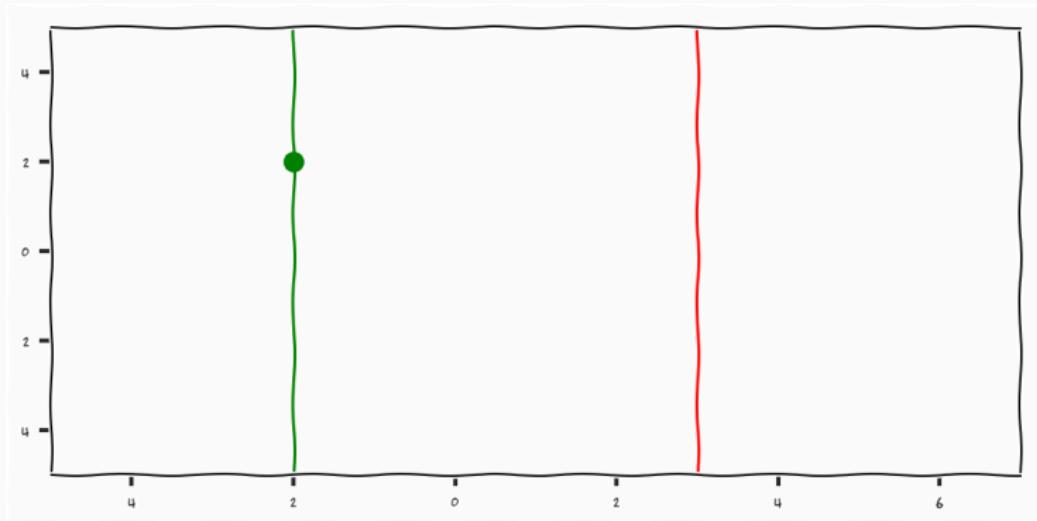
Gaussian Processes



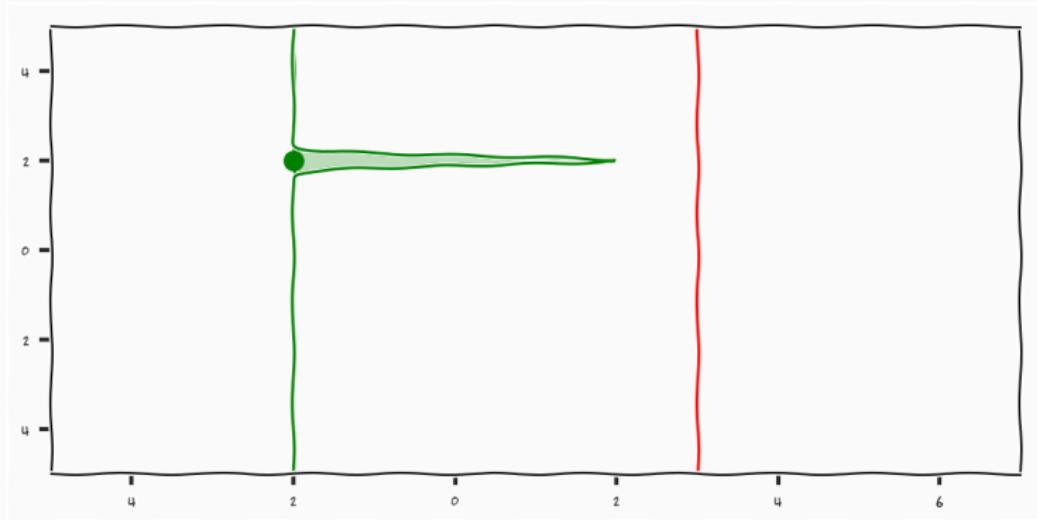
Gaussian Processes



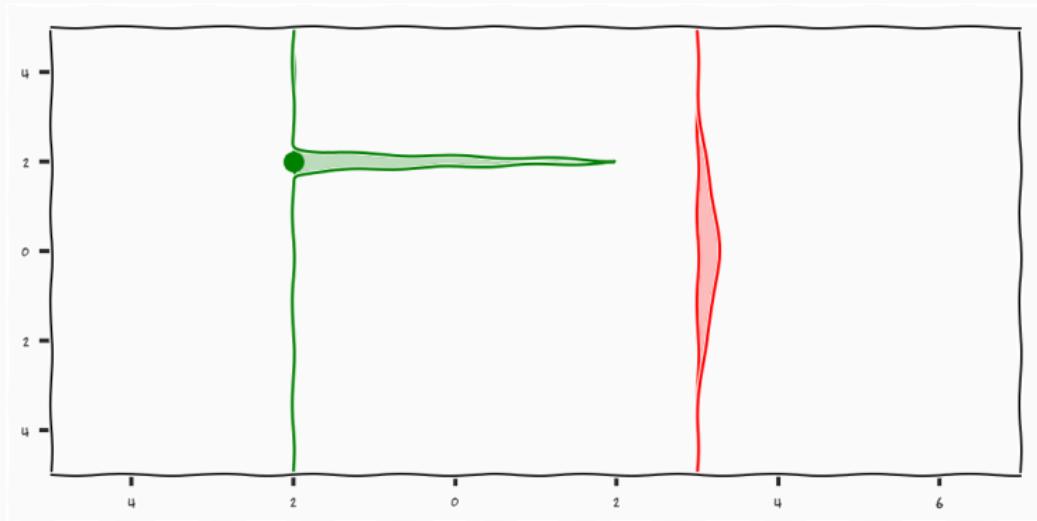
Gaussian Processes



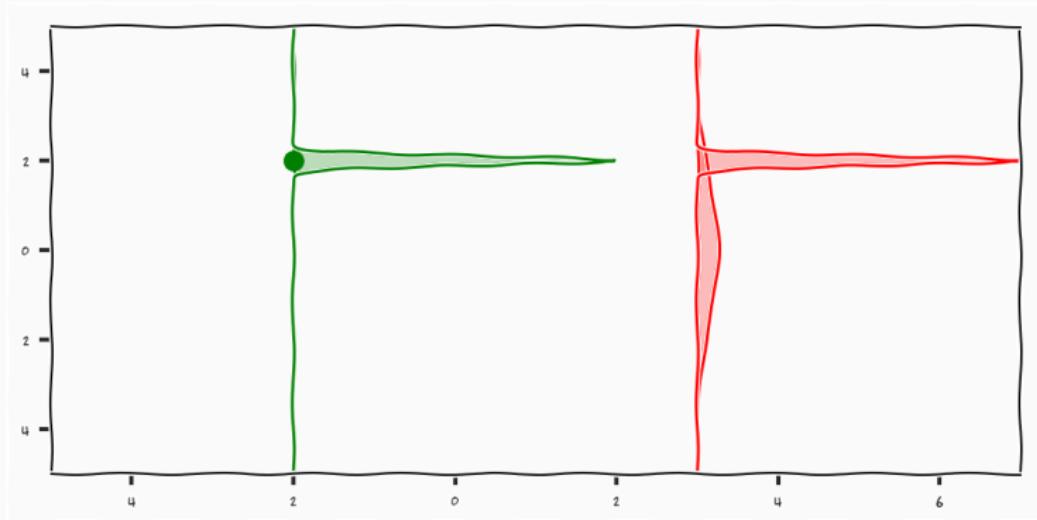
Gaussian Processes



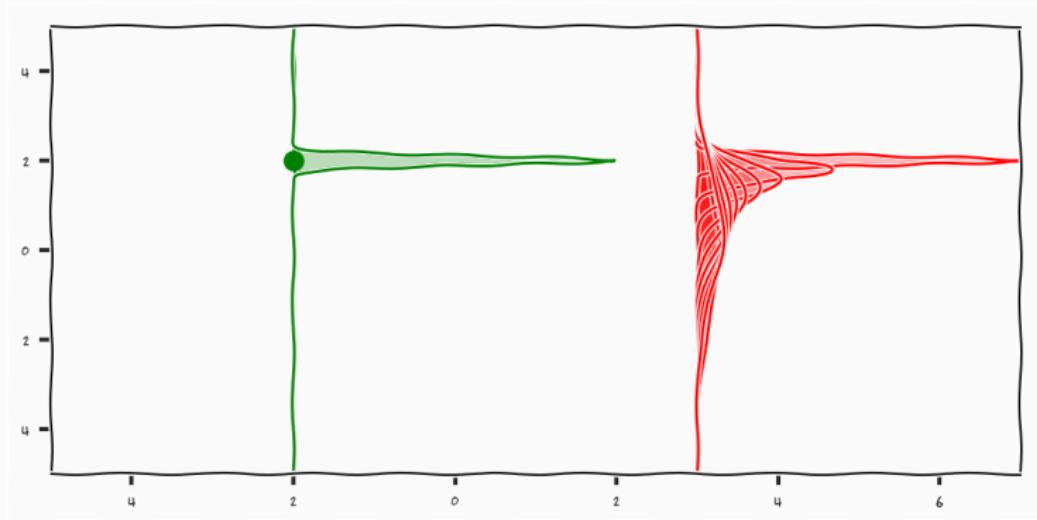
Gaussian Processes



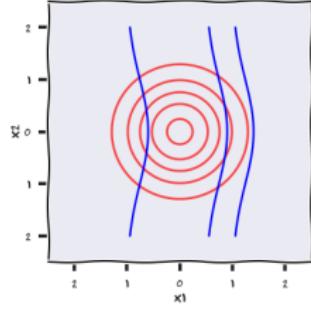
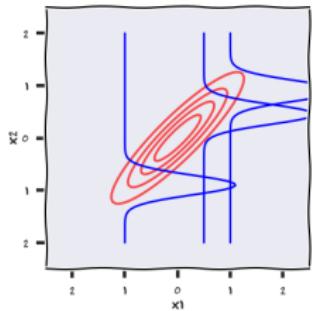
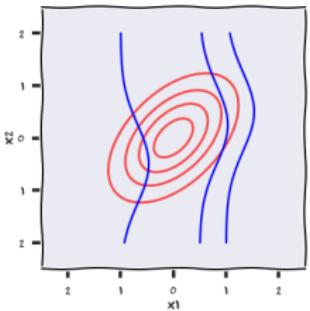
Gaussian Processes



Gaussian Processes



Conditional Gaussians

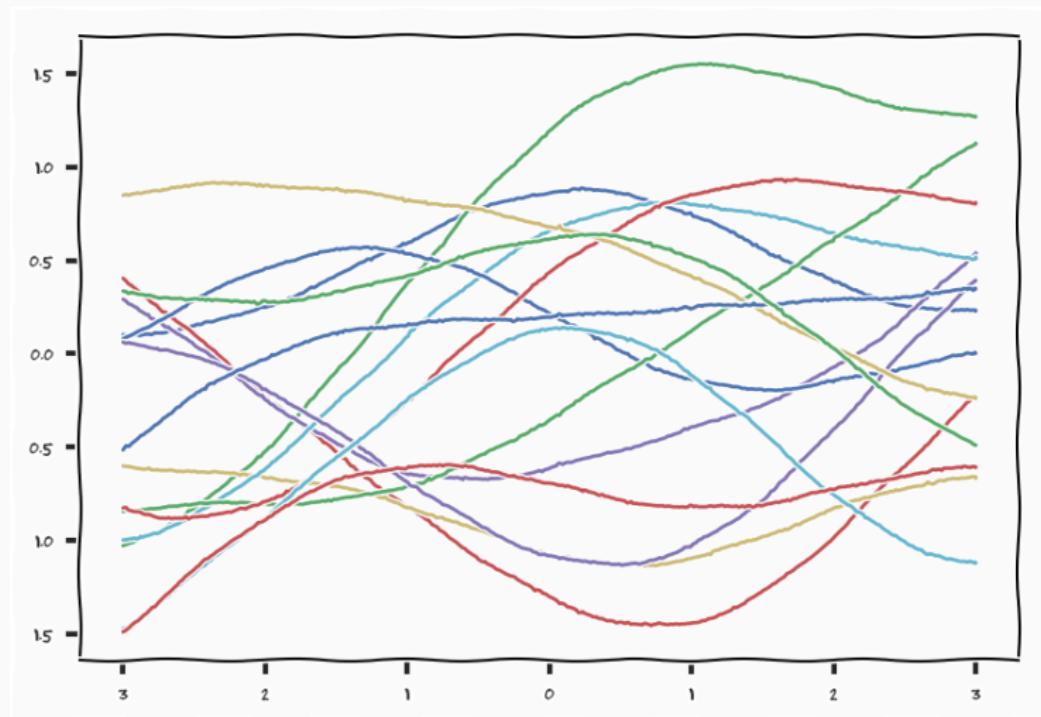


$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$$

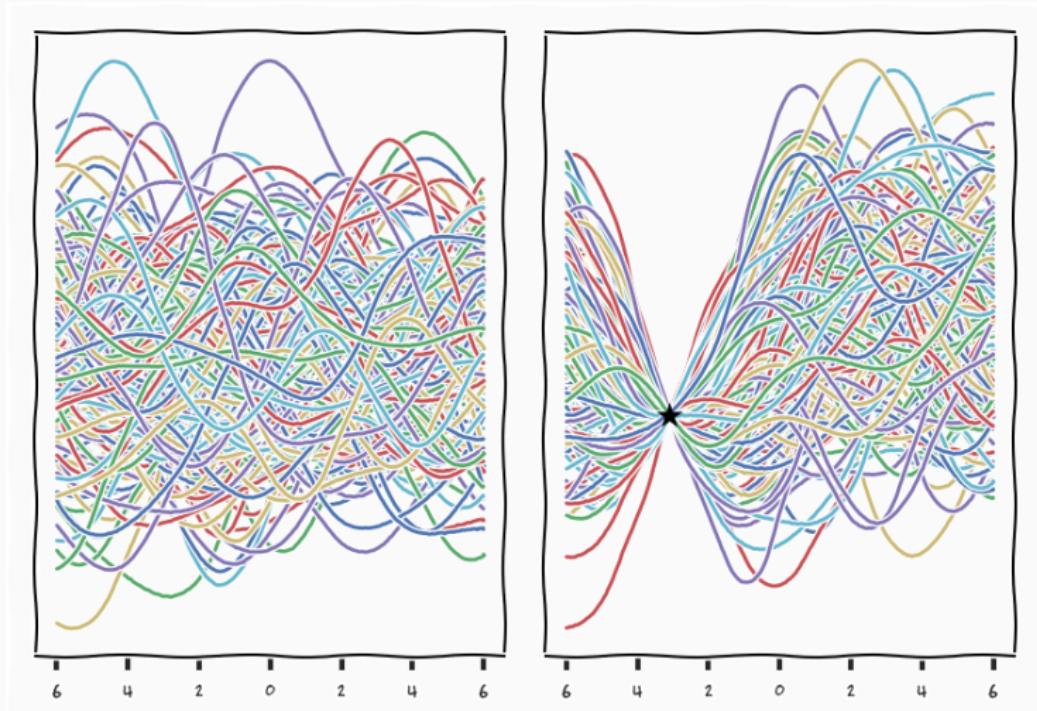
$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}\right)$$

