

Machine Learning

Gaussian Processes II

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

October 10, 2017

<http://www.carlhenrik.com>

Introduction

- We have uncertainty in our observed outputs

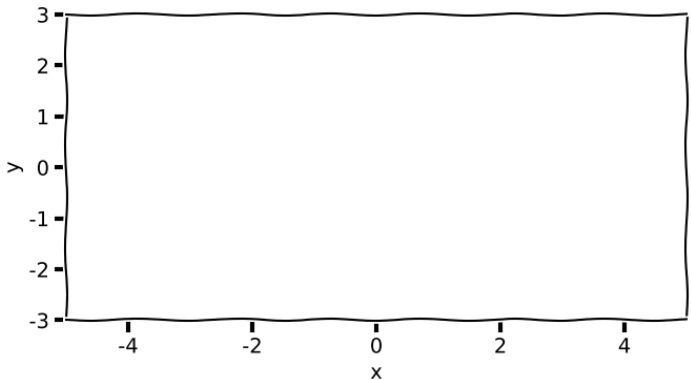
- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping

- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping
 - Linear, it is a line

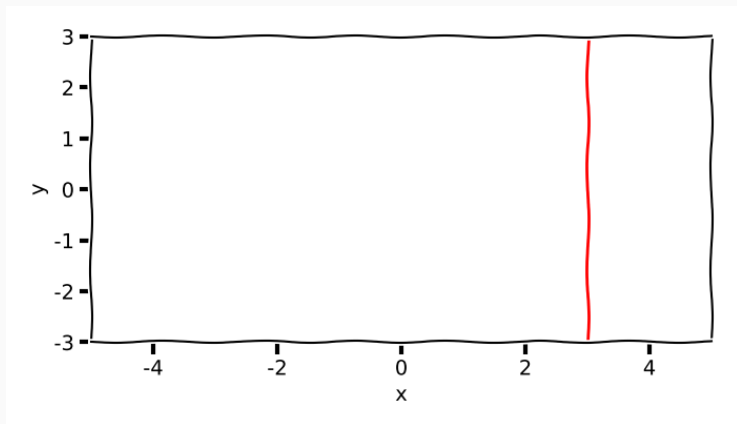
- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping
 - Linear, it is a line
 - Kernels, it is this specific basis function

- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping
 - Linear, it is a line
 - Kernels, it is this specific basis function
- *need a prior assumption over the space of functions*

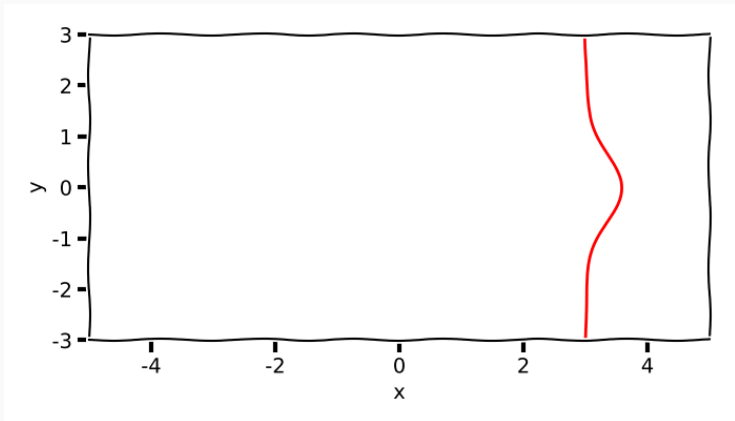
Gaussian Processes



Gaussian Processes

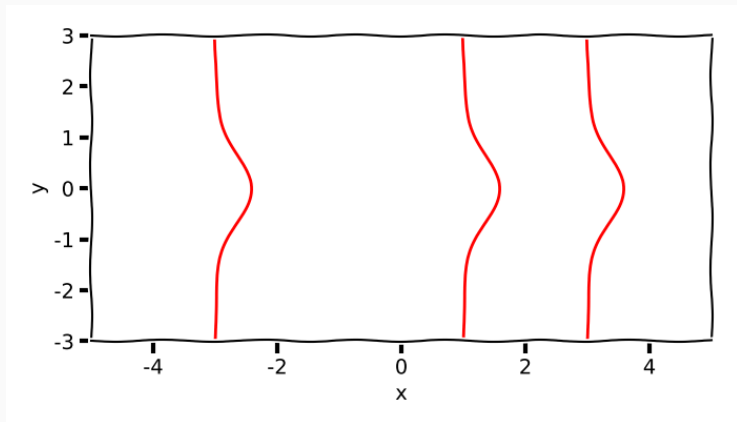


Gaussian Processes



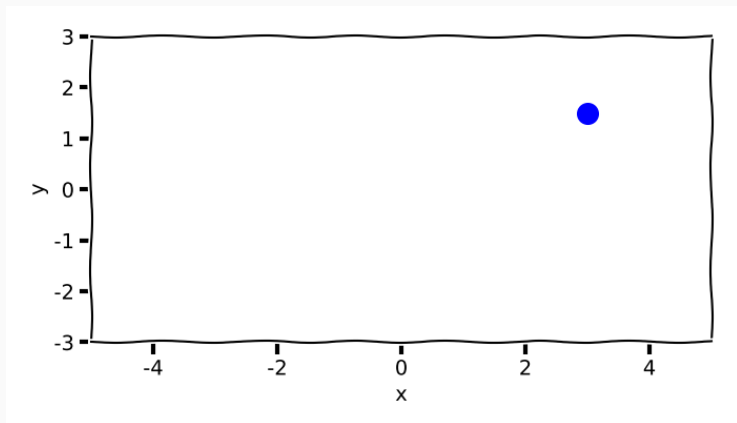
$$p(y|x) = \mathcal{N}(\mu(x), \Sigma(x))$$

