

Machine Learning

Unsupervised Learning

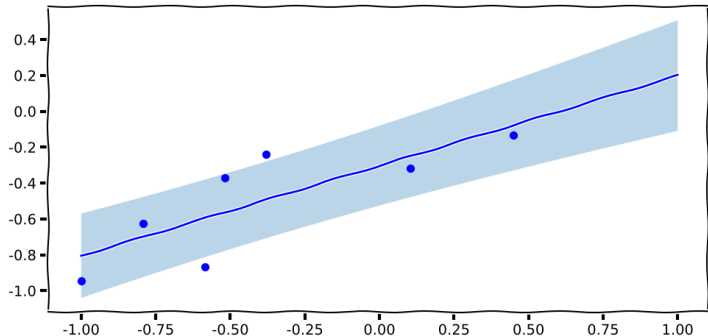
Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

October 22, 2019

<http://www.carlhenrik.com>

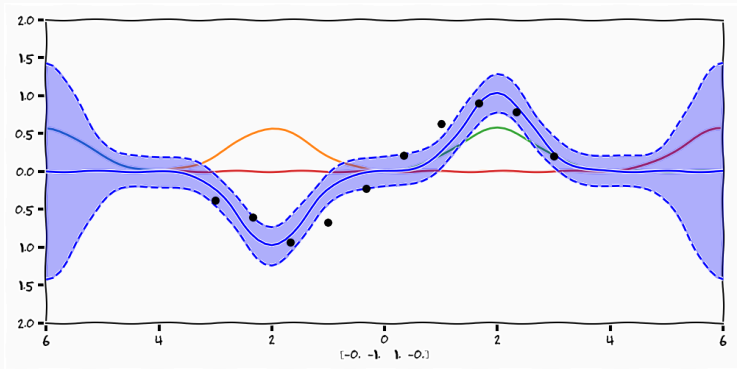
Introduction

Regression: Linear



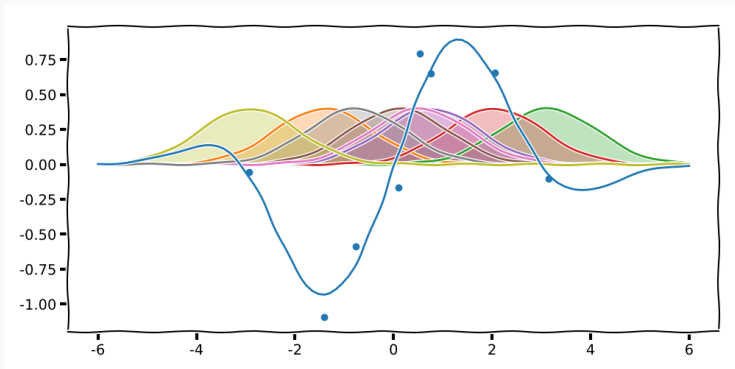
$$y_i = \mathbf{w}^T \mathbf{x}_i$$

Regression: Linear Basis



$$y_i = \sum_{k=1}^K w_k \phi_k(x_i)$$

Regression: Kernel

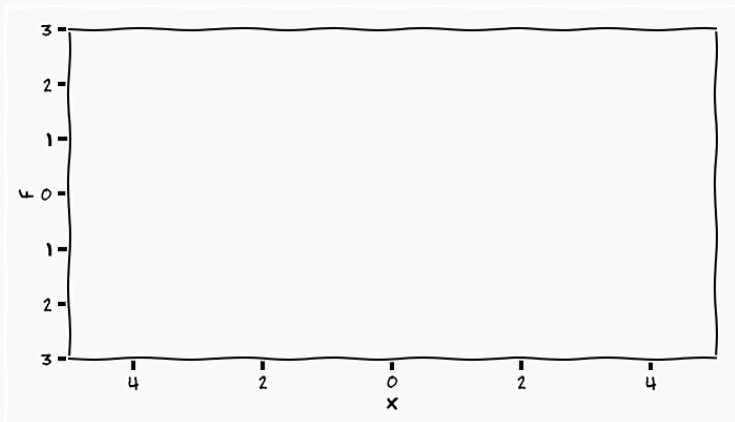


$$y_i = k(\mathbf{x}_i, \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \lambda \mathbf{I})^{-1} \mathbf{y}$$

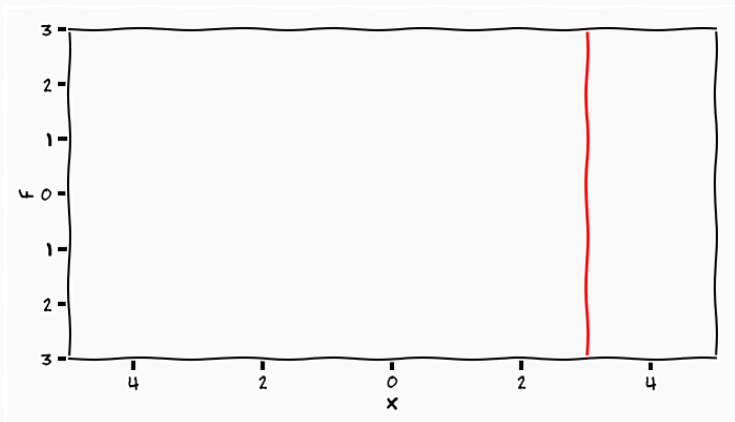
Regression

- Linear
 - we are limited by lines
- Basis functions
 - + nonlinear functions
 - how many basis functions should I have, what should they look like?
 - prior hard to interpret
- Kernel
 - + complexity set by data
 - no uncertainty in our estimate

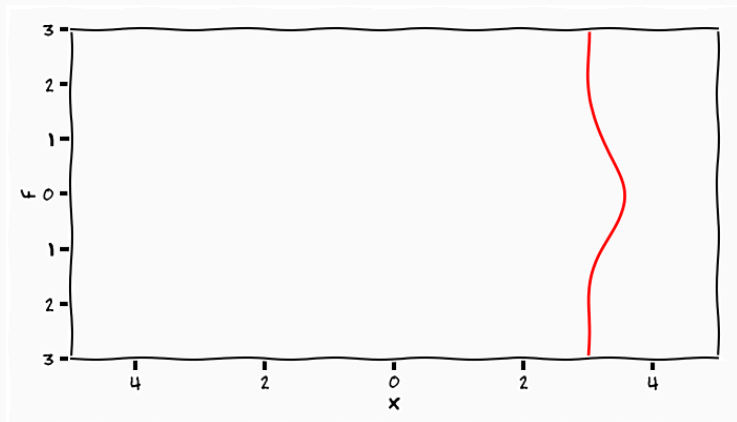
Gaussian Processes



Gaussian Processes

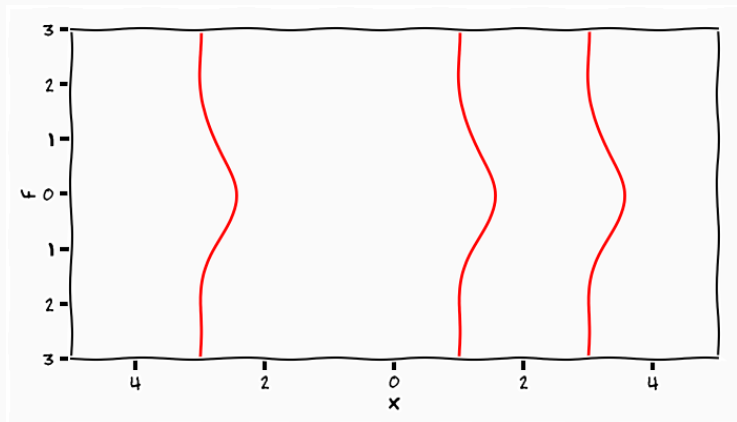


Gaussian Processes



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

Gaussian Processes



$$p(f_1, f_2, f_3 | x_1, x_2, x_3)$$

Gaussian Process: definition

$$p(f_1, f_2, \dots, f_N, \dots | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}))$$

$$= \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_N) \\ \vdots \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_N) & \cdots \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_N) & \cdots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \cdots & k(x_N, x_N) & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \right)$$