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

Gaussian Processes

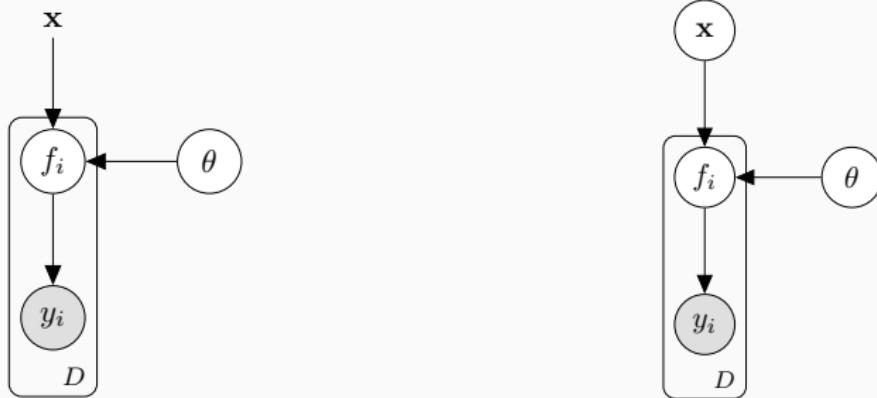


Gaussian Processes



Unsupervised Learning

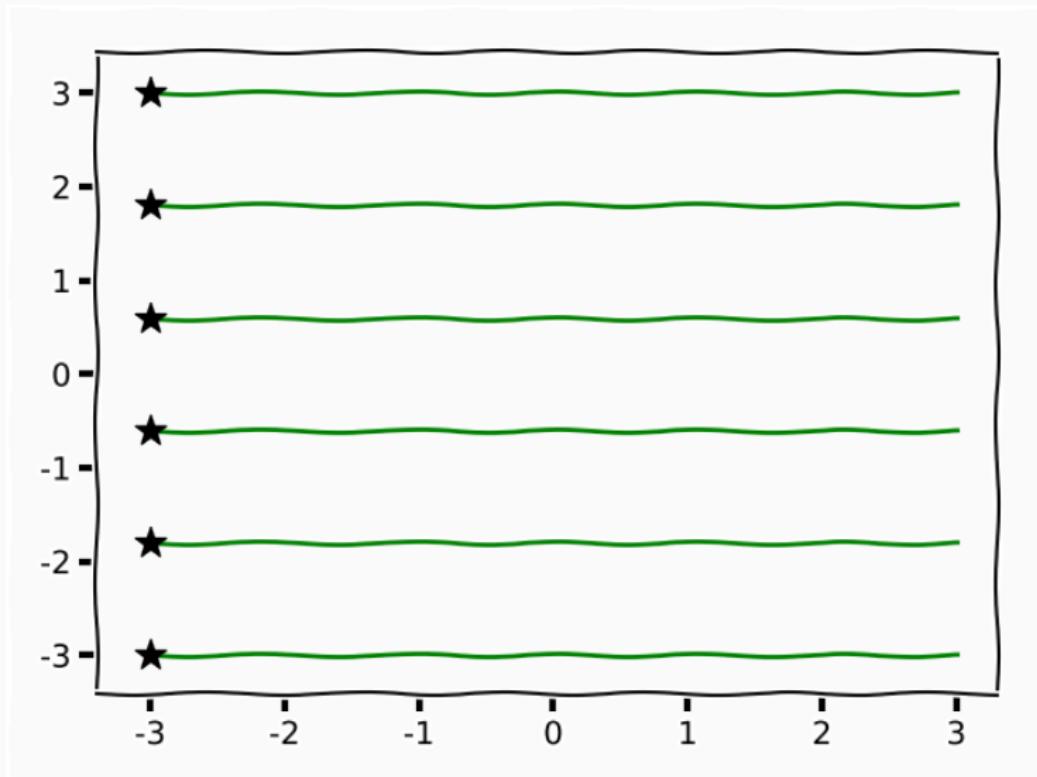
Unsupervised Learning



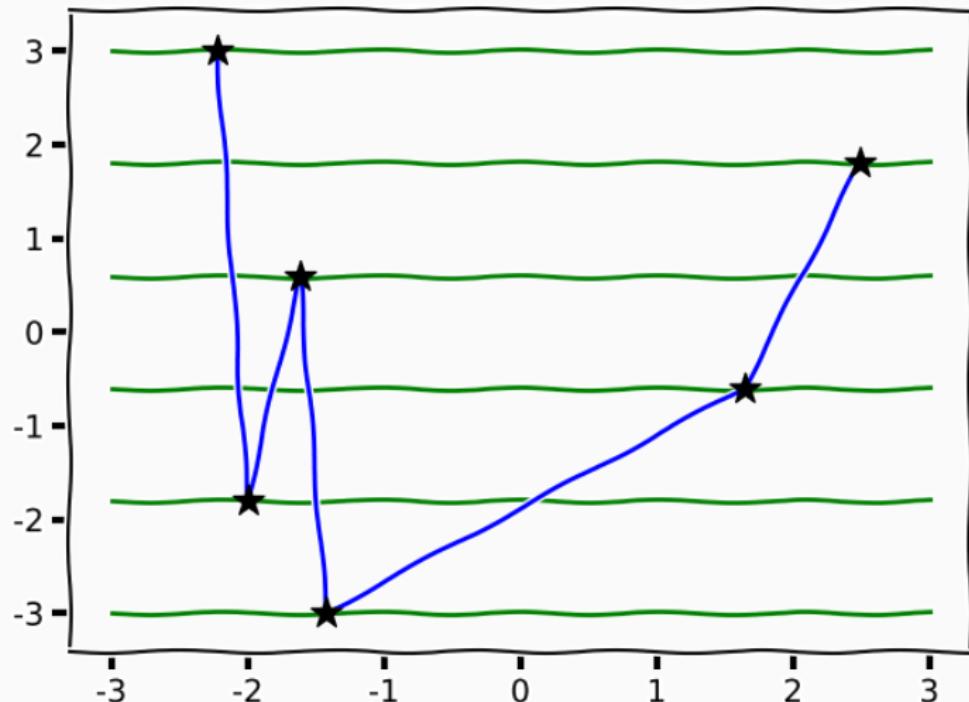
$$p(y|x)$$

$$p(y)$$

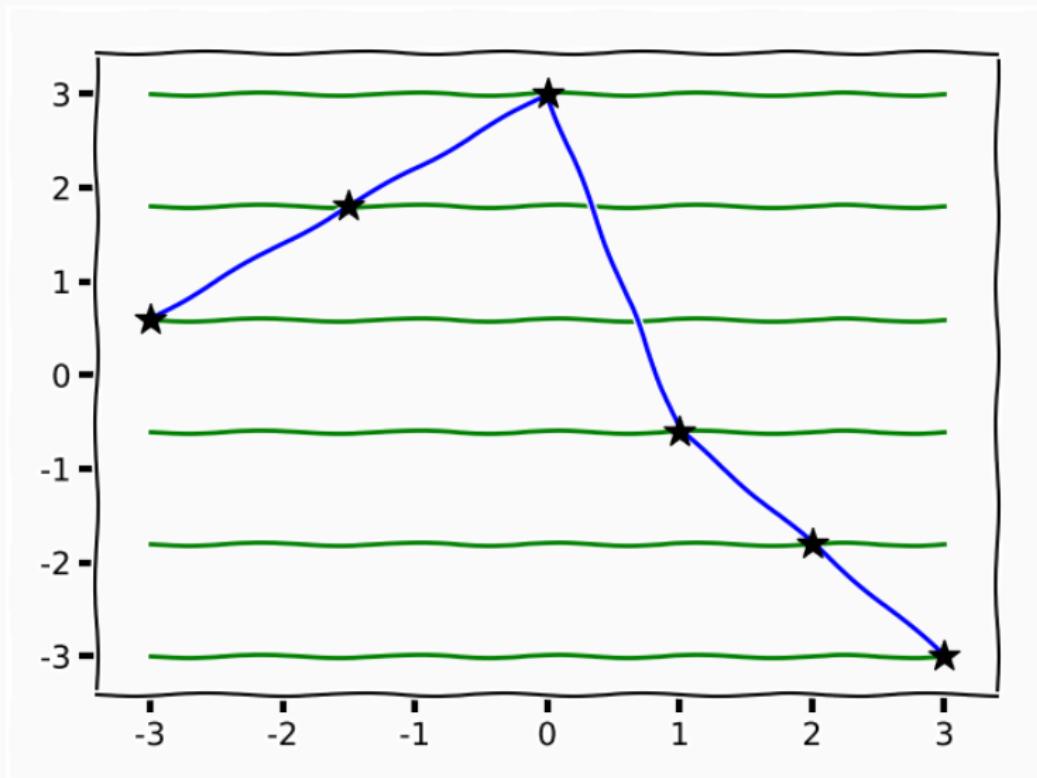
Unsupervised Learning



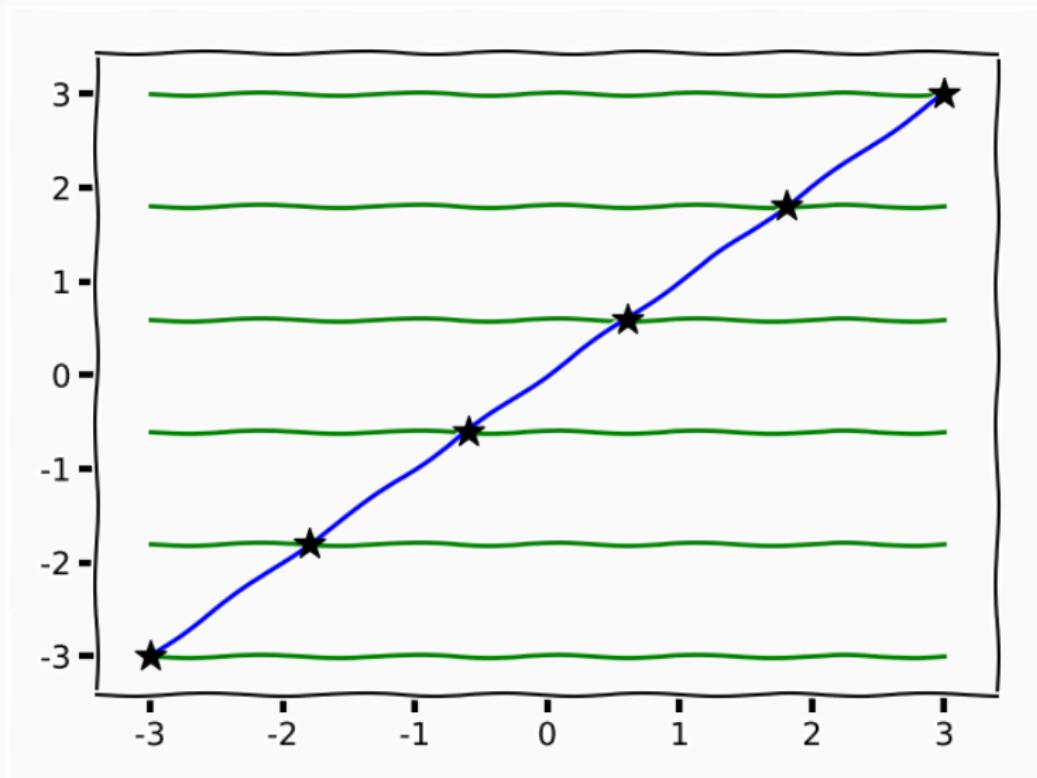
Unsupervised Learning



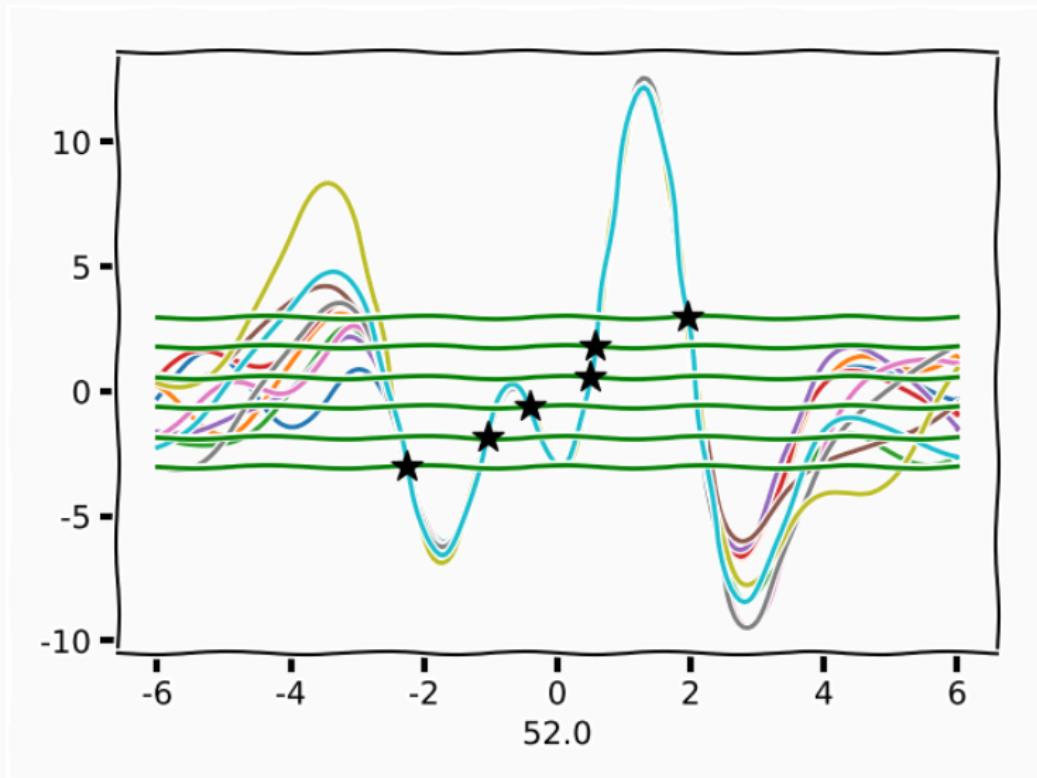
Unsupervised Learning



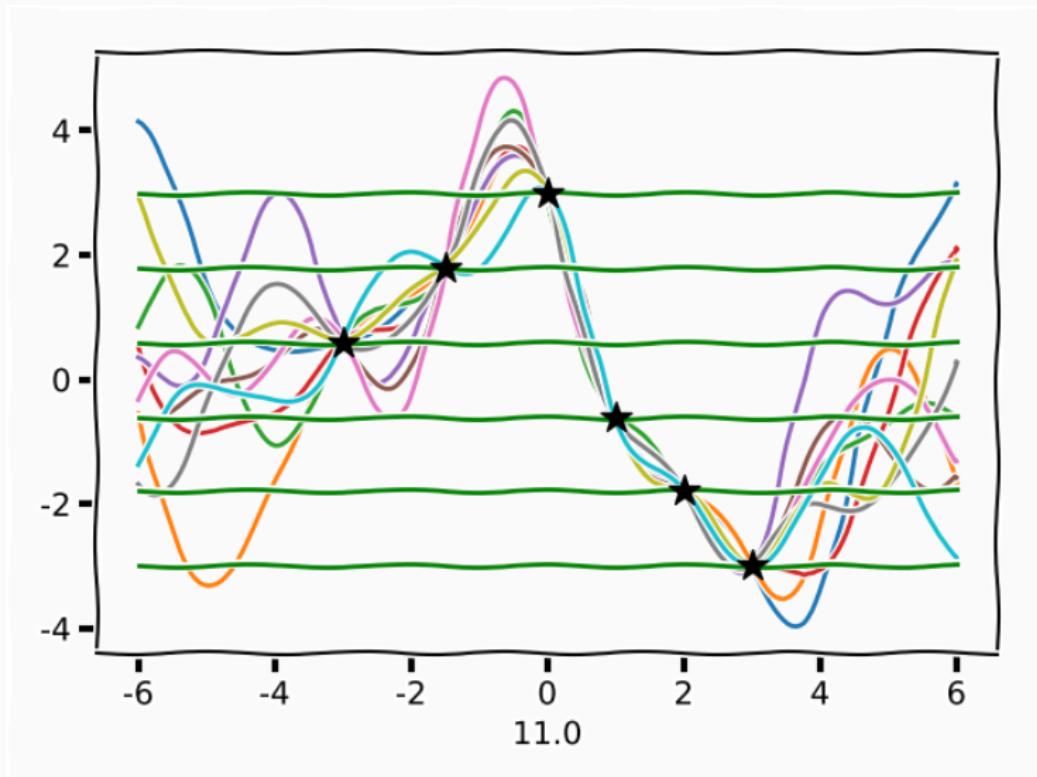
Unsupervised Learning



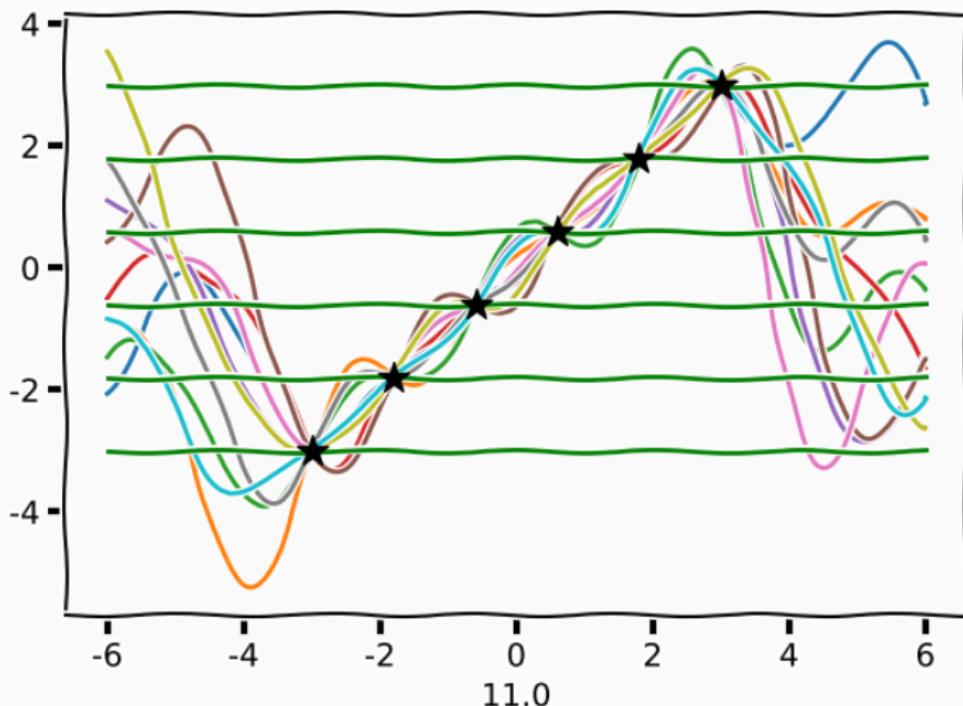
Unsupervised Learning



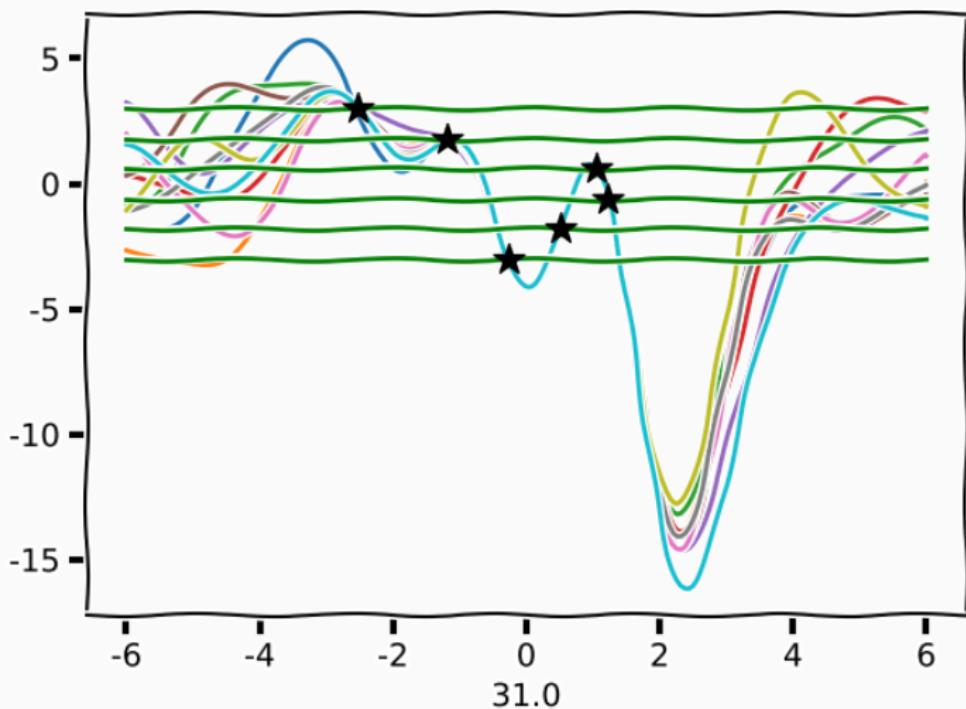
Unsupervised Learning



Unsupervised Learning



Unsupervised Learning



Priors



Priors

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

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

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. 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*

Priors



Being Bayesian¹



¹By Dieric Bouts (circa 1420-1475) - The Yorck Project: 10.000 Meisterwerke der Malerei, Public Domain, URL

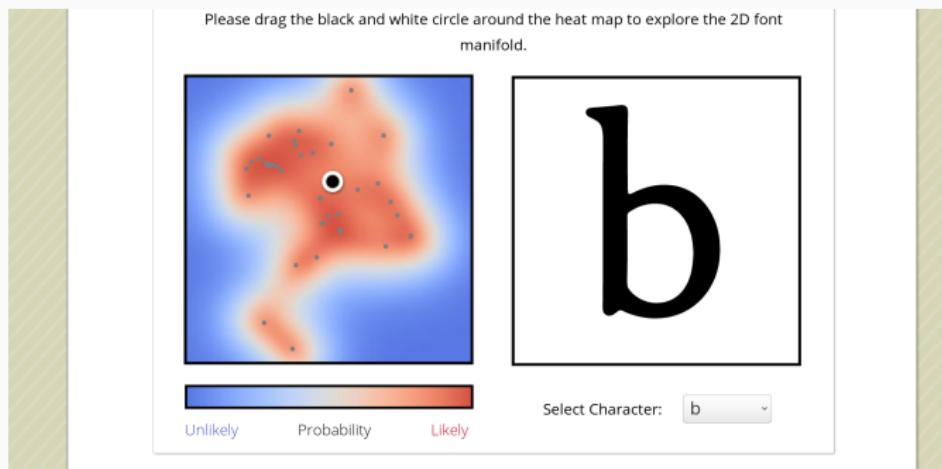
Unsupervised Learning with GPs

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

²Lawrence, 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

Font Demo



URL

Bayesian GP-LVM⁴

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

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

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

Bayesian GP-LVM⁴

- 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!³

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

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

Variational Bayes

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

Variational Bayes

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

Variational Bayes

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

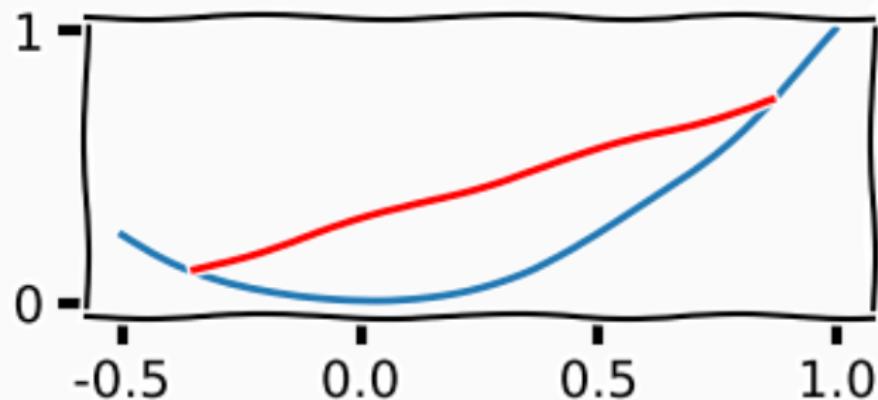
Variational Bayes

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

Variational Bayes

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

Jensen Inequality