Gaussian Processes



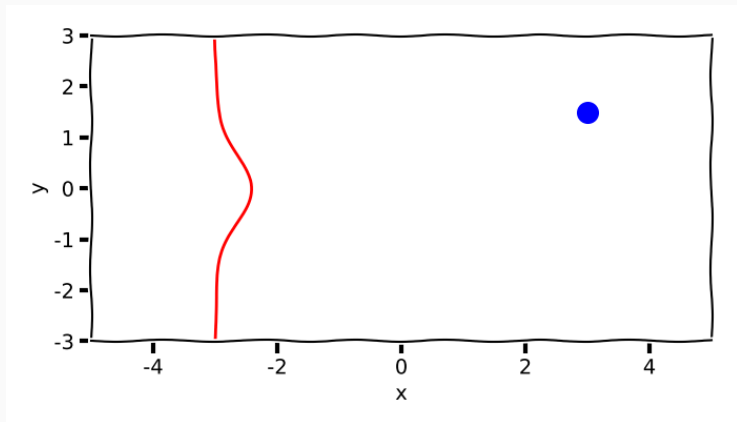
$$p(y_1, y_2, y_3 | x_1, x_2, x_3)$$

Gaussian Processes



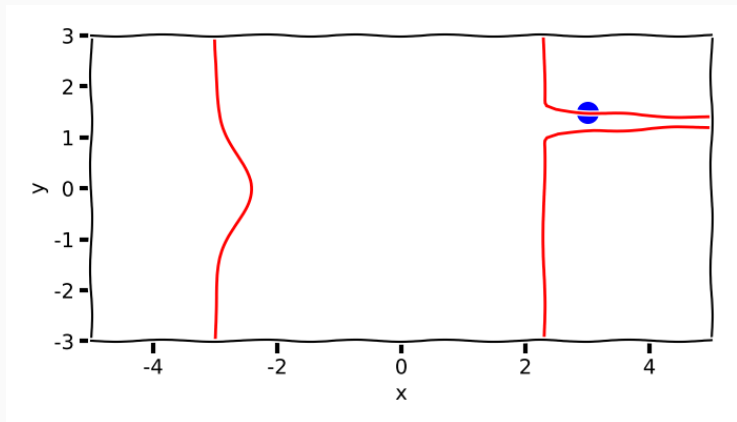
$$p(y_2|x_2, y_1, x_1) = \mathcal{N}(\mu(x_2, x_1, y_1), \Sigma(x_2, x_1, y_1))$$

Gaussian Processes



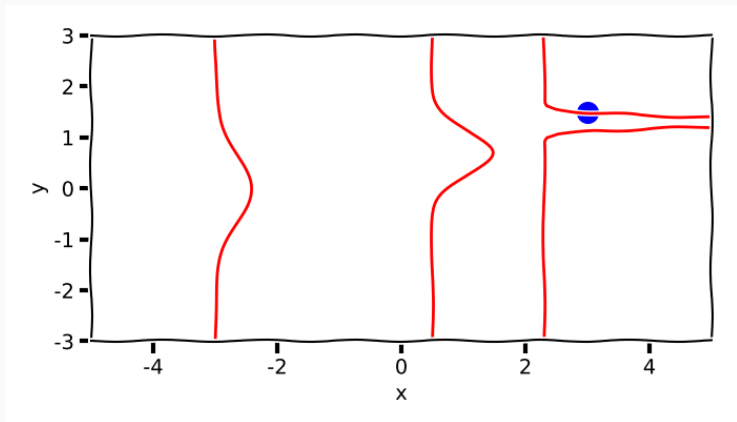
$$p(y_2|x_2, y_1, x_2) = \mathcal{N}(\mu(x_2, x_1, y_1), \Sigma(x_2, x_1, y_1))$$

Gaussian Processes



$$p(y_2|x_2, y_1, x_2) = \mathcal{N}(\mu(x_2, x_1, y_1), \Sigma(x_2, x_1, y_1))$$

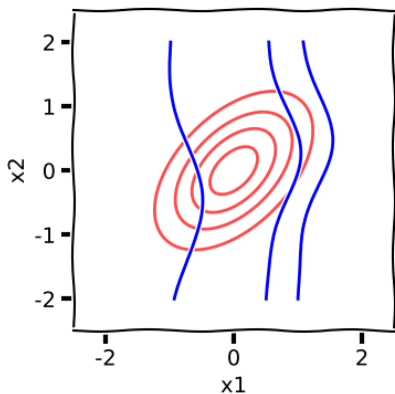
Gaussian Processes



$$p(y_2|x_2, y_1, x_2) = \mathcal{N}(\mu(x_2, x_1, y_1), \Sigma(x_2, x_1, y_1))$$