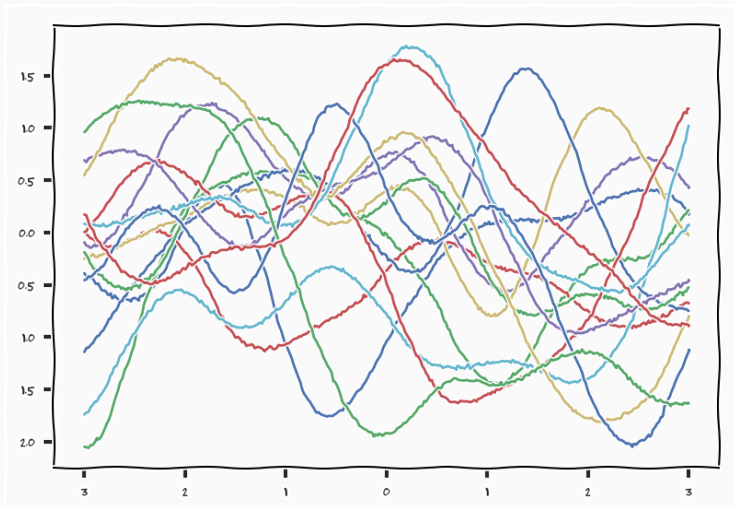
Marginal

$$p(f_1, f_2 | x_1, x_2) = \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{bmatrix} \right)$$

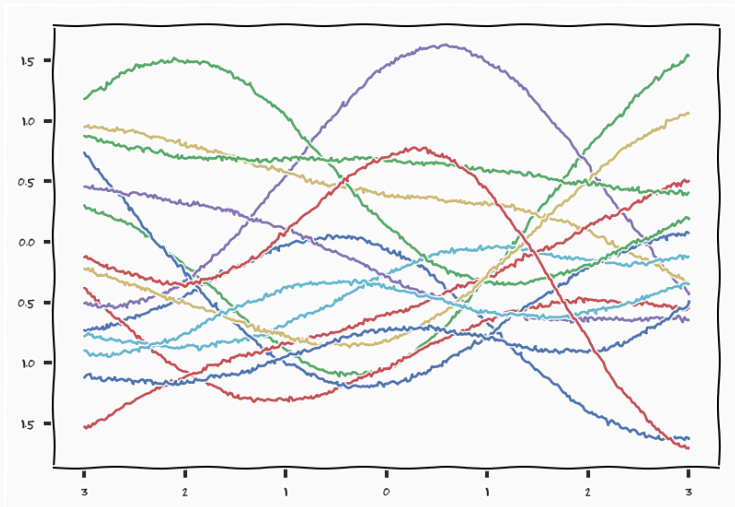
Conditional

$$p(f_1 | f_2, x_1, x_2) = \mathcal{N}(\mu(x_1) + k(x_1, x_2)k(x_2, x_2)^{-1}(f_2 - \mu(x_2)), \\ k(x_1, x_1) - k(x_1, x_2)k(x_2, x_2)^{-1}k(x_2, x_1))$$

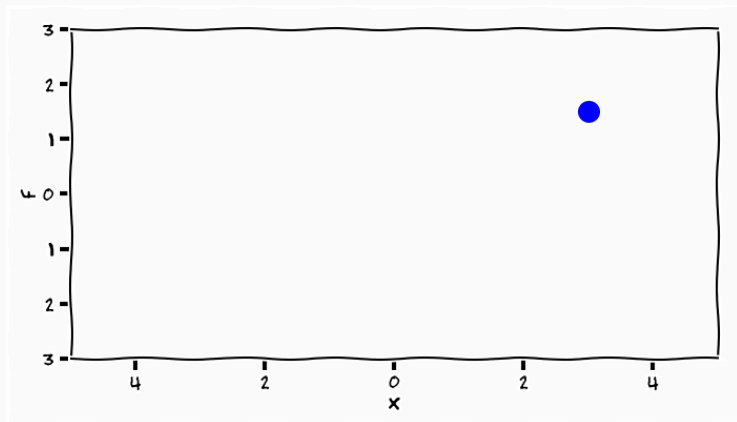
Sampling



Sampling

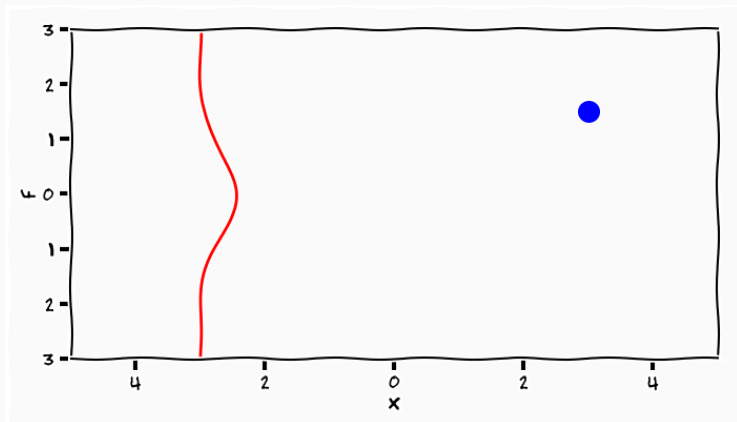


Gaussian Processes



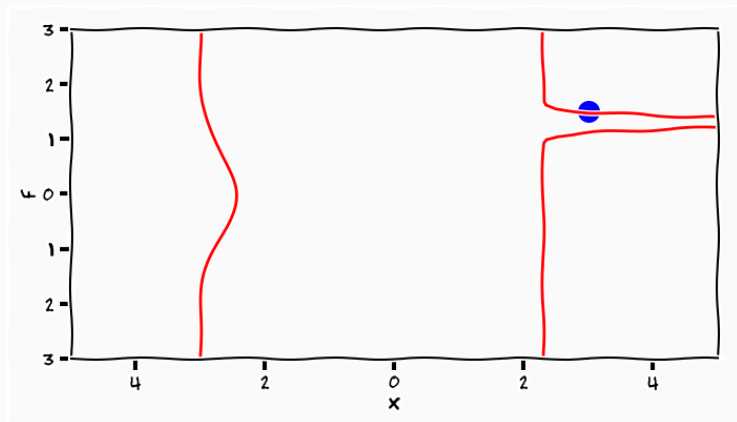
$$p(f_2 | x_2, f_1, x_1) = \mathcal{N}(\mu(x_2, x_1, f_1), \Sigma(x_2, x_1, f_1))$$

Gaussian Processes



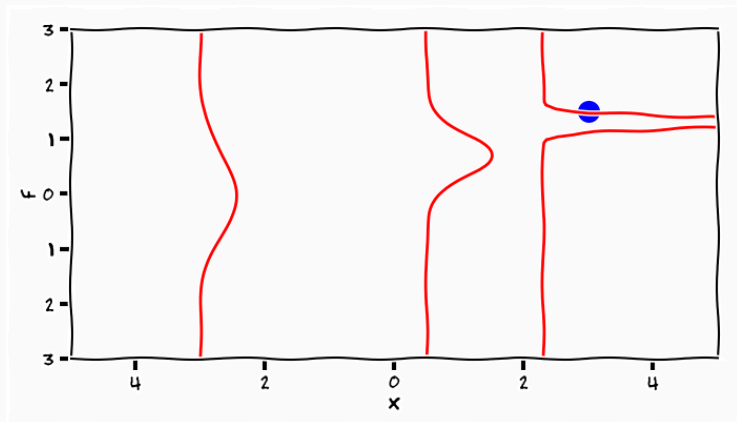
$$p(f_2|x_2, y_1, x_2) = \mathcal{N}(\mu(x_2, x_1, f_1), \Sigma(x_2, x_1, f_1))$$

Gaussian Processes



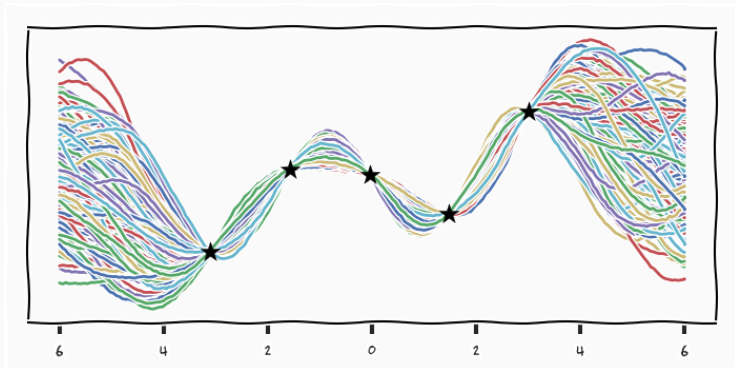
$$p(f_2|x_2, f_1, x_2) = \mathcal{N}(\mu(x_2, x_1, f_1), \Sigma(x_2, x_1, f_1))$$