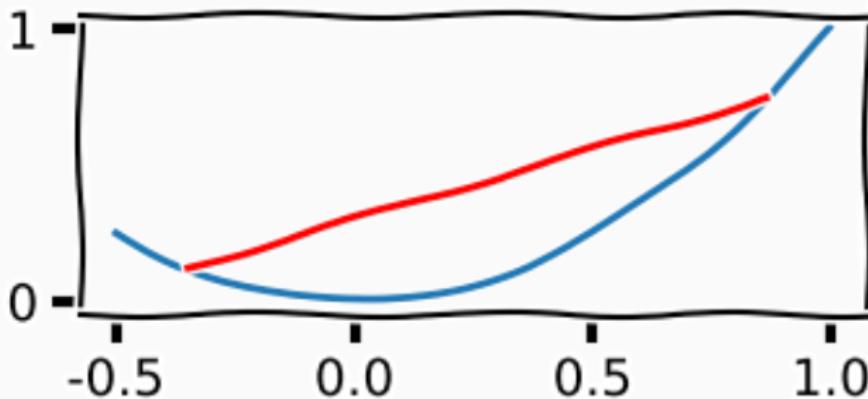
Convex Function

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

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

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

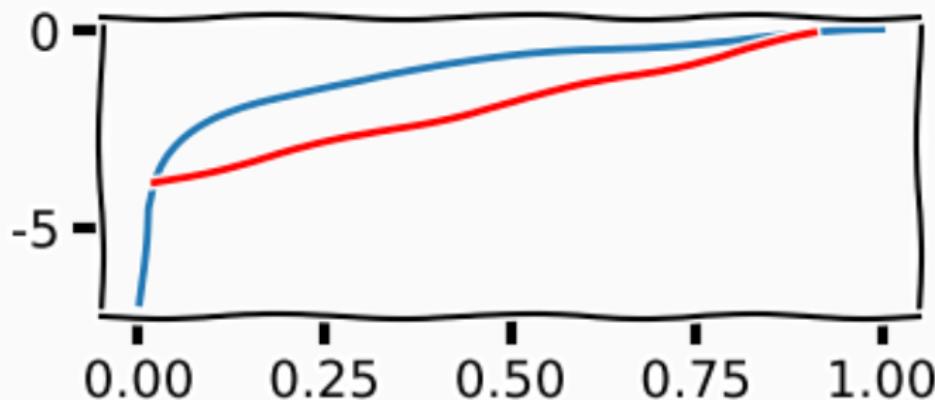
Jensen Inequality



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

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

Jensen Inequality in Variational Bayes



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

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

Variational Bayes cont.

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

Variational Bayes cont.

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

Variational Bayes cont.

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

Variational Bayes cont.

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

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

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. $\operatorname{argmin}_q \text{KL}(q(\mathbf{X}) || p(\mathbf{X}|\mathbf{Y}))$
⇒ variational distributions are approximations to intractable posteriors

ELBO

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

ELBO

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

ELBO

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

ELBO

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

ELBO

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

ELBO

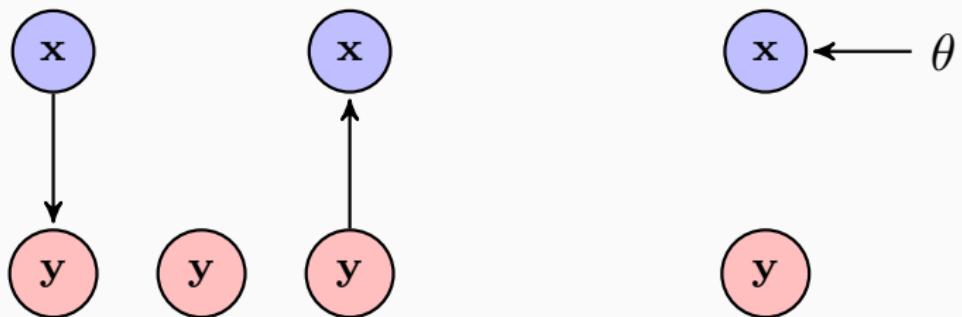
$$\begin{aligned} \log p(\mathbf{Y}) &= \text{KL}(q(\mathbf{X}) || p(\mathbf{X}|\mathbf{Y})) + \underbrace{\mathbb{E}_{q(\mathbf{X})} [\log p(\mathbf{X}, \mathbf{Y})] - H(q(\mathbf{X}))}_{\text{ELBO}} \\ &\geq \mathbb{E}_{q(\mathbf{X})} [\log p(\mathbf{X}, \mathbf{Y})] - H(q(\mathbf{X})) = \mathcal{L}(q(\mathbf{X}))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

ELBO

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

- 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

ELBO



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

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

Why is this useful?