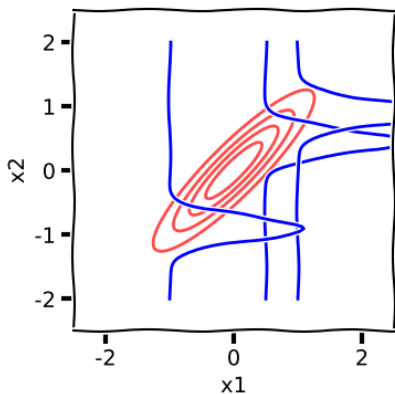
$$p(y_1, y_2, y_3 | x_1, x_2, x_3) = \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \mu(x_3) \end{bmatrix}, \begin{bmatrix} \Sigma(x_1, x_1) & \Sigma(x_1, x_2) & \Sigma(x_1, x_3) \\ \Sigma(x_2, x_1) & \Sigma(x_2, x_2) & \Sigma(x_2, x_3) \\ \Sigma(x_3, x_1) & \Sigma(x_3, x_2) & \Sigma(x_3, x_3) \end{bmatrix} \right)$$

Conditional Gaussians



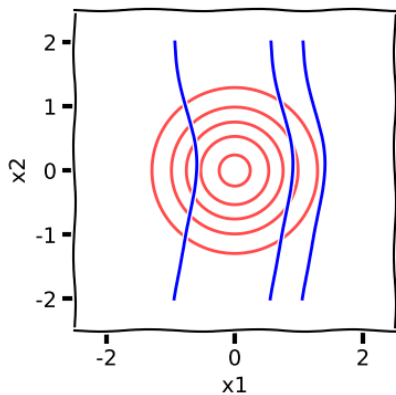
$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$$

Conditional Gaussians



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

Conditional Gaussians

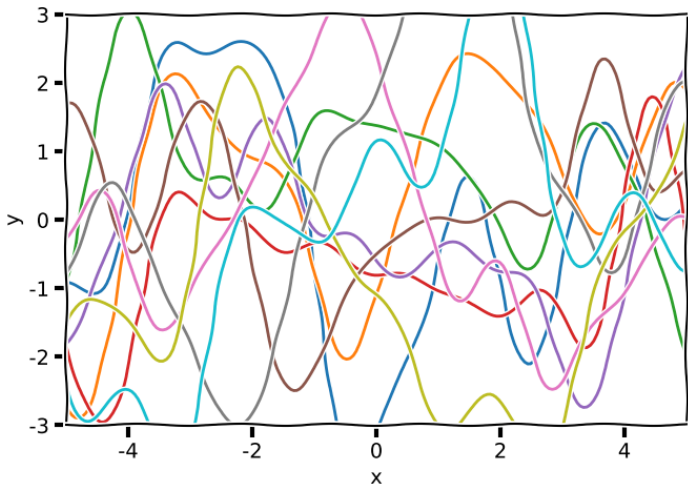


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

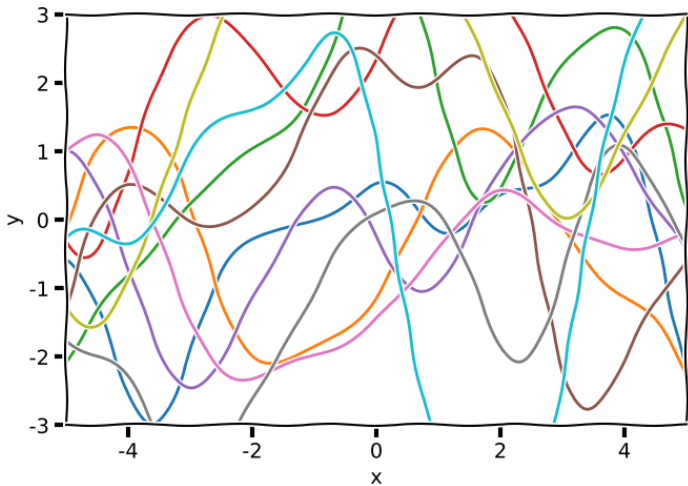
$$p(y_1, \dots, y_N | x_1, \dots, x_N) = \mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \vdots \\ \mu_N \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1N} \\ \vdots & \ddots & \vdots \\ \Sigma_{N1} & \cdots & \Sigma_{NN} \end{bmatrix} \right)$$

- $x \in \mathcal{X}$ our infinite input domain
- $\mu_i = \mu(x_i)$ a function from $\mu(x) : \mathcal{X} \rightarrow \mathbb{R}$
- $\Sigma_{ij} = k(x_i, x_j)$ a function $k(x_i, x_j) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- these two functions completely specifies a Gaussian process