Gaussian Processes

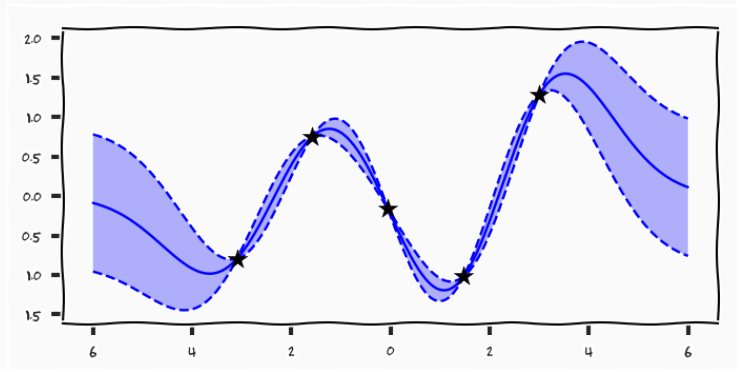


$$p(f_2|x_2, f_1, x_2) = \mathcal{N}(\mu(x_2, x_1, f_1), \Sigma(x_2, x_1, f_1))$$

Gaussian Processes Posterior



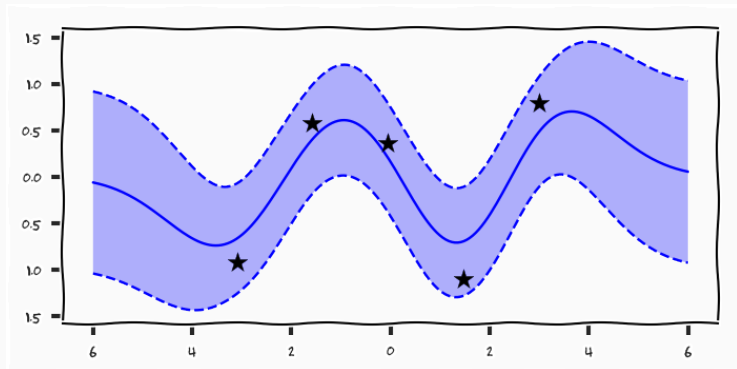
Gaussian Processes Posterior



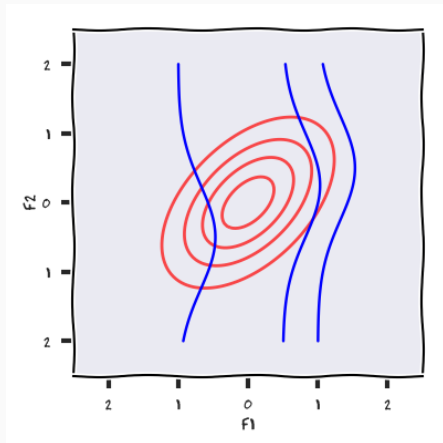
$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}) = \mathcal{N}(k(\mathbf{x}_*, \mathbf{X})^\top (K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{f}, \\ k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^\top (K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{X}, \mathbf{x}_*))$$

Gaussian Processes



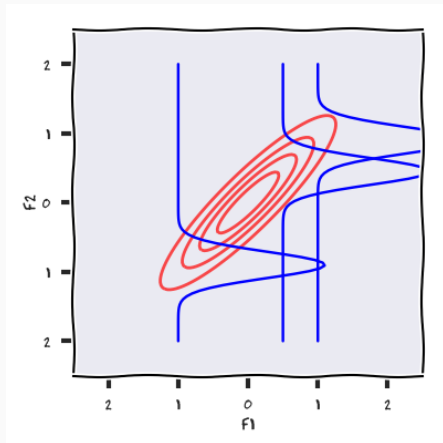
Conditional Gaussians



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

$$\left[\begin{array}{c} \mu(x_1) \\ \mu(x_2) \end{array} \right] \left[\begin{array}{cc} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{array} \right]$$

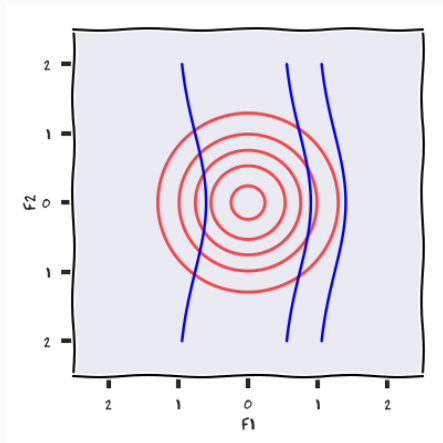
Conditional Gaussians



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

$$\left[\begin{array}{c} \mu(x_1) \\ \mu(x_2) \end{array} \right] \left[\begin{array}{cc} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{array} \right]$$

Conditional Gaussians

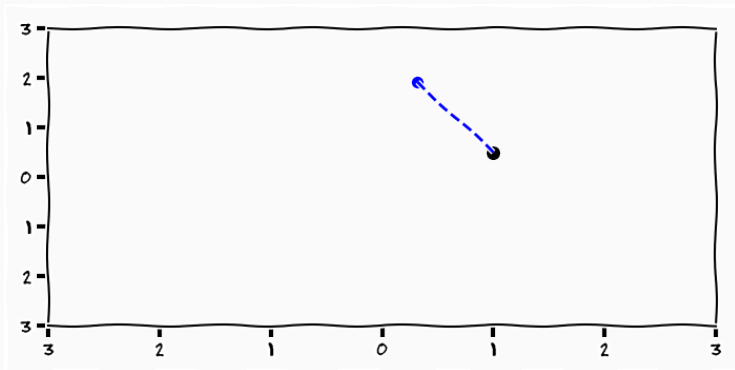


$$\mathcal{N} \left(\underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\mu}, \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}_{\Sigma} \right)$$

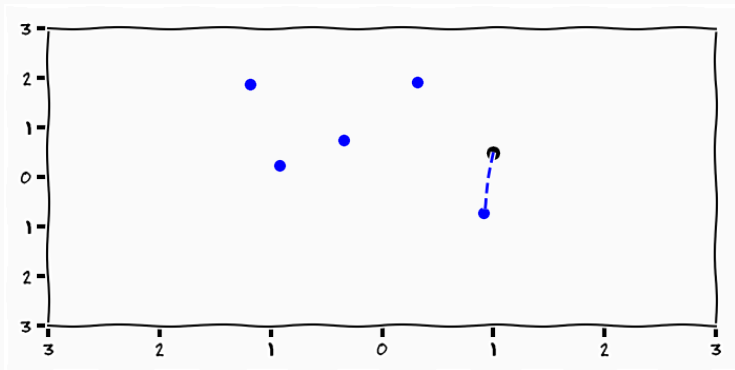
$$\left[\begin{array}{c} \mu(x_1) \\ \mu(x_2) \end{array} \right] \left[\begin{array}{cc} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{array} \right]$$

Non-Parametrics??