Why is this a sensible thing to do?

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

– Ryan Adams⁵

⁵Talking 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 Adams⁵

⁵Talking 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 Adams⁵

⁵Talking Machines Season 2, Episode 5

Lower Bound⁶

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

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

Lower Bound⁶

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

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

Lower Bound⁶

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

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

Lower Bound⁶

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

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

Lower Bound

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

- Has not eliviate the problem at all, X still needs to go through F to reach the data
- Idea of sparse approximations⁷

⁷Quinonero-Candela, Joquin, & Rasmussen, C. E. (2005). A unifying view of sparse approximate Gaussian process regression & Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs

Lower Bound

- Add another set of samples from the same prior

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

Lower Bound

- 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

$$p(\mathbf{f}_{:,j}, \mathbf{u}_{:,j} | \mathbf{X}, \mathbf{Z}) = p(\mathbf{f}_{:,j} | \mathbf{u}_{:,j}, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_{:,j} | \mathbf{Z})$$

$$= \mathcal{N}(\mathbf{f}_{:,j} | \mathbf{K}_{fu}(\mathbf{K}_{uu})^{-1} \mathbf{u}_{:,j}, \mathbf{K}_{ff} - \mathbf{K}_{fu}(\mathbf{K}_{uu})^{-1} \mathbf{K}_{uf}) \mathcal{N}(\mathbf{u}_{:,j} | \mathbf{0}, \mathbf{K}_{uu}),$$

Lower Bound

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

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

Lower Bound

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

- 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

Lower Bound

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

- 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)

Lower Bound

- Variational distributions are approximations to intractable posteriors,

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

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

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

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} ,

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

Lower Bound

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

Lower Bound

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

Lower Bound

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

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

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

Lower Bound

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

Lower Bound

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

Summary

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

- 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

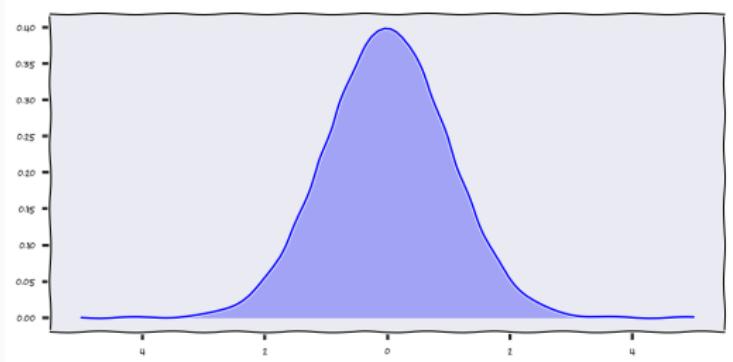
Latent space priors⁸

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

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

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

The Gaussian blob



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

Automatic Relevance Determination

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

GPy

RBF(..., ARD=True)

Matern32(..., ARD=True)

Change of Variables

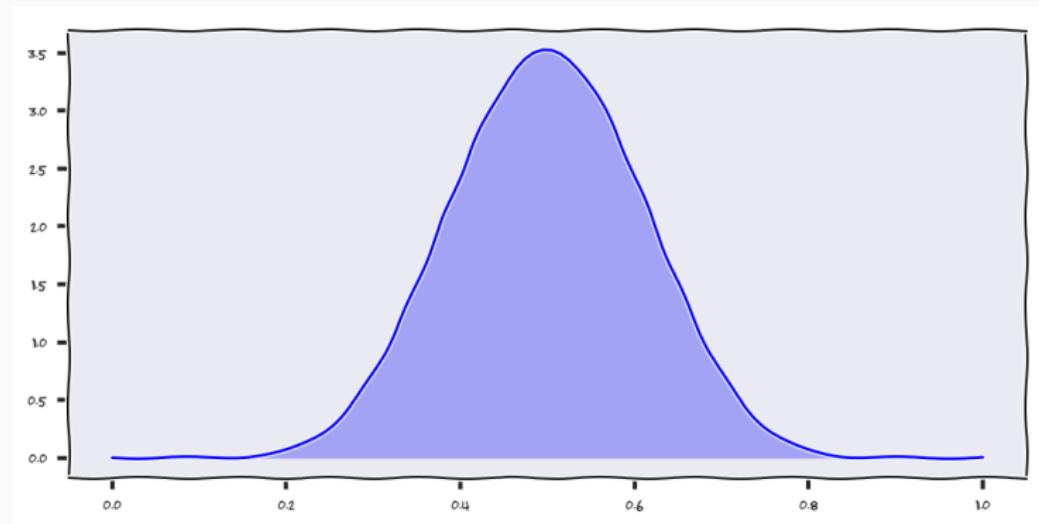
Theorem (Change of Variable)

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

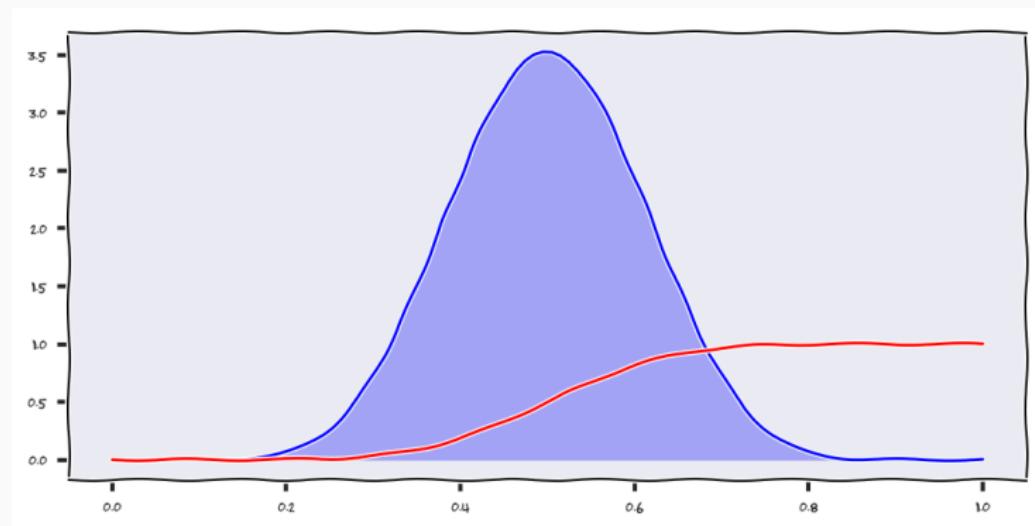
$$p_y(y) = p_x(\psi(y)) |\nabla \psi(y)|$$

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

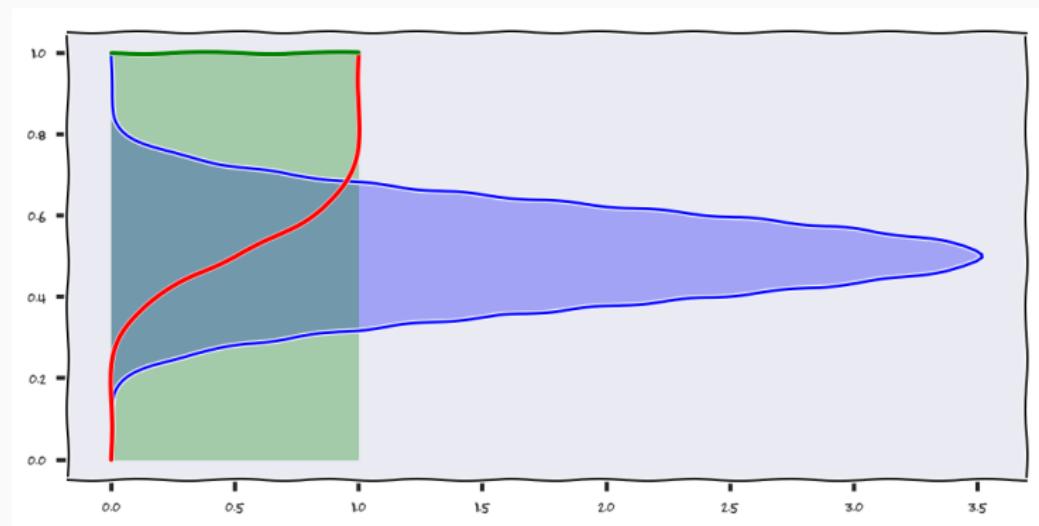
Sampling



Sampling



Sampling



Change of Variables

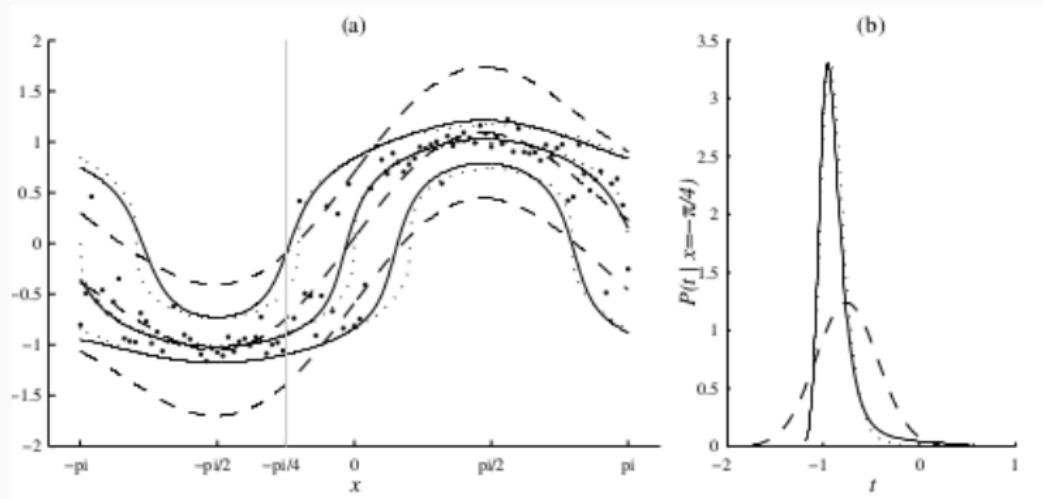
Theorem (Change of Variable)

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.

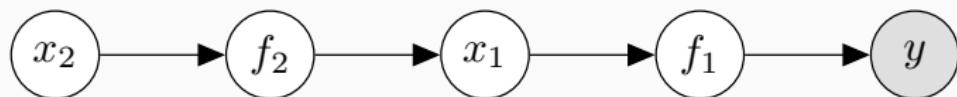
Warped Gaussian Processes^{9, 10}



⁹Snelson, E., & Ghahramani, Z. (2004). Warped Gaussian Processes

¹⁰Lazaro-Gredilla, Miguel (2012). Bayesian Warped Gaussian Processes. In , Advances in Neural Information Processing Systems

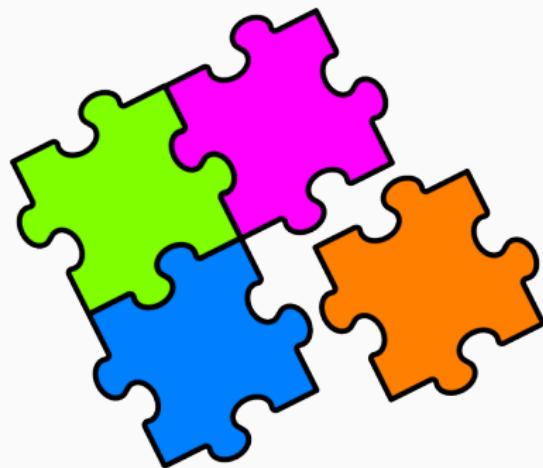
Deep Gaussian Processes¹¹



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

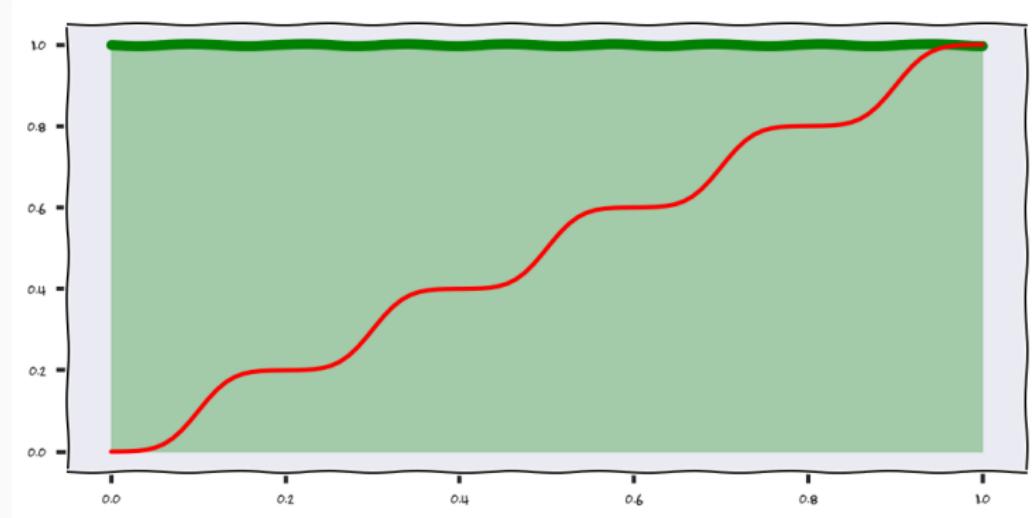
¹¹Damianou, A. C., & Lawrence, N. D. (2013). Deep Gaussian Processes