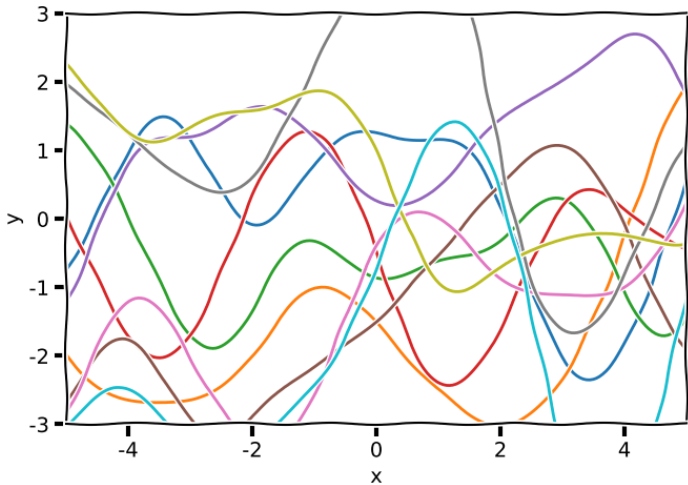
Sampling



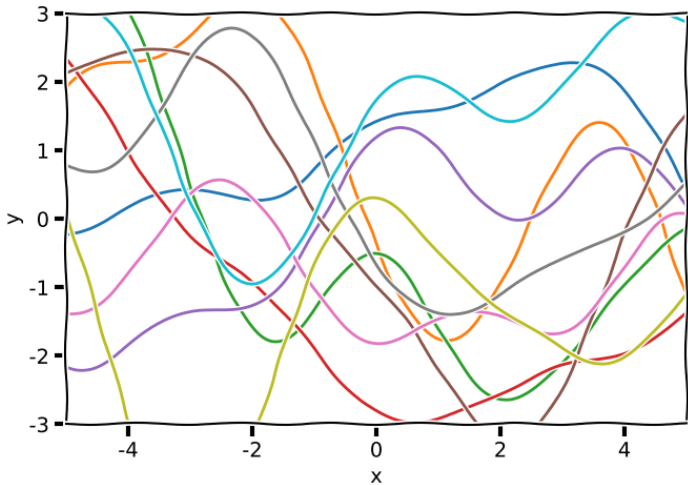
Sampling



Sampling



Sampling



$$p(y_1, \dots, y_N | x_1, \dots, x_N) = \mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \vdots \\ \mu_N \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1N} \\ \vdots & \ddots & \vdots \\ \Sigma_{N1} & \cdots & \Sigma_{NN} \end{bmatrix} \right)$$

- $x \in \mathcal{X}$ our infinite input domain
- $\mu_i = \mu(x_i)$ a function from $\mu(x) : \mathcal{X} \rightarrow \mathbb{R}$
- $\Sigma_{ij} = k(x_i, x_j)$ a function $k(x_i, x_j) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- these two functions completely specifies a Gaussian process

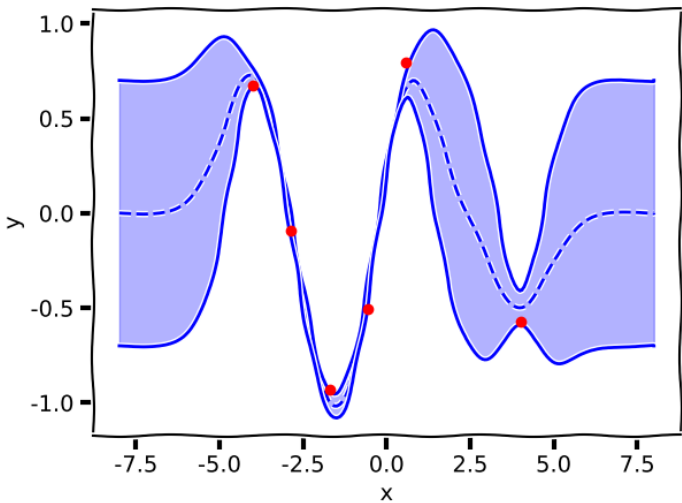
$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

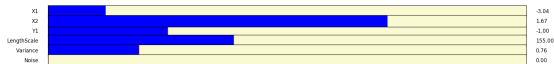
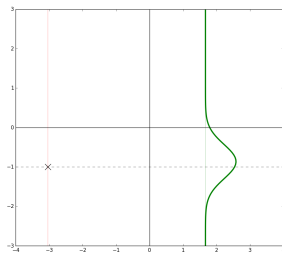
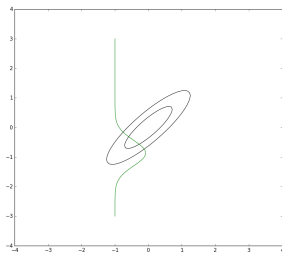
$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}, \theta) = \mathcal{N}(k(\mathbf{x}_*, \mathbf{X})^\top K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f}, \\ k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^\top K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{x}_*))$$

$$p(y_1, y_2 | x_1, x_2) = p(y_1 | y_2, x_1, x_2) p(y_2 | x_2)$$

- As soon as we evaluate the GP at a finite number of locations its a simple Gaussian distribution
- We used the relationship above to derive the posterior over y_1

Posterior





Reset

Learning

$$p(f|\mathbf{x}) = \mathcal{N}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 e^{-\frac{2}{\ell^2} \sin^2\left(\pi \frac{|\mathbf{x}_i - \mathbf{x}_j|}{p}\right)}$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \frac{2}{\pi} \sin^{-1} \left(\frac{2\mathbf{x}_i^T \Sigma \mathbf{x}_j}{\sqrt{(1 + 2\mathbf{x}_i^T \Sigma \mathbf{x}_i)(1 + 2\mathbf{x}_j^T \Sigma \mathbf{x}_j)}} \right)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 e^{\|\mathbf{x}_i - \mathbf{x}_j\|^2 / l^2}$$

- how do we set the parameters of the co-variance function?