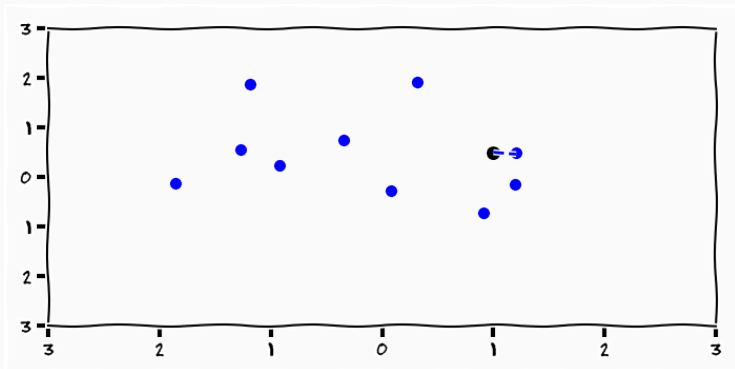
Nearest Neighbour



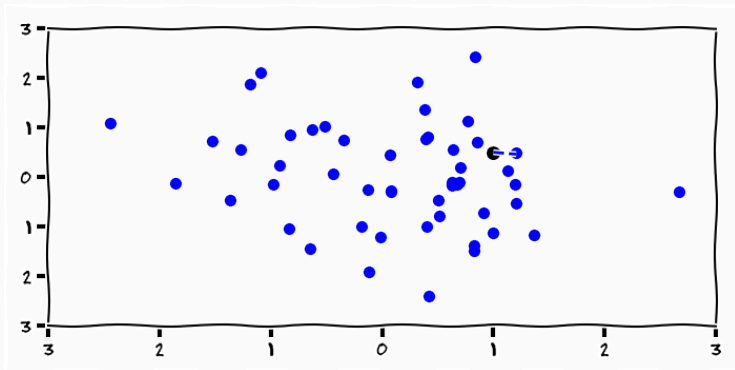
Nearest Neighbour



Nearest Neighbour



Nearest Neighbour



Learning

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

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

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

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

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

- 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

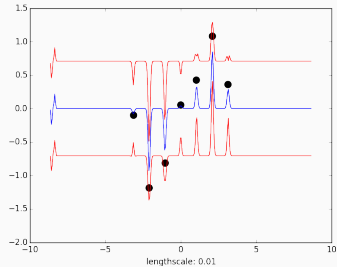
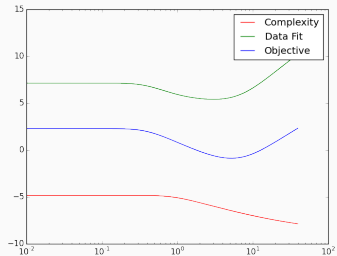
$$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathbf{Y}|\mathbf{X}, \theta)$$

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

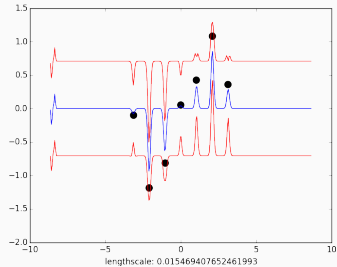
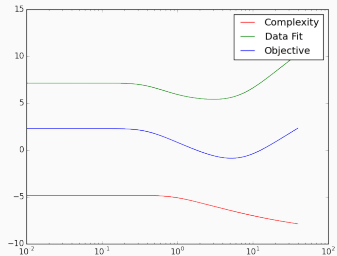
$$\operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} -\log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})) = \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$$

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

Learning

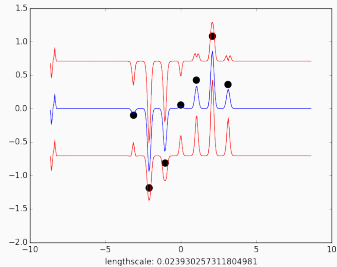
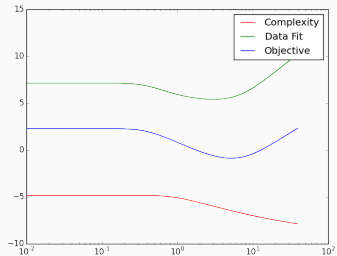


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

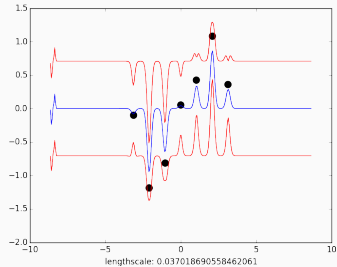
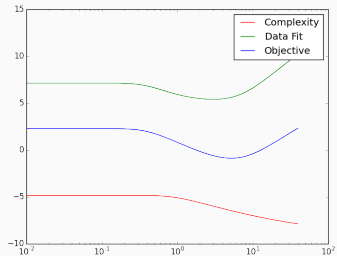


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

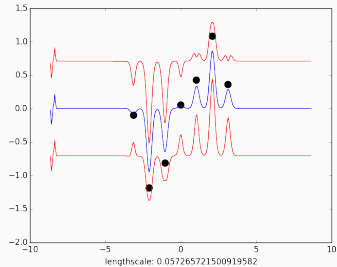
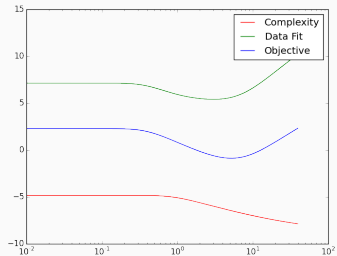
Learning



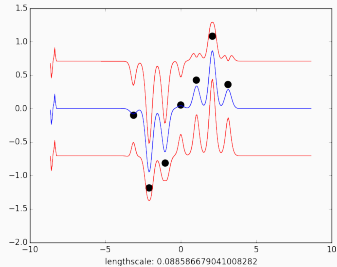
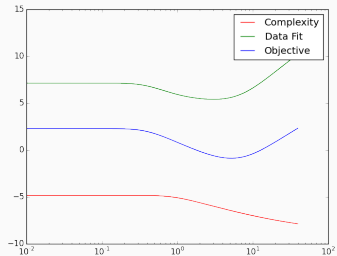
$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$



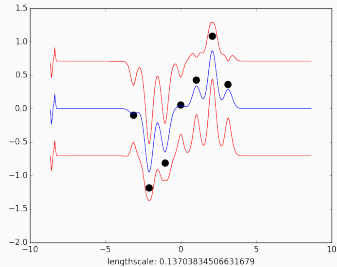
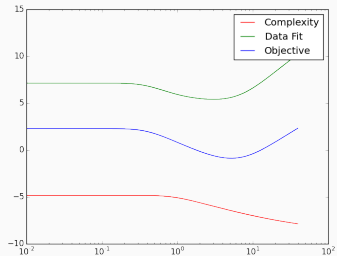
$$\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \quad \frac{1}{2} \log |\mathbf{K}|$$



$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

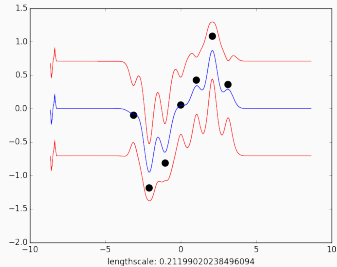
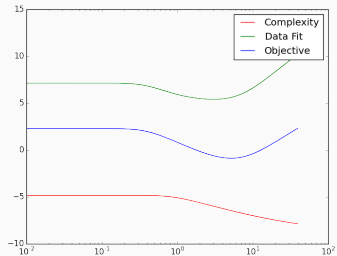


$$\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \quad \frac{1}{2} \log |\mathbf{K}|$$

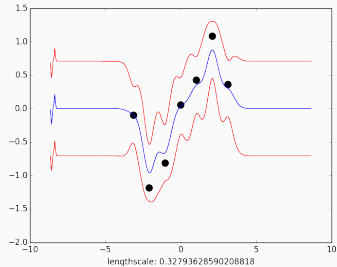
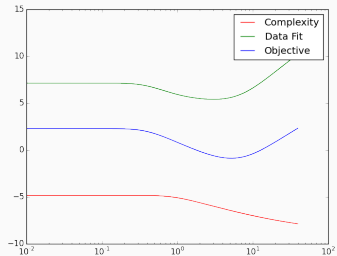


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning

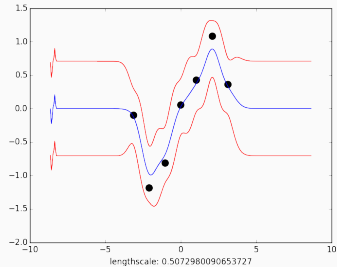
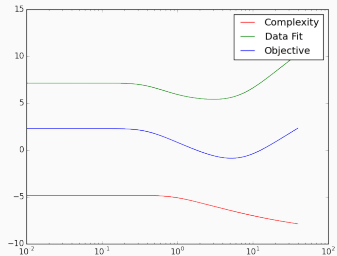


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

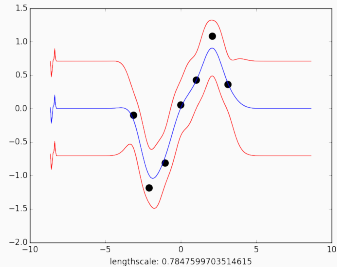
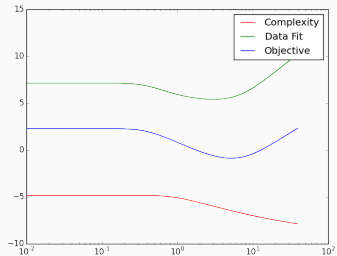


$$\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \quad \frac{1}{2} \log |\mathbf{K}|$$