Composite Functions

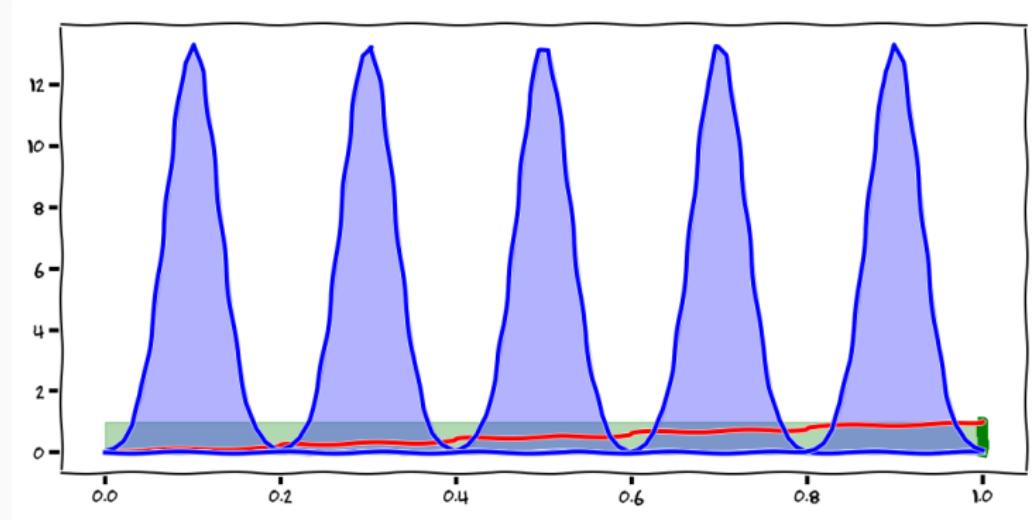


$$y = f_k(f_{k-1}(\dots f_0(x))) = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$$

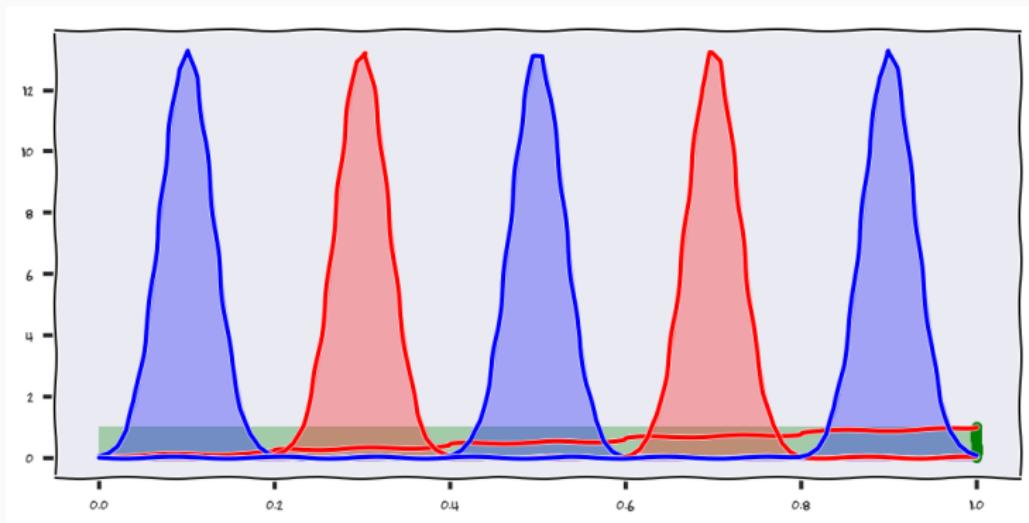
Change of Variables



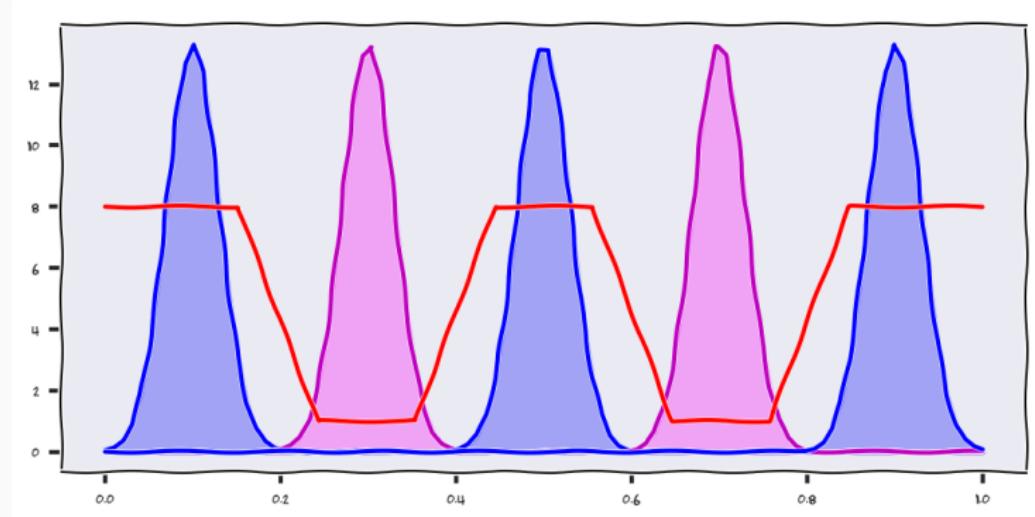
Change of Variables



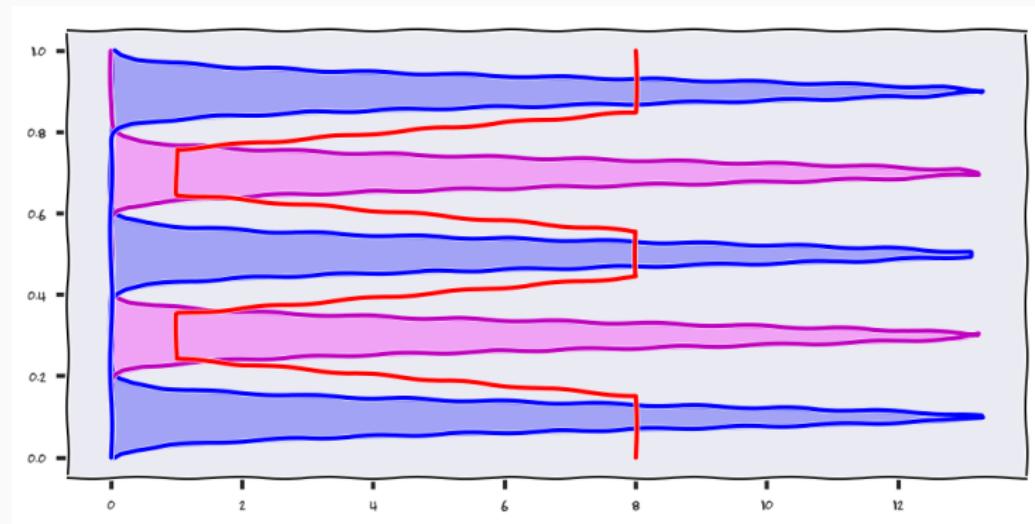
Change of Variables



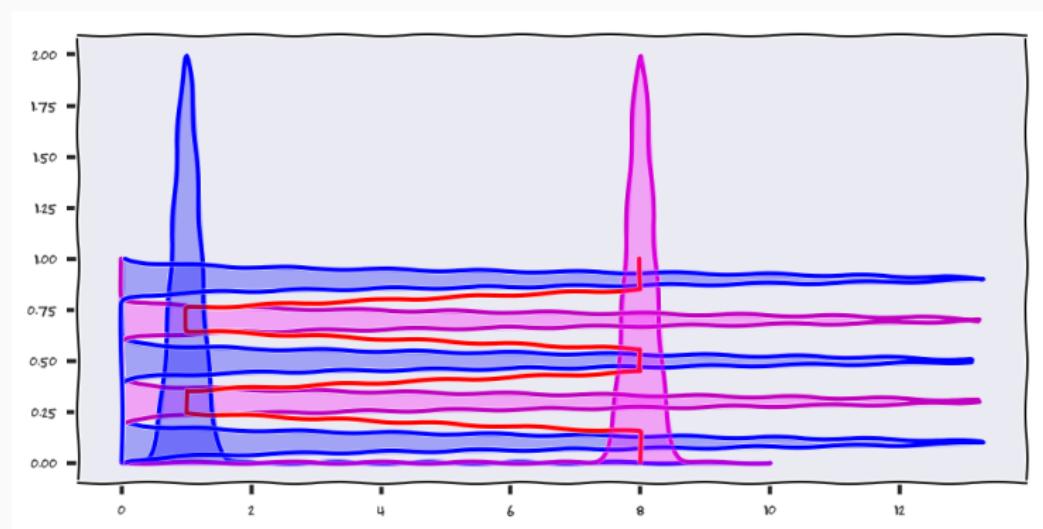
Change of Variables



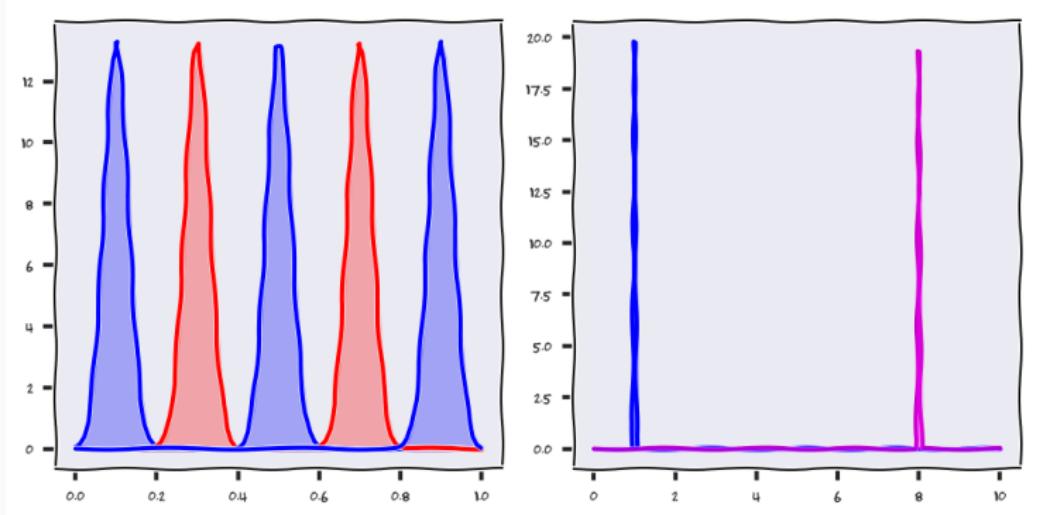
Change of Variables



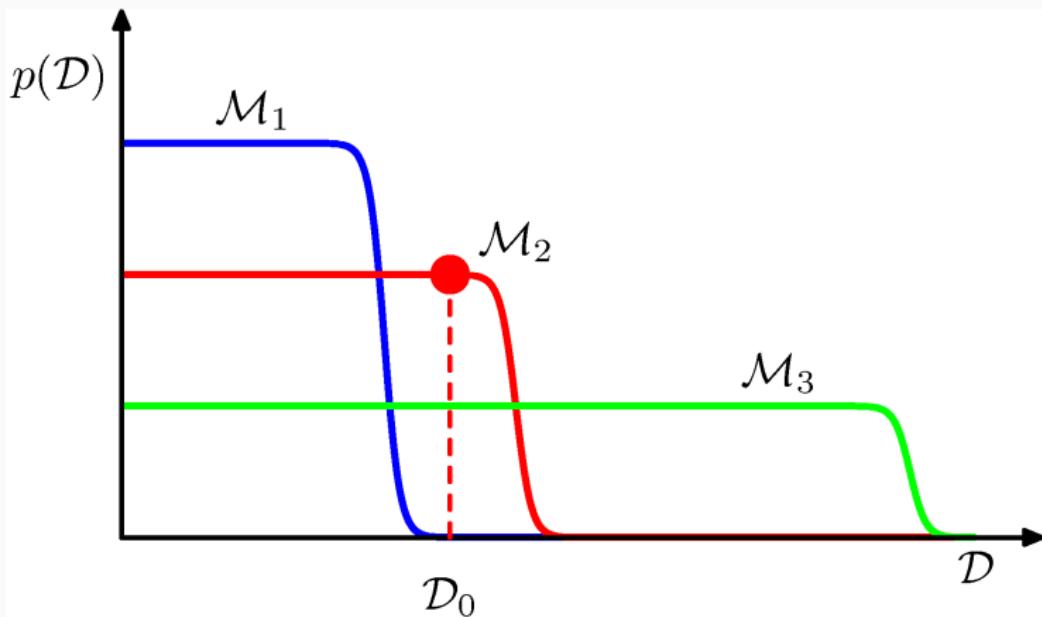
Change of Variables



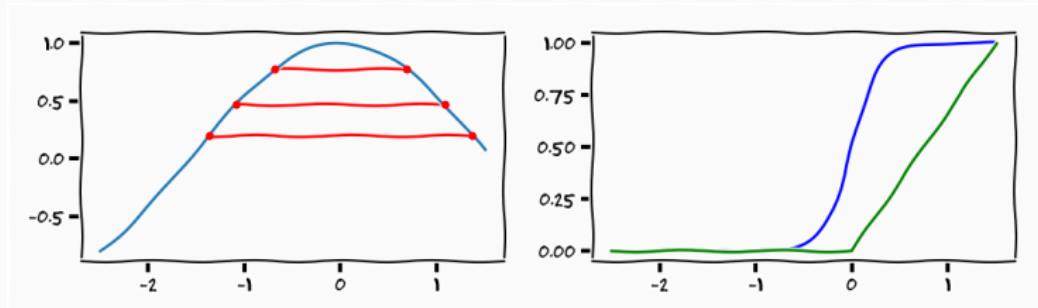
Change of Variables



MacKay plot



Composition functions

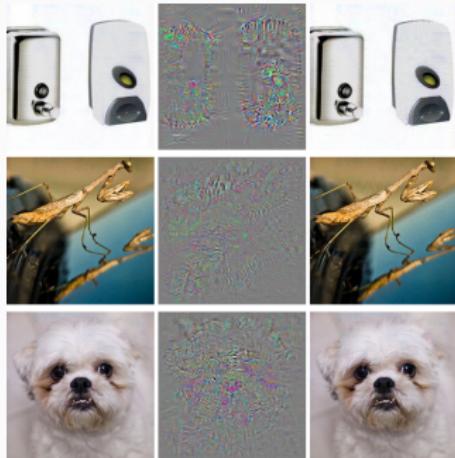
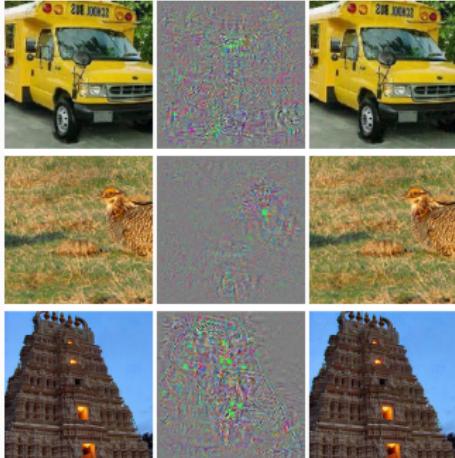


$$y = f_k(f_{k-1}(\dots f_0(x))) = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$$