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

- We are not interested in \mathbf{f} directly
- Marginalise out \mathbf{f}
- Gaussian likelihood and Gaussian prior \rightarrow Gaussian marginal

- Deterministic world

$$\mathbb{E}[y] = \int yp(y)dy$$

- Deterministic world

$$\mathbb{E}[y] = \int yp(y)dy$$

- Stochastic world

$$\mathbb{E}[p(y)] = \int p(y|x)p(x)dx$$

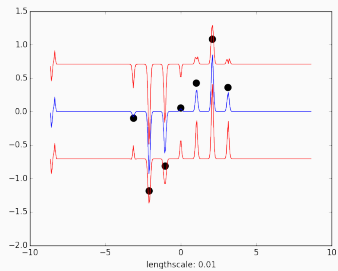
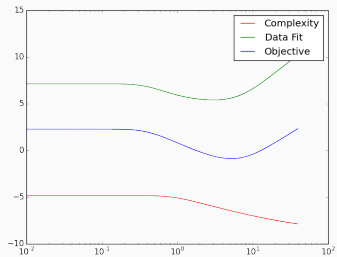
$$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathbf{Y}|\mathbf{X}, \theta)$$

- Type-II Maximum likelihood [1] 3.5.0
- minimise logarithm of marginal likelihood

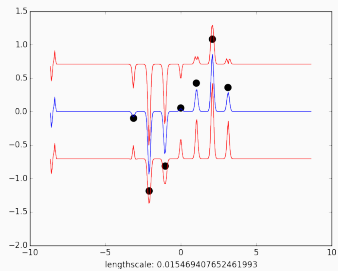
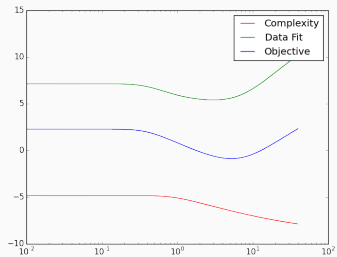
$$\operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} -\log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})) = \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$$

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

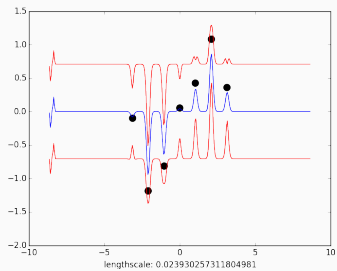
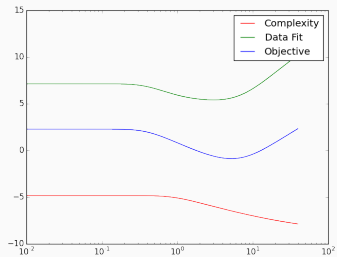
Learning



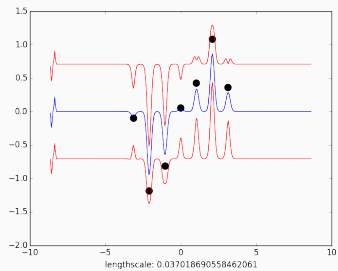
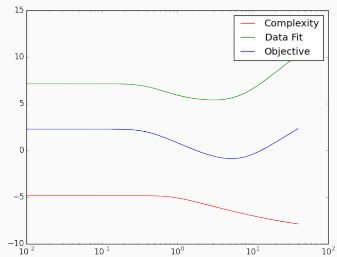
Learning



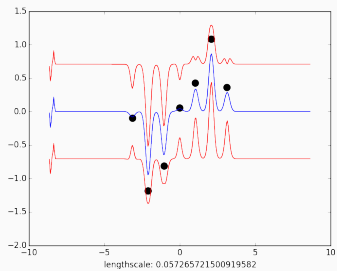
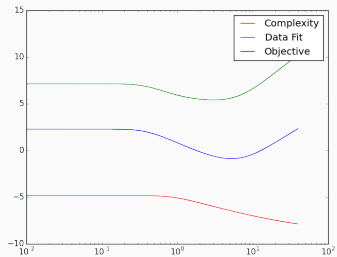
Learning



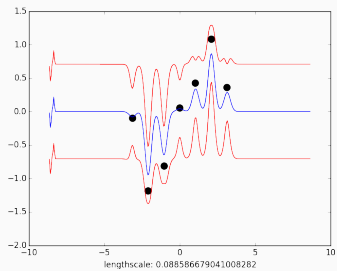
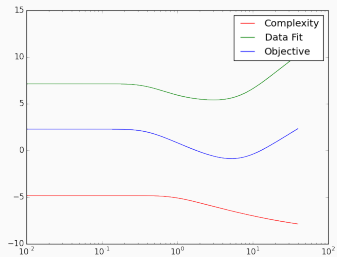
Learning



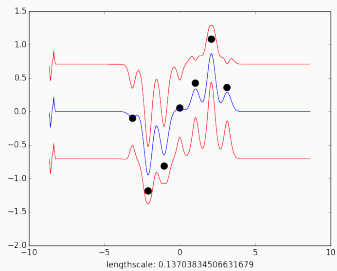
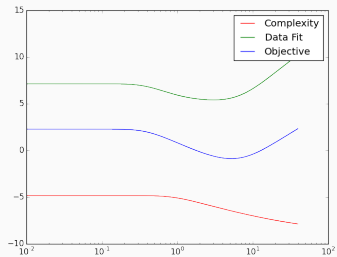
Learning



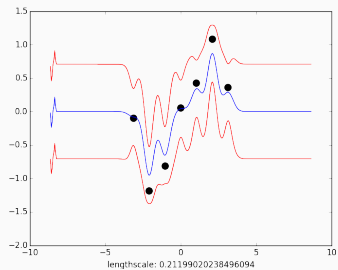
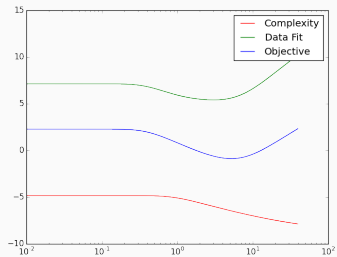
Learning