Learning

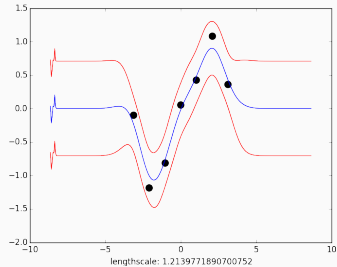
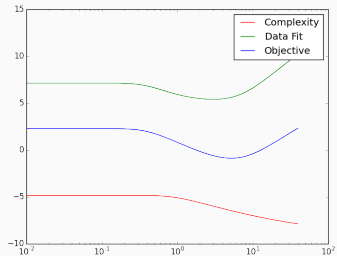


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

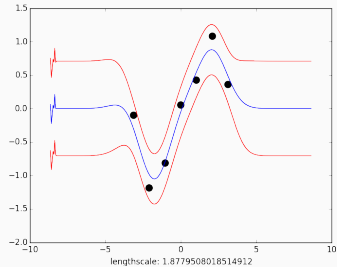
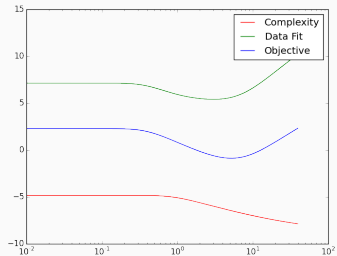


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning

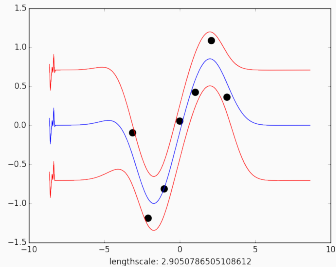
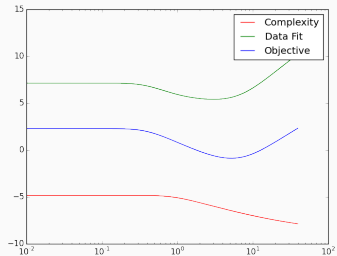


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$



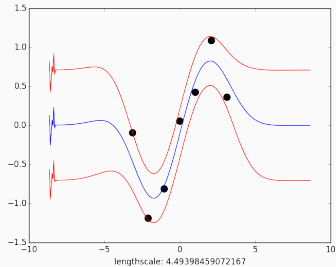
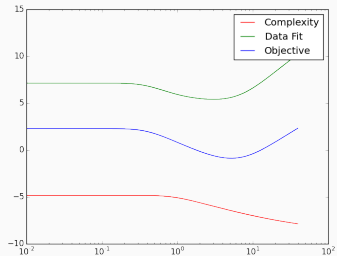
$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning



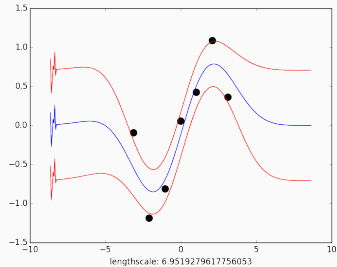
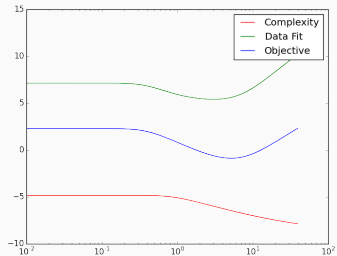
$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning

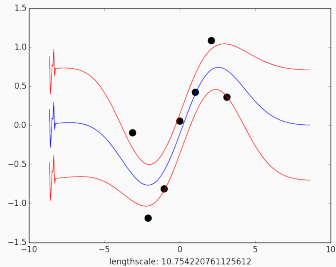
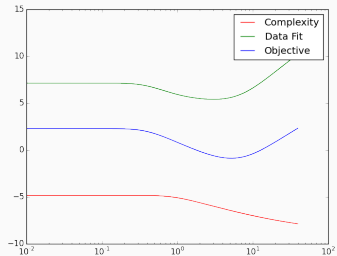


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning

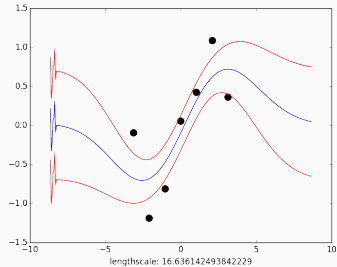
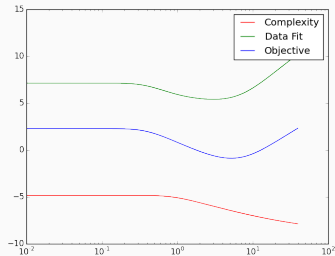


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$



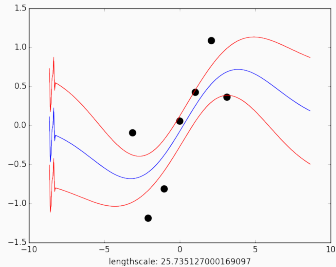
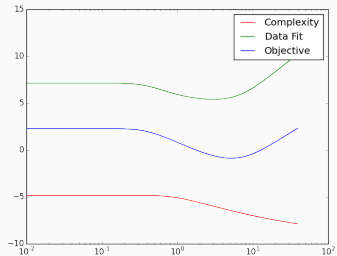
$$\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \quad \frac{1}{2} \log |\mathbf{K}|$$

Learning



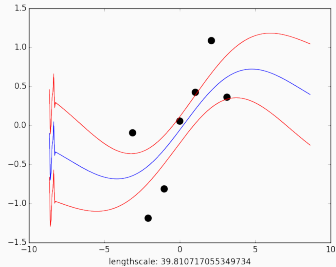
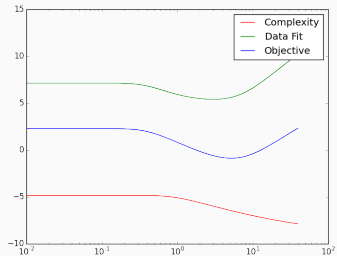
$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning



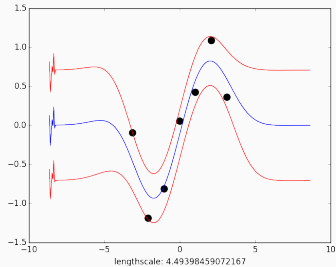
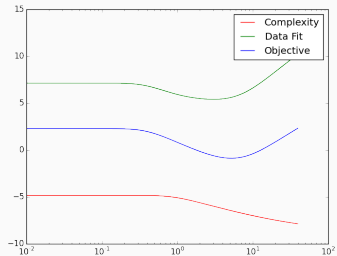
$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning



$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

Learning

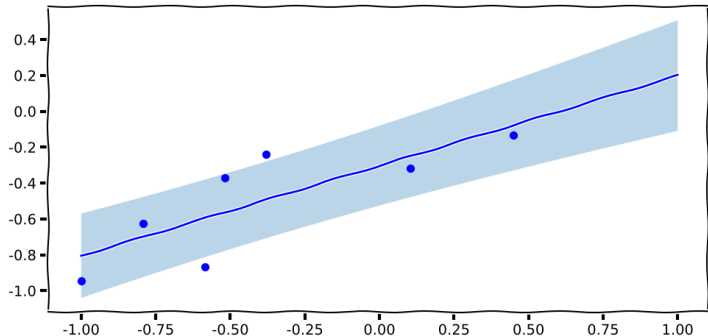


$$\frac{1}{2}\mathbf{y}^T\mathbf{K}^{-1}\mathbf{y} \quad \frac{1}{2}\log|\mathbf{K}|$$

- completely specified by mean and covariance function
- mean and covariance are functions of **input** variable
- every instantiation of the function is jointly Gaussian
 - conditional and marginal distribution trivial
- very flexible
 - covariance function can encode any behaviour
- infer parameters through Type-II maximum likelihood