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

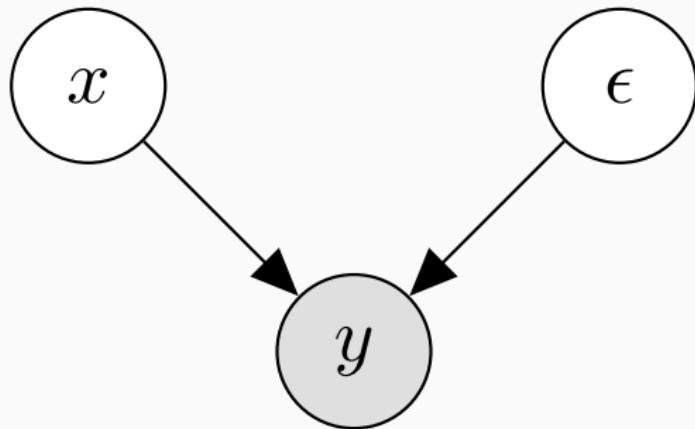
$$\text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2) \subseteq \dots \subseteq \text{Im}(f_k)$$

Data inefficiency¹²



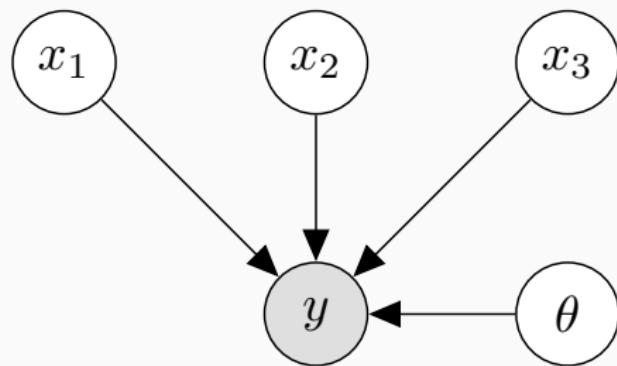
¹²Nguyen, A. M., Yosinski, J., & Clune, J., Deep neural networks are easily fooled: high confidence predictions for unrecognizable images, CoRR, abs/1412.1897(), (2014).

Explaining Away



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

Factor Analysis

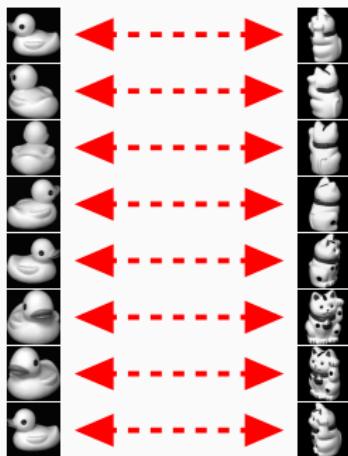


$$y = f(x_1, x_2, x_3) + \epsilon$$

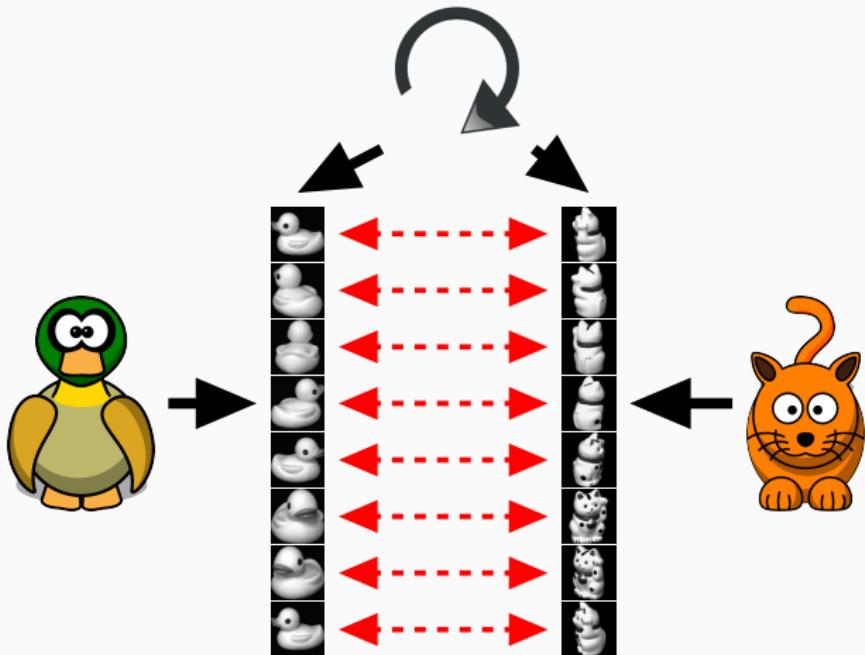
Alignments



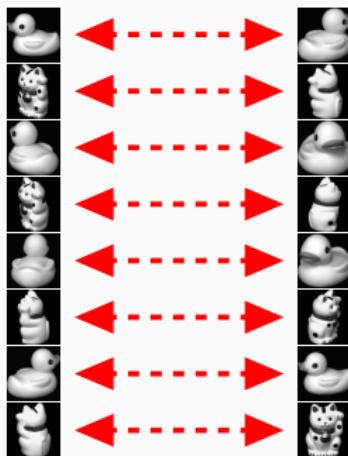
Alignments



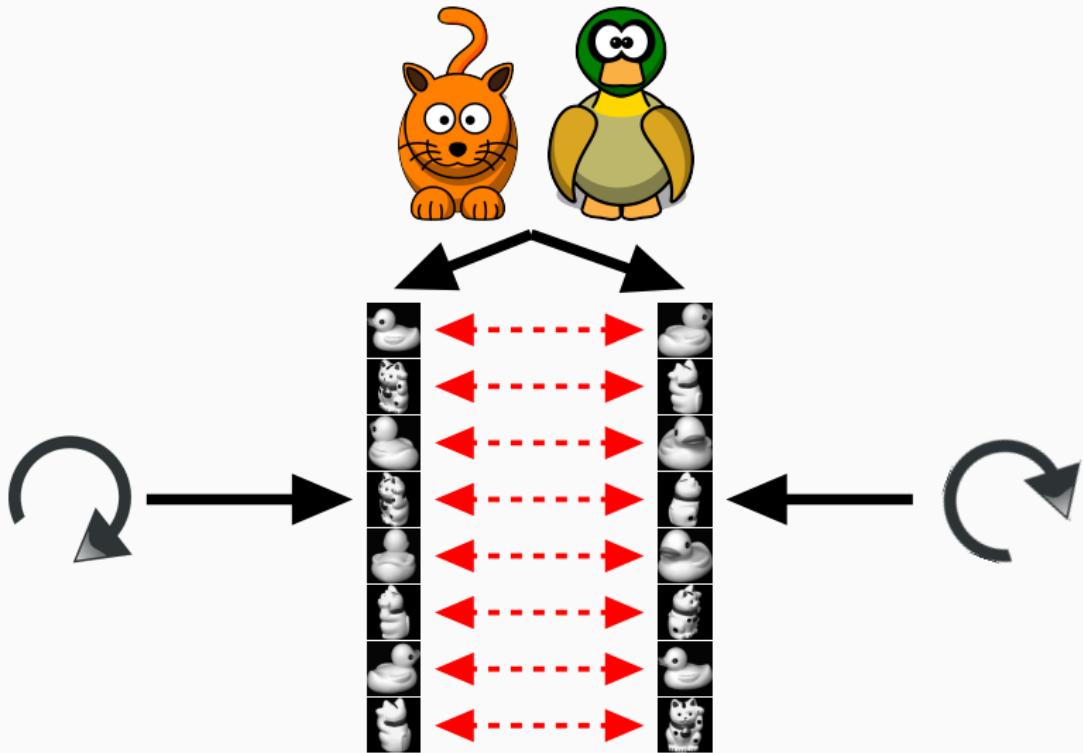
Alignments



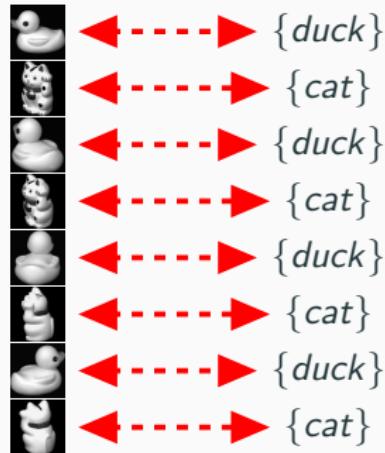
Alignments



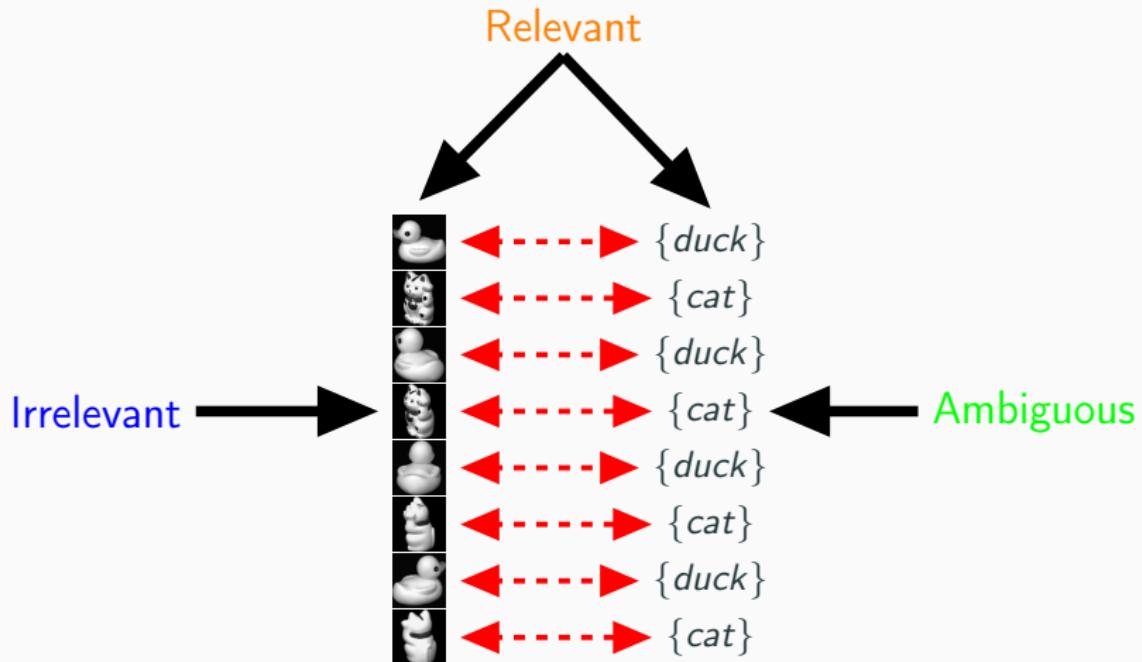
Alignments



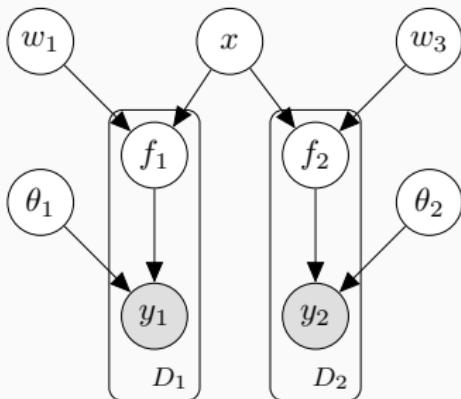
Alignments



Alignments



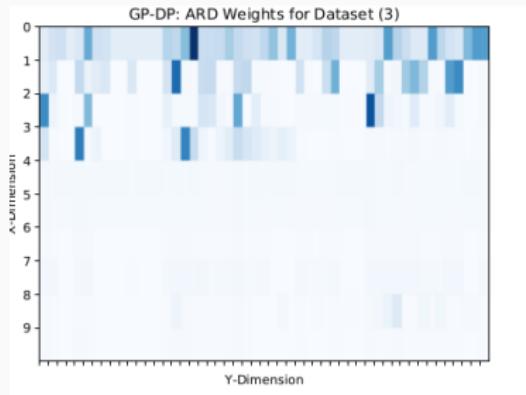
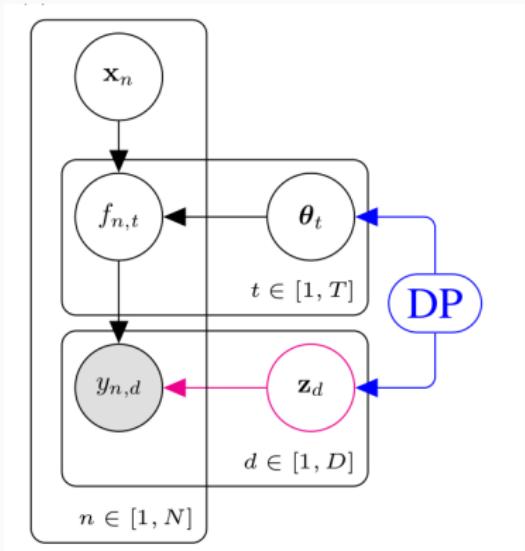
IBFA with GP-LVM¹³