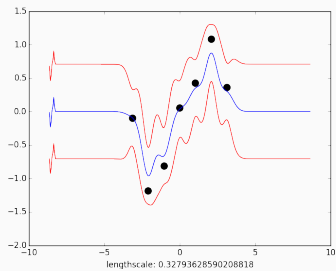
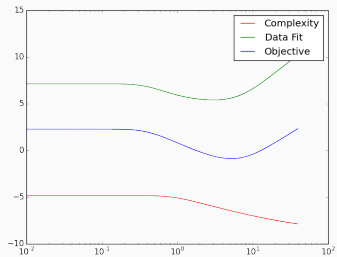
Learning



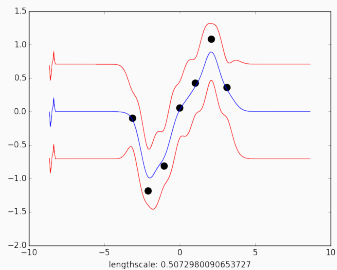
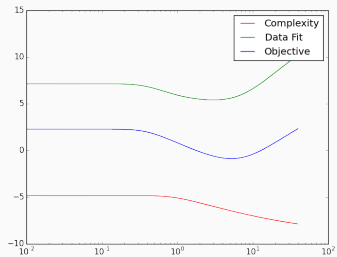
Learning



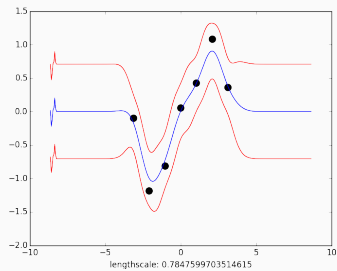
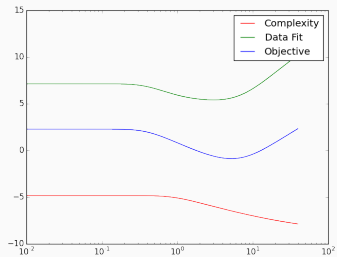
Learning



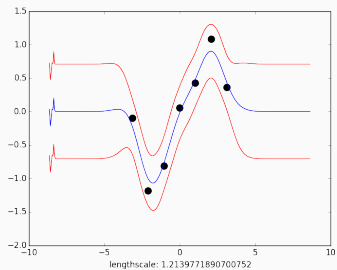
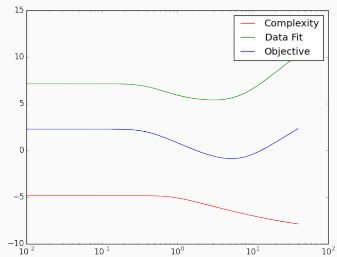
Learning



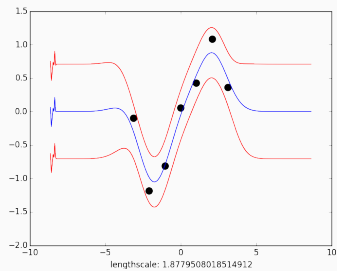
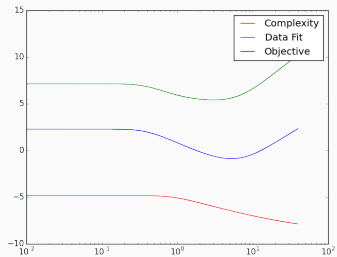
Learning



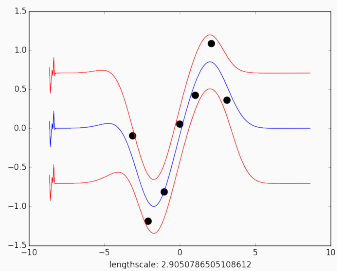
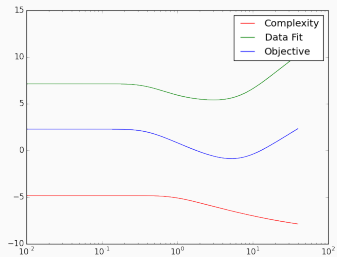
Learning



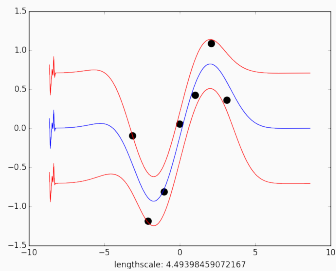
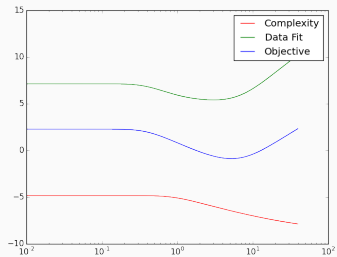
Learning



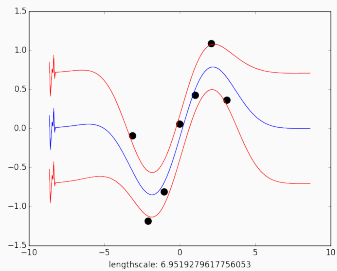
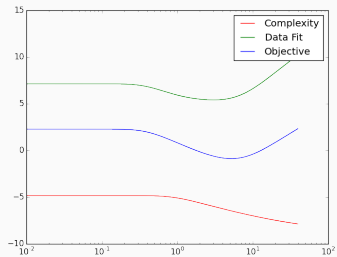
Learning