Unsupervised Learning

Regression: Linear



$$y_i = \mathbf{w}^T \mathbf{x}_i$$

Supervised Learning

$$y_i = f(x_i)$$

- learn relationship $f(\cdot)$ between pairs of data x_i and y_i

Supervised Learning

$$y_i = f(x_i)$$

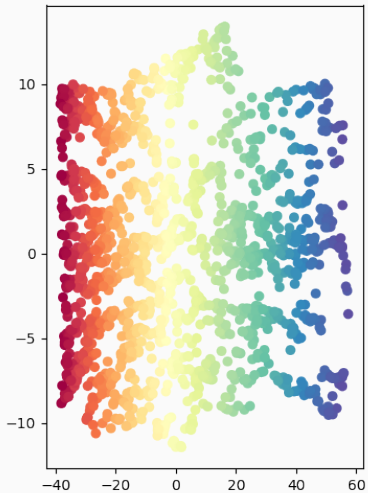
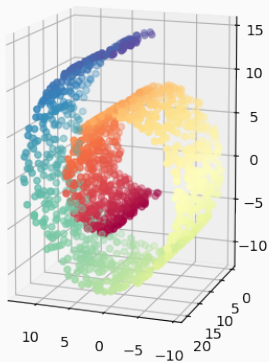
- learn relationship $f(\cdot)$ between pairs of data x_i and y_i

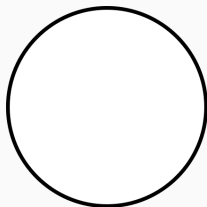
Unsupervised Learning

$$y_i = f(x_i)$$

- learn a representation \mathbf{X} from data \mathbf{Y}

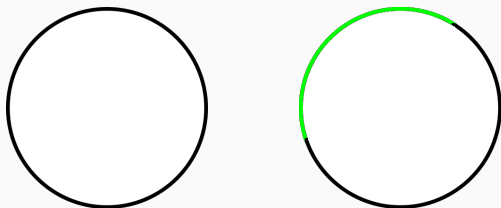
Manifold





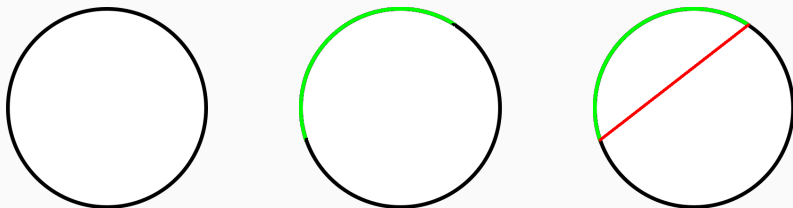
"In mathematics, a manifold is a topological space that near each point resembles Euclidean space"¹

¹<http://en.wikipedia.org/wiki/Manifold>



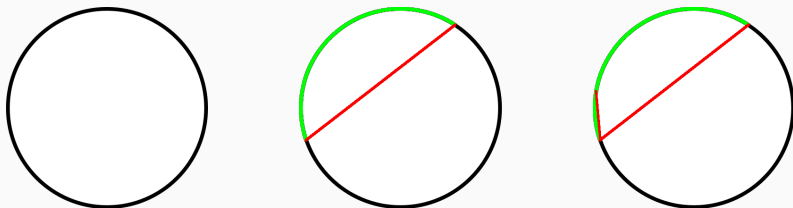
"In mathematics, a manifold is a topological space that near each point resembles Euclidean space"¹

¹<http://en.wikipedia.org/wiki/Manifold>



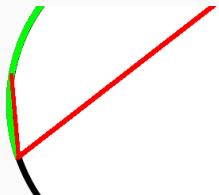
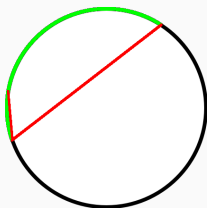
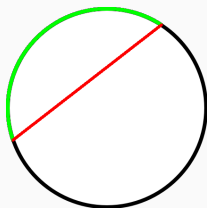
"In mathematics, a manifold is a topological space that near each point resembles Euclidean space"¹

¹<http://en.wikipedia.org/wiki/Manifold>



"In mathematics, a manifold is a topological space that near each point resembles Euclidean space"¹

¹<http://en.wikipedia.org/wiki/Manifold>



"In mathematics, a manifold is a topological space that near each point resembles Euclidean space"¹