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

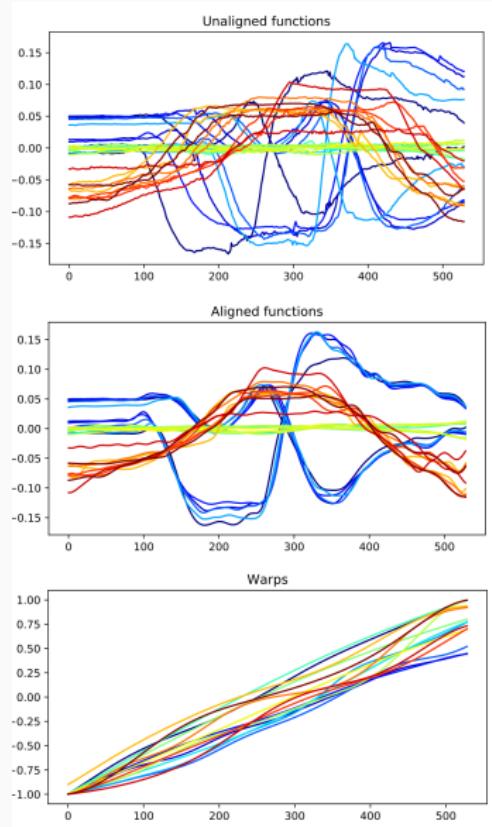
¹³Damianou, A., Lawrence, N. D., & Ek, C. H. (2016). Multi-view learning as a nonparametric nonlinear inter-battery factor analysis

GP-DP¹⁴

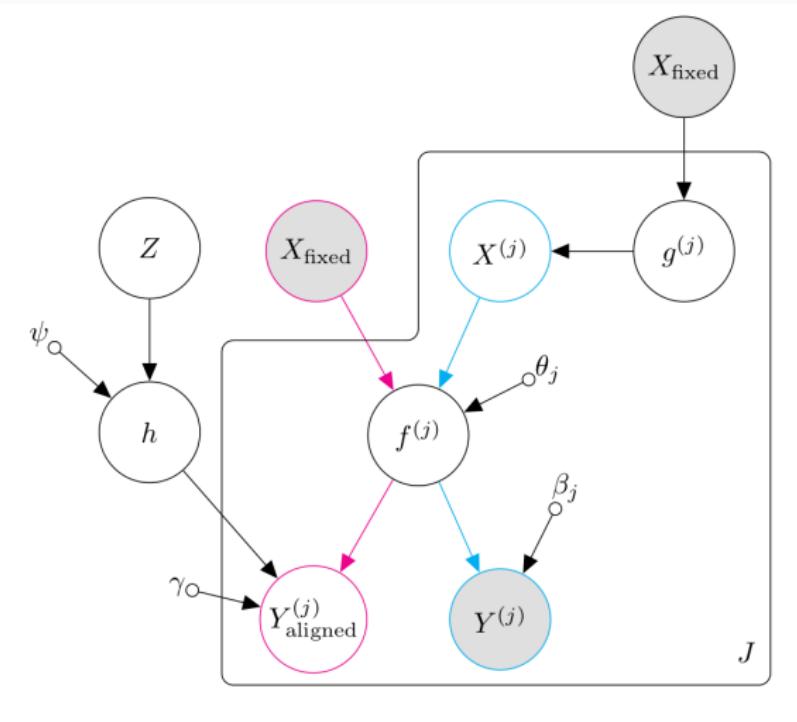


¹⁴ Lawrence, A. R., Ek, C. H., & Campbell, N. D. F., Dp-gp-lvm: a bayesian non-parametric model for learning multivariate dependency structures, CoRR, (), (2018).

Alignment Learning

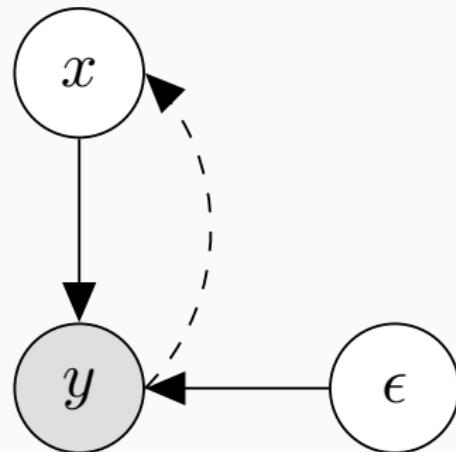


Alignment Learning¹⁵



¹⁵Kazlauskaitė, I., Ek, C. H., & Campbell, N. D. F., Gaussian Process Latent Variable Alignment Learning, CoRR, (), (2018).

Constrained Latent Space



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

Constrained Latent Space

- Dai, Z., Damianou, A., Gonzalez, Javier, & Lawrence, N., Variational auto-encoded deep Gaussian processes, International Conference on Learning Representations (ICLR), (2016).
- Snoek, J., Adams, R. P., & Larochelle, H., Nonparametric guidance of autoencoder representations using label information, Journal of Machine Learning Research, 13(), 2567–2588 (2012).
- Ek, C. H., Torr, P. H. S., & Lawrence, N. D., Gaussian process latent variable models for human pose estimation, International conference on Machine learning for multimodal interaction, (), 132–143 (2007).

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 such as 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