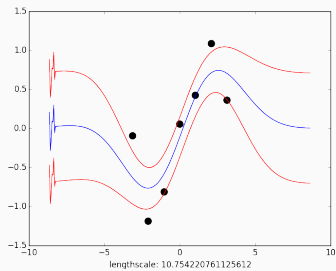
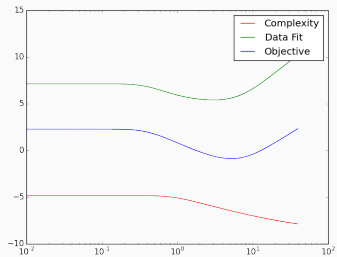
Learning



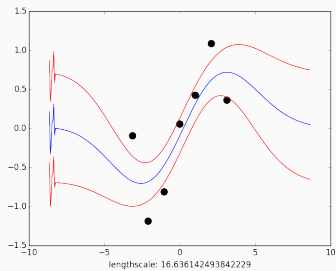
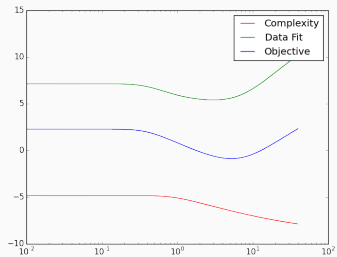
Learning



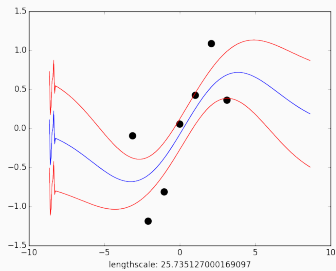
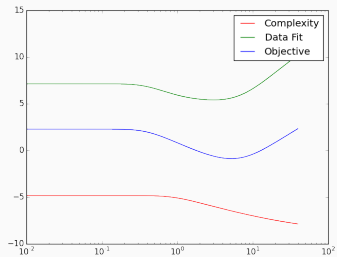
Learning



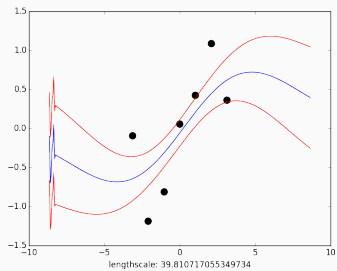
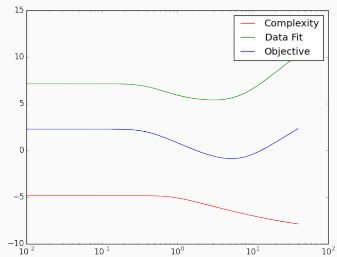
Learning



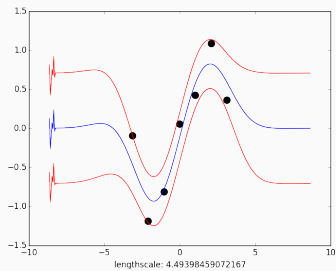
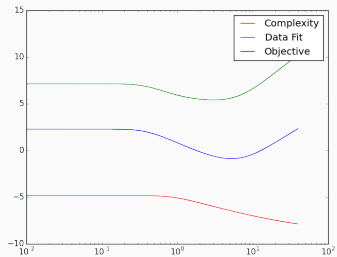
Learning



Learning



Learning



Gaussian Processes

- completely specified by mean and covariance function
- mean and covariance takes **input** variable
- every instantiation of the function is jointly Gaussian
 - conditional and marginal distribution trivial
- very flexible
 - covariance function can encode any behaviour

Unsupervised Learning

$$y = f(x)$$

- given input output pairs we have made assumptions about f
- from data we can update our assumption
- can we push this further?

$$y = f(x)$$

- In unsupervised learning we are given **only** output
- Input is *latent*
- Task: recover both f and x

Latent Variable Models



Latent Variable Models

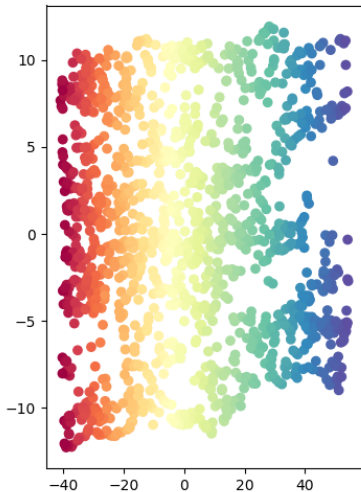
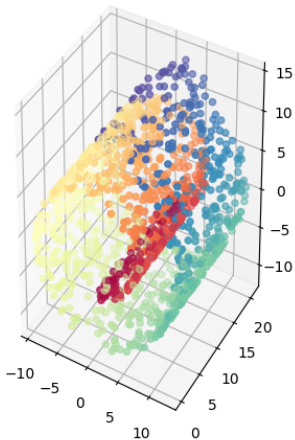


output data $y \in \mathbb{R}^{256 \times 256} \rightarrow 65536$ dimensions

input location on sphere $\rightarrow 3$ dimensions

manifold images lie on a 3 dimensional surface in 65536 dimensions

Manifold



Linear Latent Variable Models [1] 12.2

- Observed data

$$\mathbf{x} \in \mathbb{R}^D$$

- Latent variable

$$\mathbf{z} \in \mathbb{R}^M$$

- Mapping

$$\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \epsilon$$

- Likelihood: make noise assumption $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