¹<http://en.wikipedia.org/wiki/Manifold>

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

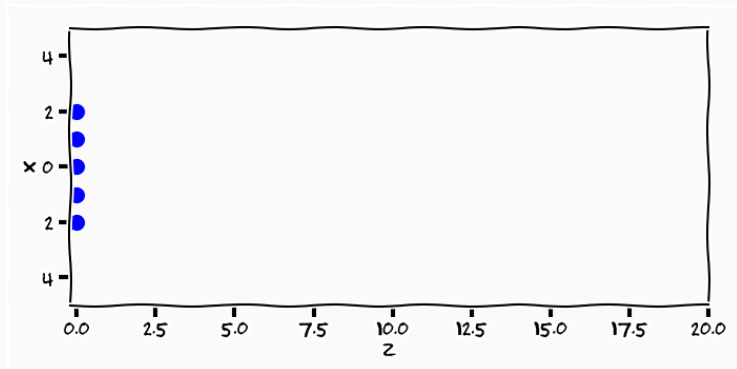
$$y = f(x)$$

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

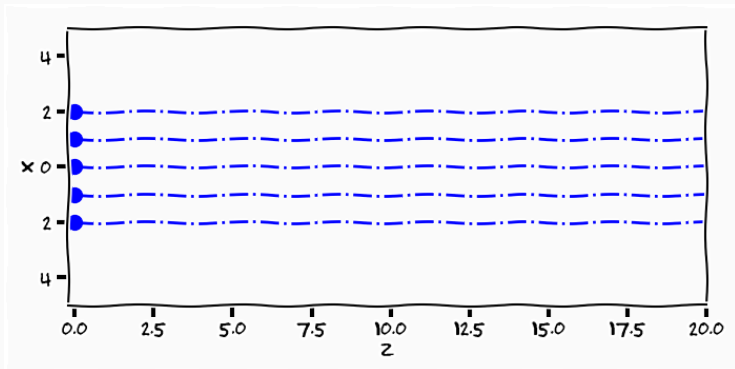
$$y = f(x)$$

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

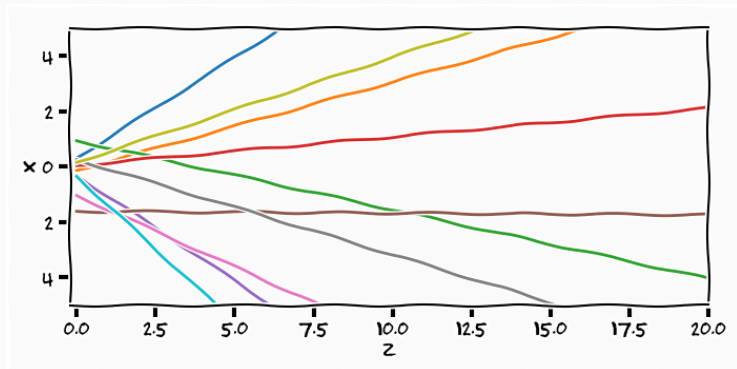
Unsupervised Learning



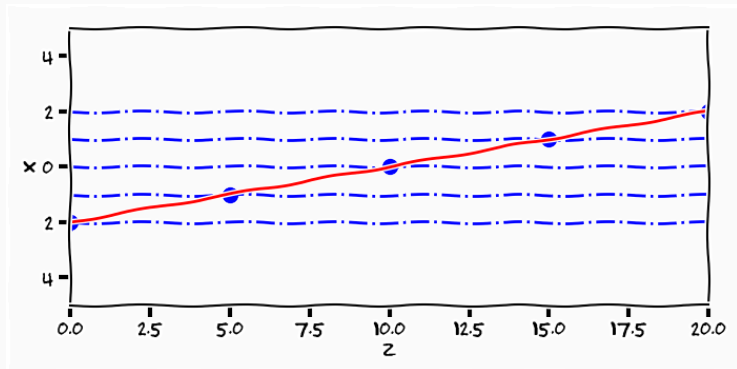
Unsupervised Learning



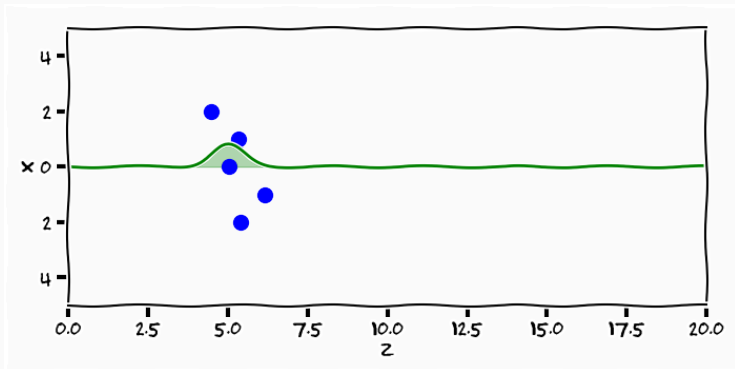
Unsupervised Learning



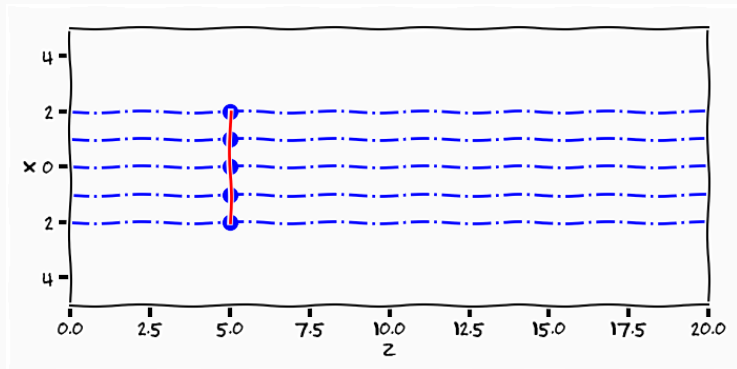
Unsupervised Learning



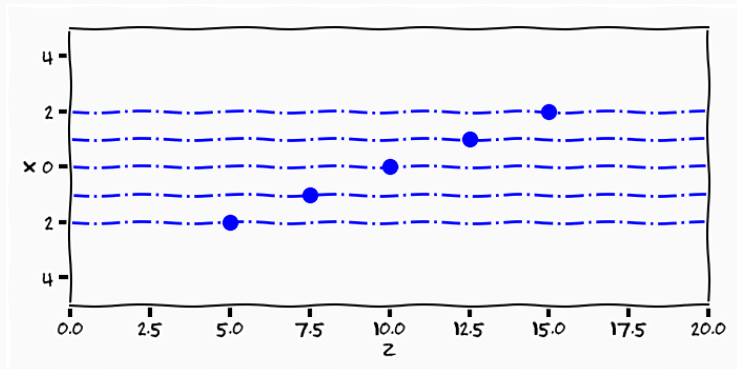
Unsupervised Learning



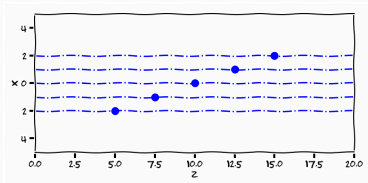
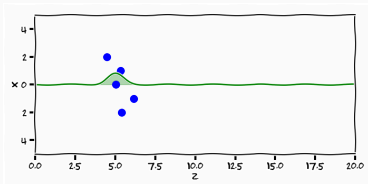
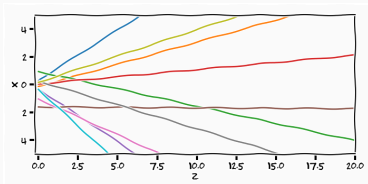
Unsupervised Learning



Unsupervised Learning



Unsupervised Learning



- Linear Regression

$$p(\mathbf{W}|\mathbf{t}, \mathbf{X}) \propto p(\mathbf{t}|\mathbf{W}, \mathbf{X})p(\mathbf{W})$$

- Linear Regression

$$p(\mathbf{W}|\mathbf{t}, \mathbf{X}) \propto p(\mathbf{t}|\mathbf{W}, \mathbf{X})p(\mathbf{W})$$

- Linear Unsupervised Learning

$$p(\mathbf{W}, \mathbf{X}|\mathbf{t}) \propto p(\mathbf{t}|\mathbf{W}, \mathbf{X})p(\mathbf{W})p(\mathbf{X})$$

- Linear Regression

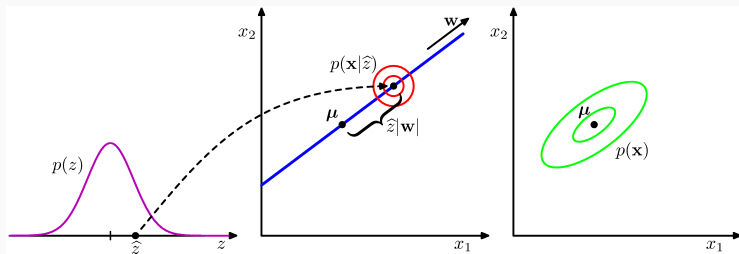
$$p(\mathbf{W}|\mathbf{t}, \mathbf{X}) \propto p(\mathbf{t}|\mathbf{W}, \mathbf{X})p(\mathbf{W})$$

- Linear Unsupervised Learning

$$p(\mathbf{W}, \mathbf{X}|\mathbf{t}) \propto p(\mathbf{t}|\mathbf{W}, \mathbf{X})p(\mathbf{W})p(\mathbf{X})$$

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

Principal Component Analysis [1] Figure 12.9



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

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

$$p(\mathbf{z}, \mathbf{W}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z}, \mathbf{W})p(\mathbf{w})p(\mathbf{z})}{p(\mathbf{x})}$$
$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{W}, \mathbf{z})p(\mathbf{W})p(\mathbf{z})d\mathbf{W}d\mathbf{z}$$

- Intractable to reach posterior distribution of both variables

Linear Latent Variable Model

1. Formulate joint distribution

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

Linear Latent Variable Model

1. Formulate joint distribution

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

2. Formulate marginal distribution over \mathbf{x}

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

Linear Latent Variable Model

1. Formulate joint distribution

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

2. Formulate marginal distribution over \mathbf{x}

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

3. Find maximum-likelihood solution to \mathbf{W}

$$\hat{\mathbf{W}} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{x} | \mathbf{W})$$

Linear Latent Variable Model

1. Formulate joint distribution

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

2. Formulate marginal distribution over \mathbf{x}

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

3. Find maximum-likelihood solution to \mathbf{W}

$$\hat{\mathbf{W}} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{x} | \mathbf{W})$$

4. Formulate posterior over the latent space

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

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

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

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

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

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

Linear Latent Variable Model

$$p(\mathbf{x}, \mathbf{z} | \mathbf{W}) = \mathcal{N} \left(\begin{bmatrix} \mathbb{E}[\mathbf{x}] \\ \mathbb{E}[\mathbf{z}] \end{bmatrix}, \begin{bmatrix} \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] & \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{z} - \mathbb{E}[\mathbf{z}])^T] \\ \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] & \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z}])(\mathbf{z} - \mathbb{E}[\mathbf{z}])^T] \end{bmatrix} \right)$$

- We can compute all the expectations above to figure out what the joint is
- Once we have the joint we can
 - pick out the marginal $p(\mathbf{x} | \mathbf{W})$
 - get the conditional $p(\mathbf{z} | \mathbf{x}, \mathbf{W})$

$$\mathbb{E}[\mathbf{z}] = \mathbf{0}$$

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W}\mathbf{z} + \mu + \epsilon] = \mathbb{E}[\mathbf{W}\mathbf{z}] + \mathbb{E}[\mu] + \mathbb{E}[\epsilon] \\ &= \mathbf{W}\mathbb{E}[\mathbf{z}] + \mathbb{E}[\mu] + \mathbb{E}[\epsilon] = \mathbf{W}\mathbf{0} + \mu + \mathbf{0} \\ &= \mu\end{aligned}$$

$$\begin{aligned}\mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] &= \mathbb{E}[(\mathbf{z} - \mathbf{0})(\mathbf{x} - \boldsymbol{\mu})^T] \\&= \mathbb{E}[\mathbf{z}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} - \boldsymbol{\mu})^T] \\&= \mathbb{E}[\mathbf{z}(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})^T] = \mathbb{E}[\mathbf{z}(\mathbf{W}\mathbf{z})^T + \mathbf{z}\boldsymbol{\epsilon}^T] \\&= \mathbb{E}[\mathbf{z}\mathbf{z}^T\mathbf{W}^T] + \mathbb{E}[\mathbf{z}]\mathbb{E}[\boldsymbol{\epsilon}] = \mathbb{E}[(\mathbf{z} - \mathbf{0})(\mathbf{z} - \mathbf{0})^T]\mathbf{W}^T + \mathbf{0} \cdot \mathbf{0} \\&= \mathbf{I}\mathbf{W}^T = \mathbf{W}^T\end{aligned}$$