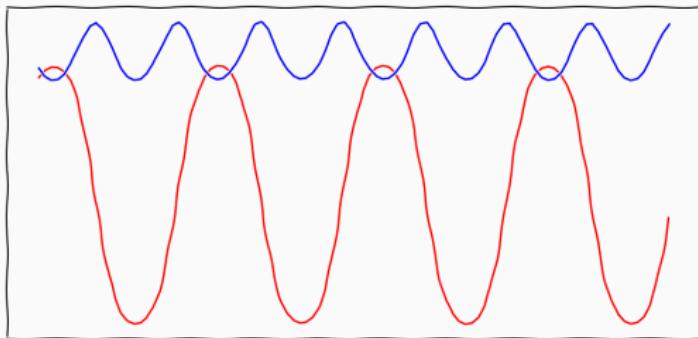
Summary II

- 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

Appendix

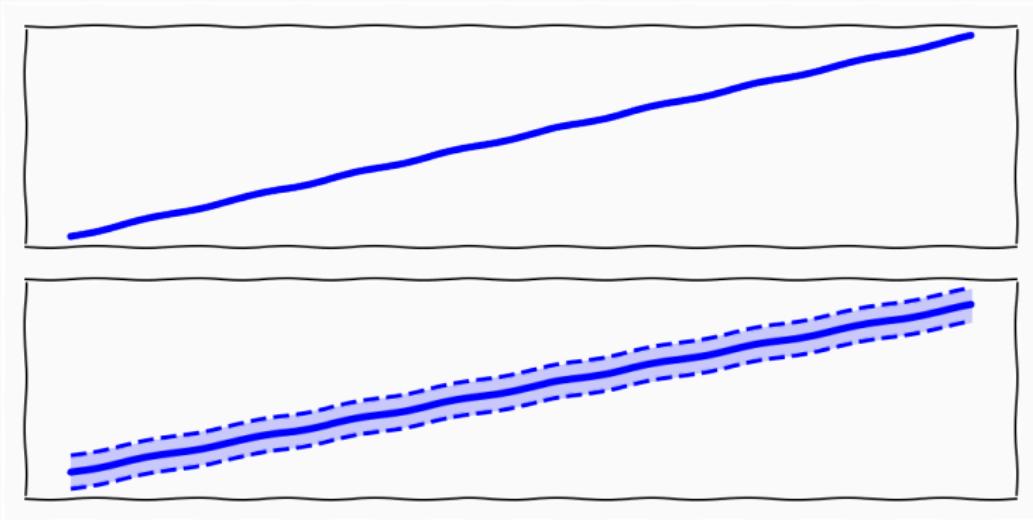
Composition: priors



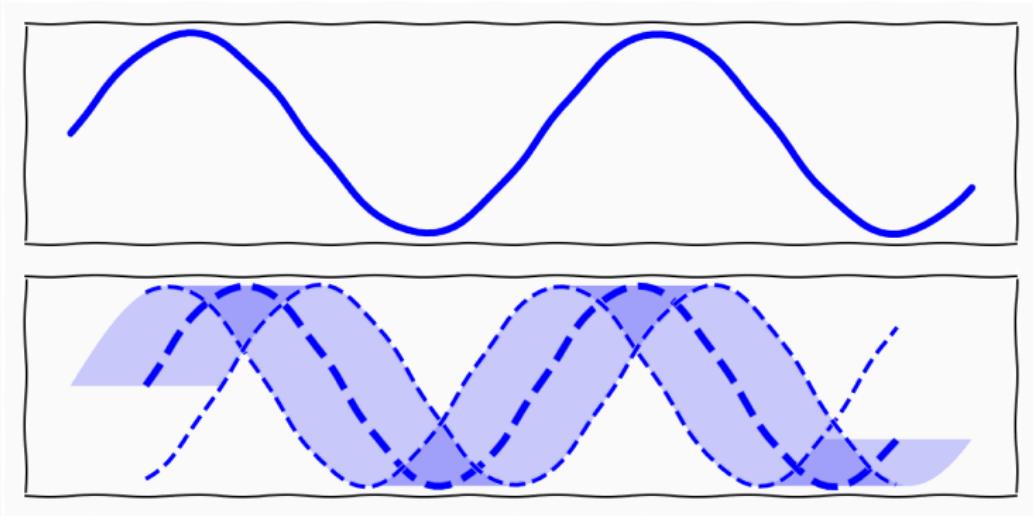
Composition: priors

`./bin/composition2.png`

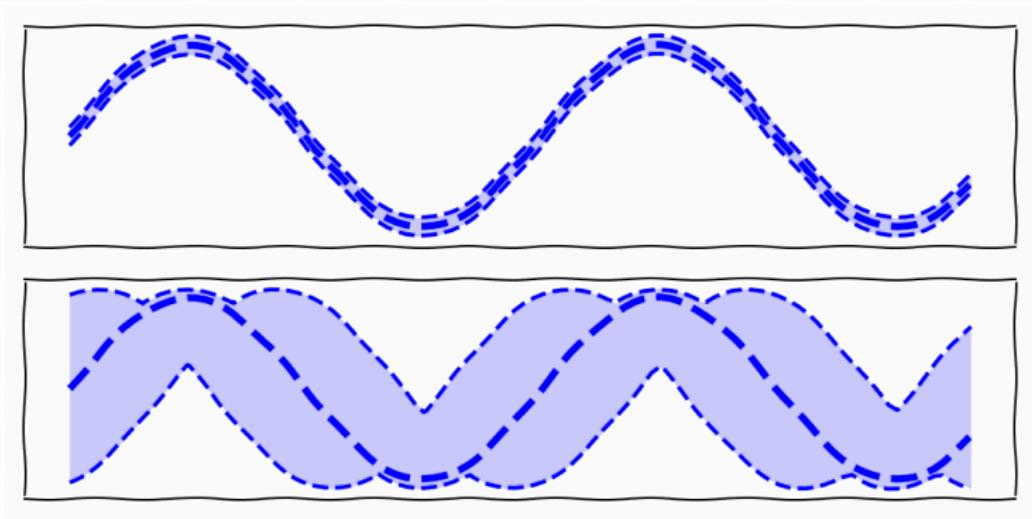
Composition: uncertainty



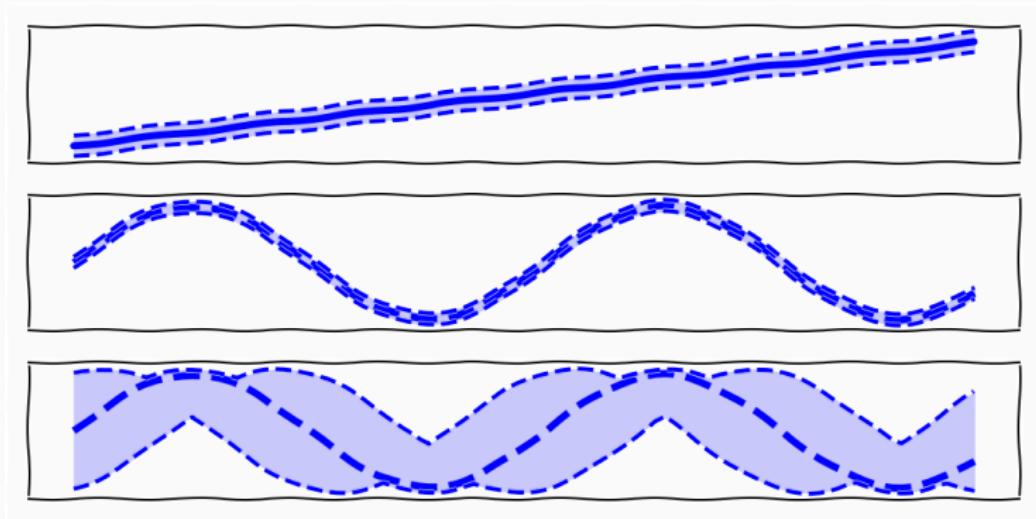
Composition: uncertainty



Composition: uncertainty



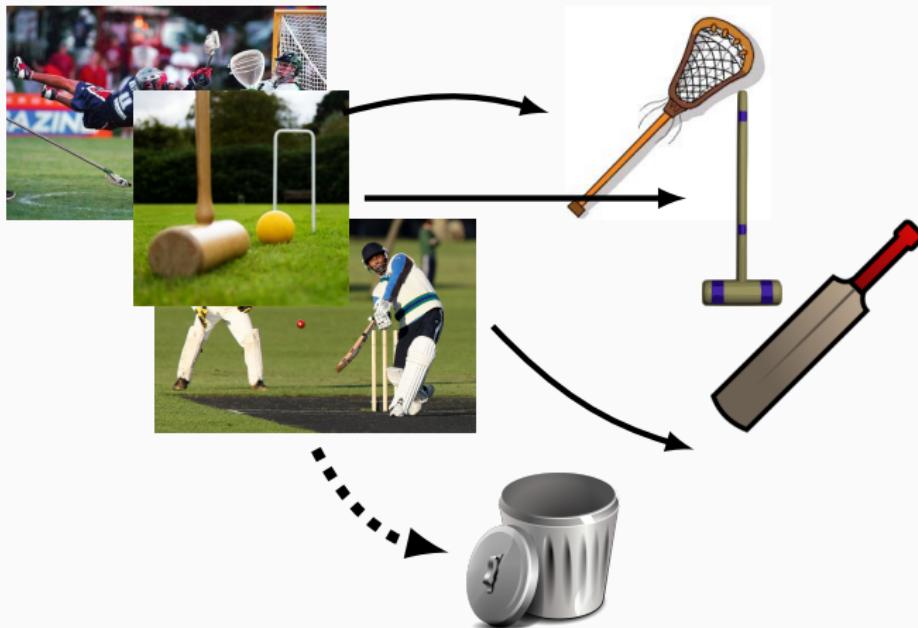
Composition: uncertainty



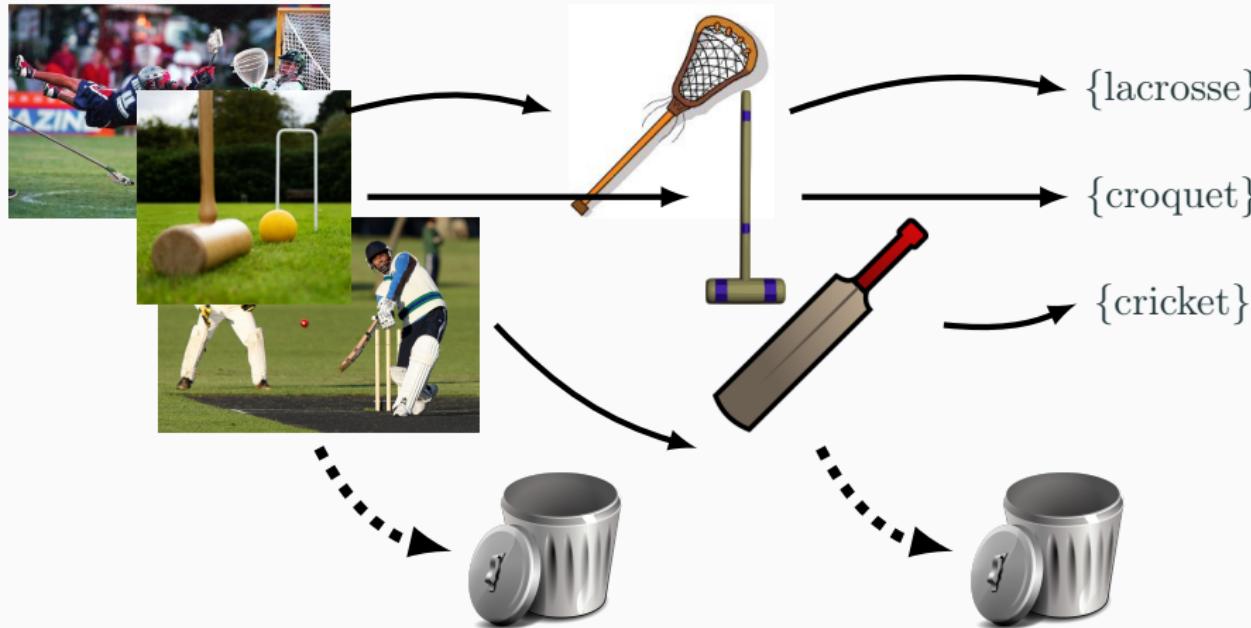
Latent Structure



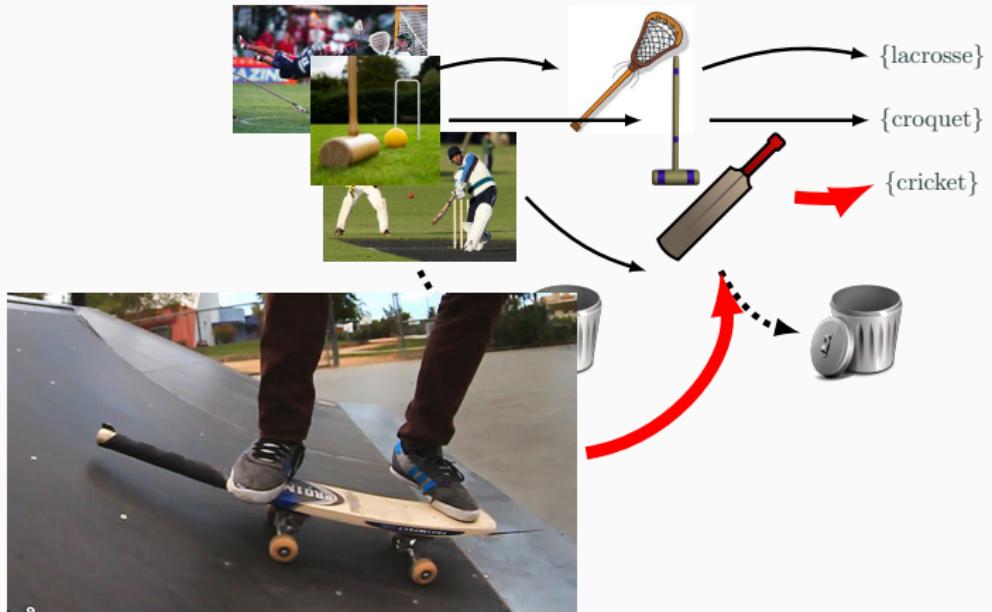
Latent Structure



Latent Structure

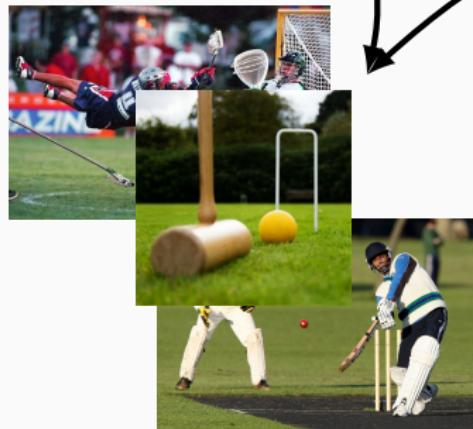


Latent Structure



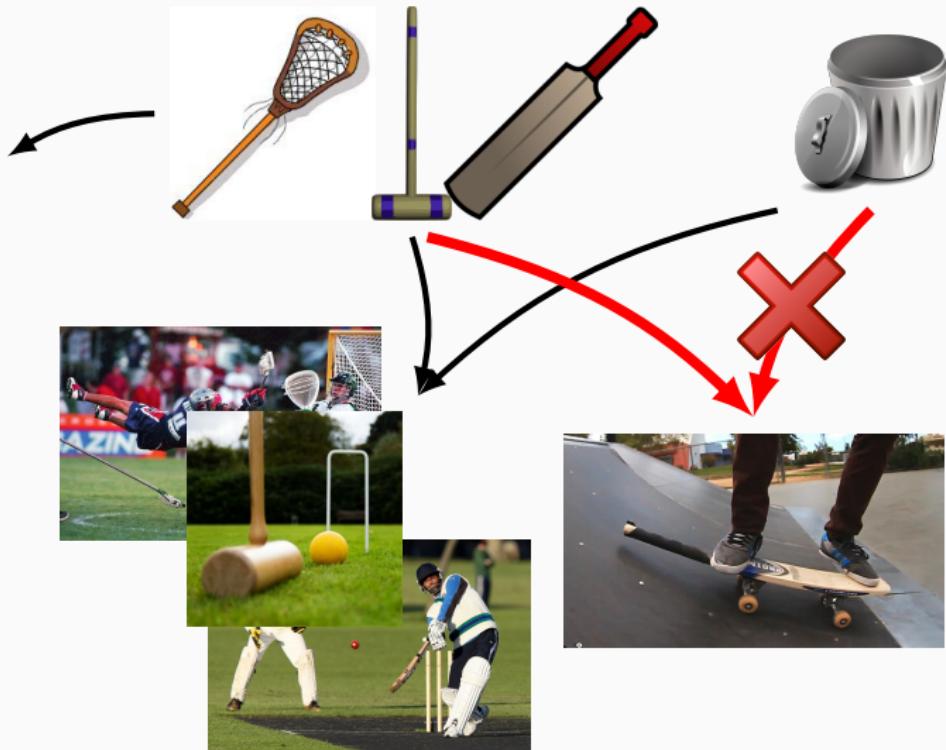
Latent Structure

{lacrosse}
{croquet}
{cricket}



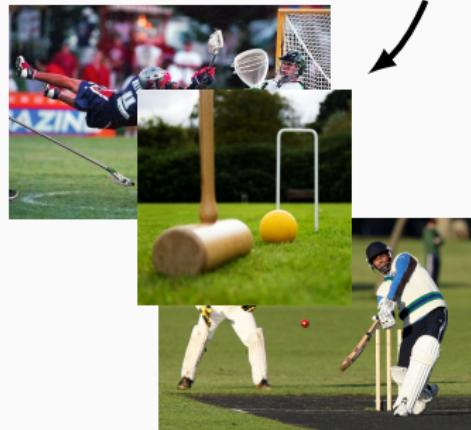
Latent Structure

{lacrosse}
{croquet}
{cricket}

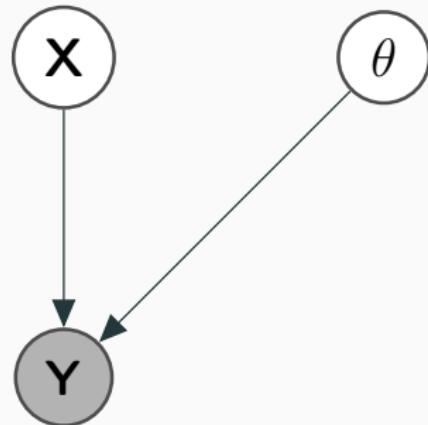


Latent Structure

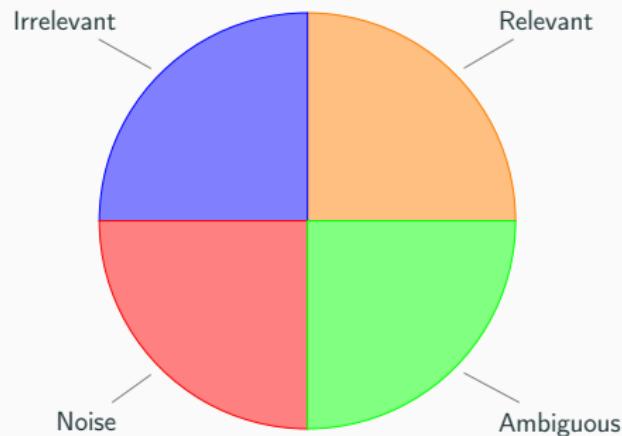
{lacrosse}
{croquet}
{cricket}



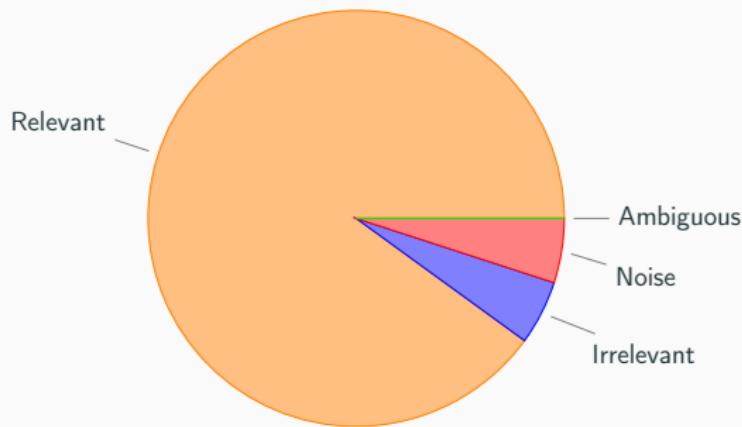
Latent Structure



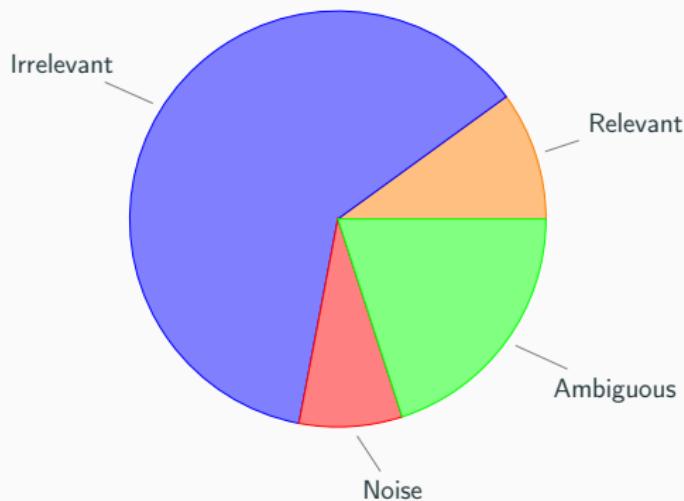
Latent Structure



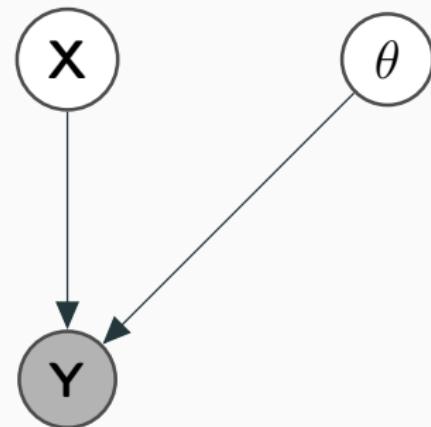
Latent Structure



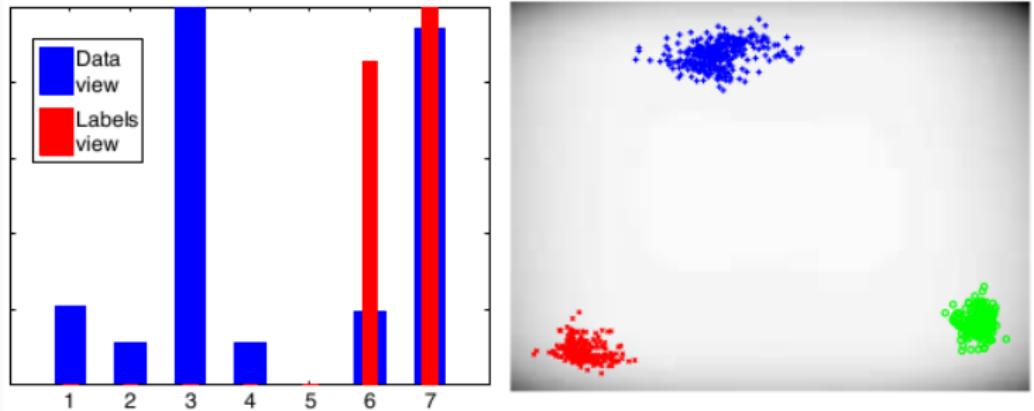
Latent Structure



Latent Structure



IBFA with GP-LVM



IBFA with GP-LVM