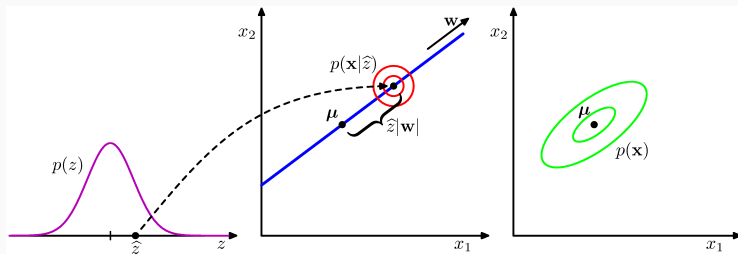
$$p(\mathbf{x}|\mathbf{z}, \mathbf{W}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

- Prior ?

$$\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \epsilon$$

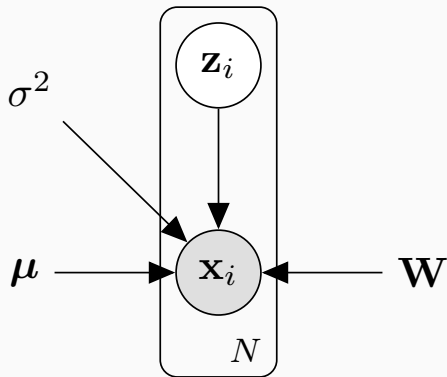
- marginalise out both \mathbf{W} and \mathbf{z} is intractable
- marginalise out one and infer the other
- $\mathbf{W} \in \mathbb{R}^{D \times M}$ and $\mathbf{z} \in \mathbb{R}^{M \times N}$
- N commonly larger than $D \Rightarrow$ integrate out \mathbf{z}

Principal Component Analysis [1] Figure 12.9



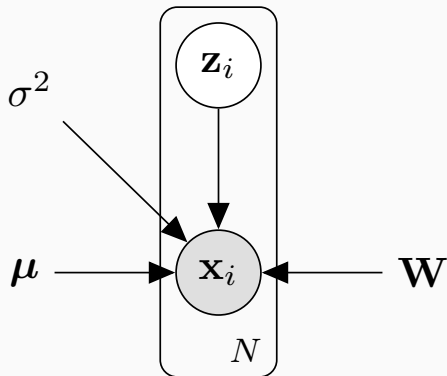
$$p(z) = \mathcal{N}(z|0, I)$$

Graphical Model



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

Graphical Model



$$p(x, z | W, \mu, \sigma^2) = p(x | z, W, \mu, \sigma^2) p(z)$$

$$p(\mathbf{x}|\mathbf{W}) = \int p(\mathbf{x}|\mathbf{z}, \mathbf{W})p(\mathbf{z})d\mathbf{z} = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$$

- Gaussian distribution closed under linear transformation (interesting proof)

$$\begin{aligned} p(\mathbf{x}|\mathbf{W}) &= \int p(\mathbf{x}|\mathbf{z}, \mathbf{W})p(\mathbf{z})d\mathbf{z} = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \\ &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}) \end{aligned}$$

- Gaussian distribution closed under linear transformation (interesting proof)

$$\log p(\mathbf{X} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)$$

- find stationary with respect to each variable gives Maximum likelihood solution to \mathbf{W} , $\boldsymbol{\mu}$ and σ^2
- *In the assignment we make it easier and take derivatives instead and optimise*

$$\log p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \log p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)$$

- find stationary with respect to each variable gives Maximum likelihood solution to \mathbf{W} , $\boldsymbol{\mu}$ and σ^2
- *In the assignment we make it easier and take derivatives instead and optimise*

$$\begin{aligned}\log p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \log p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \log(2\pi) - \frac{N}{2} \log|\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}| \\ &\quad - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\end{aligned}$$

- find stationary with respect to each variable gives Maximum likelihood solution to \mathbf{W} , $\boldsymbol{\mu}$ and σ^2
- *In the assignment we make it easier and take derivatives instead and optimise*

$$p(\mathbf{z}|\mathbf{x}) \propto p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^2(\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})^{-1})$$

- Gaussian likelihood and Gaussian prior \rightarrow Gaussian posterior

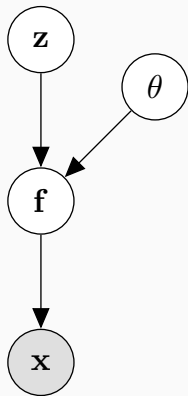
Principal Component Analysis

- You might have seen this explained in a different way
 - *Retain variance*
 - *Error minimisation*
- These provides the same solution as the maximum likelihood but solved by an eigenvalue problem
- Do not provide intuition as it doesn't state assumptions

Question 15-21

You now have all the material to finish the assignment!

Non-linear Latent variable model



$$p(\mathbf{x}|\mathbf{z}, \theta) = \int p(\mathbf{x}|\mathbf{f})p(\mathbf{f}|\mathbf{z}, \theta)d\mathbf{f}$$

Font Demo

Summary

- Type II Maximum likelihood
- As long as I make assumptions I can learn from data
- Unsupervised learning, just the same, just a prior instead of observations
- Next 3 lectures

eof

References



Christopher M. Bishop.

Pattern Recognition and Machine Learning (Information Science and Statistics).

Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.