Linear Latent Variable Model

$$\begin{aligned}\mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] &= \mathbb{E}[(\mathbf{W}\mathbf{z} + \mu + \epsilon - \mu)(\mathbf{W}\mathbf{z} + \mu + \epsilon - \mu)^T] \\&= \mathbb{E}[(\mathbf{W}\mathbf{z} + \epsilon)(\mathbf{W}\mathbf{z} + \epsilon)^T] = \mathbb{E}[\mathbf{W}\mathbf{z}(\mathbf{W}\mathbf{z})^T + \mathbf{W}\mathbf{z}\epsilon^T + \epsilon\mathbf{W}\mathbf{z}^T + \epsilon\epsilon^T] \\&= \mathbb{E}[\mathbf{W}\mathbf{z}\mathbf{z}^T\mathbf{W}^T] + \mathbb{E}[\mathbf{W}\mathbf{z}\epsilon^T] + \mathbb{E}[\epsilon(\mathbf{W}\mathbf{z})^T] + \mathbb{E}[\epsilon\epsilon^T] \\&= \mathbf{W}\mathbb{E}[\mathbf{z}\mathbf{z}^T]\mathbf{W}^T + \mathbf{W}\mathbb{E}[\mathbf{z}]\mathbb{E}[\epsilon] + \mathbb{E}[\epsilon]\mathbb{E}[\mathbf{z}^T]\mathbf{W}^T + \mathbb{E}[(\epsilon - 0)(\epsilon - 0)^T] \\&= \mathbf{W}\mathbf{I}\mathbf{W}^T + \mathbf{W}\mathbf{0} + \mathbf{0}\mathbf{W}^T + \sigma^2\mathbf{I} \\&= \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}\end{aligned}$$

Linear Latent Variable Model

Joint

$$p(\mathbf{x}, \mathbf{z} | \mathbf{W}) = \mathcal{N} \left(\begin{bmatrix} \mu \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I} & \mathbf{W} \\ \mathbf{W}^T & \mathbf{I} \end{bmatrix} \right)$$

Marginal

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

Posterior

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

$$\log p(\mathbf{x}|\mathbf{W})$$

- find stationary with respect to each variable gives Maximum likelihood solution to \mathbf{W}
- you can also do the same with μ and σ^2

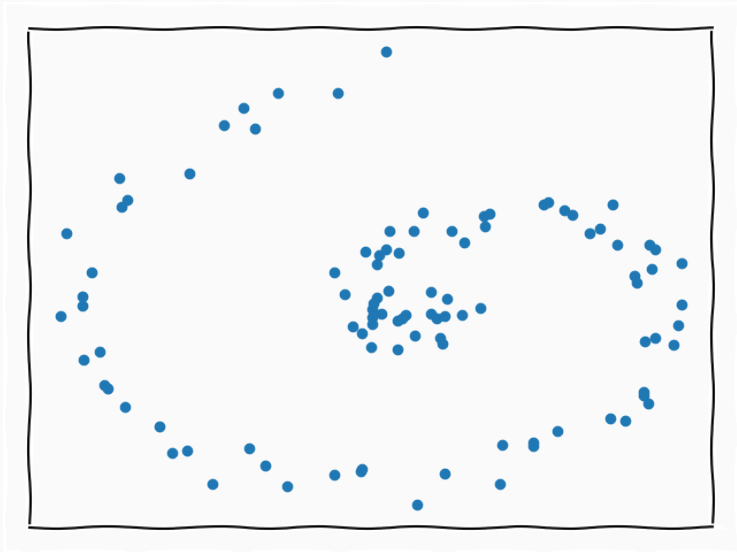
$$\log p(\mathbf{x}|\mathbf{W}) = \sum_{n=1}^N \log p(\mathbf{x}_n|\mathbf{W})$$

- find stationary with respect to each variable gives Maximum likelihood solution to \mathbf{W}
- you can also do the same with μ and σ^2

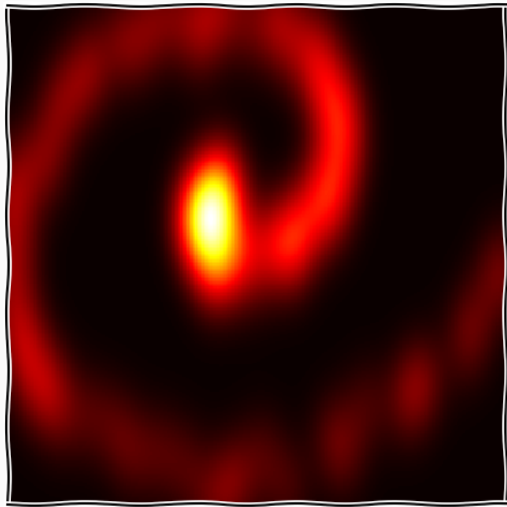
$$\begin{aligned}\log p(\mathbf{x}|\mathbf{W}) &= \sum_{n=1}^N \log p(\mathbf{x}_n|\mathbf{W}) \\ &= -\frac{ND}{2} \log(2\pi) - \frac{N}{2} \log |\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}| \\ &\quad - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})^{-1}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

- find stationary with respect to each variable gives Maximum likelihood solution to \mathbf{W}
- you can also do the same with $\boldsymbol{\mu}$ and σ^2

Example



Example II



$$\mathbf{V}\mathbf{V}^T = \mathbf{x}^T\mathbf{x}$$

$$\mathbf{z} = \sum_i^d \mathbf{x}\mathbf{V}_i$$

- The above is the solution if $\sigma^2 \rightarrow 0$

²[\[2\]](#)

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

1. Specify your statistical model over sample space \mathcal{Y} , relationship between "parameters" and observations

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

1. Specify your statistical model over sample space \mathcal{Y} , relationship between "parameters" and observations

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

2. Formulate your likelihood

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

1. Specify your statistical model over sample space \mathcal{Y} , relationship between "parameters" and observations

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

2. Formulate your likelihood

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

3. Formulate your belief in the "setting" of the model

$$p(\theta)$$

1. Specify your statistical model over sample space \mathcal{Y} , relationship between "parameters" and observations

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

2. Formulate your likelihood

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

3. Formulate your belief in the "setting" of the model

$$p(\theta)$$

4. Acquire data

1. Specify your statistical model over sample space \mathcal{Y} , relationship between "parameters" and observations

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

2. Formulate your likelihood

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

3. Formulate your belief in the "setting" of the model

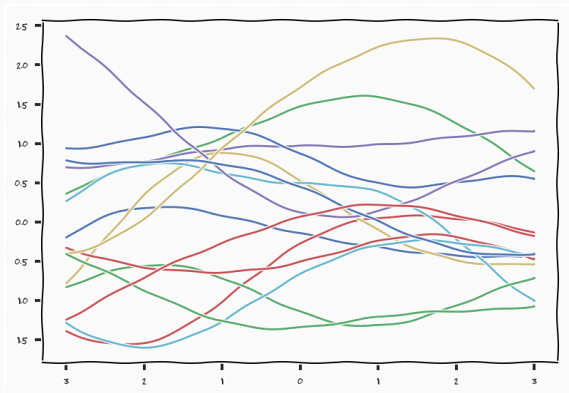
$$p(\theta)$$

4. Acquire data
5. Derive your updated belief, derive knowledge from data

$$p(\theta|\mathbf{Y}) = \frac{p(\mathbf{Y}|\theta)p(\theta)}{\int p(\mathbf{Y}|\theta)p(\theta)}$$

- Unsupervised learning is a misnomer, there is no such thing, you have to have beliefs in order to learn.
- Think about unsupervised learning as "less supervised" learning, you have to have stronger beliefs

More Interesting Priors



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

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
- Tomorrow and next three lectures

eof

References



Christopher M. Bishop.

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

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



Charles Spearman.

" General Intelligence," Objectively Determined and Measured.

The American Journal of Psychology, 15(2):201–292, 1904.