

# DD2434 - Advanced Machine Learning

## Representation Learning

Carl Henrik Ek  
{chek}@csc.kth.se

Royal Institute of Technology

November 11th, 2015



## Last Lecture

- Gaussian Processes
  - ▶ Prior over the space of functions
  - ▶ Posterior
  - ▶ Marginal Likelihood
  - ▶ Learning





# Regression

Regression model,

$$\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon \quad (1)$$

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

Introduce  $f_i$  as *instansiation* of function,

$$f_i = f(\mathbf{x}_i), \quad (3)$$

as a new random variable.

# Regression

Model,

$$p(\mathbf{Y}, \mathbf{f}, \mathbf{X}, \boldsymbol{\theta}) = p(\mathbf{Y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})p(\mathbf{X})p(\boldsymbol{\theta}) \quad (4)$$

Want to “push”  $\mathbf{X}$  through a mapping  $f$  of which we are uncertain,

$$p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}), \quad (5)$$

prior over instantiations of function.

# Gaussian Processes<sup>1</sup>

$$p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) \sim \mathcal{GP}(\mu(\mathbf{X}), k(\mathbf{X}, \mathbf{X})) \quad (6)$$

## Defenition

A Gaussian Process is an infinite collection of random variables who **any** subset is jointly gaussian. The process is specified by a mean function  $\mu(\cdot)$  and a co-variance function  $k(\cdot, \cdot)$

$$f \sim \mathcal{GP}(\mu(\cdot), k(\cdot, \cdot)) \quad (7)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

$$p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) \sim \mathcal{GP}(\mu(\mathbf{X}), k(\mathbf{X}, \mathbf{X})) \quad (8)$$

$$\mathbf{y}_i = f_i + \epsilon \quad (9)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (10)$$

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{Y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} \quad (11)$$

## Connection to Distribution

$\mathcal{GP}$  is infinite, but we only observe finite amount of data. This means conditioning on a subset of the data, the  $\mathcal{GP}$  is a just a Gaussian distribution, which is self-conjugate.

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

## The mean function

- Function of only the input location
- What do I expect the function value to be **only** accounting for the input location
- We will assume this to be constant

## The co-variance function

- Function of **two** input locations
- How should the information from other locations with **known** function value observations effect my estimate
- Encodes the behavior of the function

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

## The mean function

- Function of only the input location
- What do I expect the function value to be **only** accounting for the input location
- We will assume this to be constant

## The co-variance function

- Function of **two** input locations
- How should the information from other locations with **known** function value observations effect my estimate
- Encodes the behavior of the function

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

## The mean function

- Function of only the input location
- What do I expect the function value to be **only** accounting for the input location
- We will assume this to be constant

## The co-variance function

- Function of **two** input locations
- How should the information from other locations with **known** function value observations effect my estimate
- Encodes the behavior of the function

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

## The Prior

$$p(f|\mathbf{X}, \boldsymbol{\theta}) = \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (12)$$

$$\mu(\mathbf{x}) = \mathbf{0} \quad (13)$$

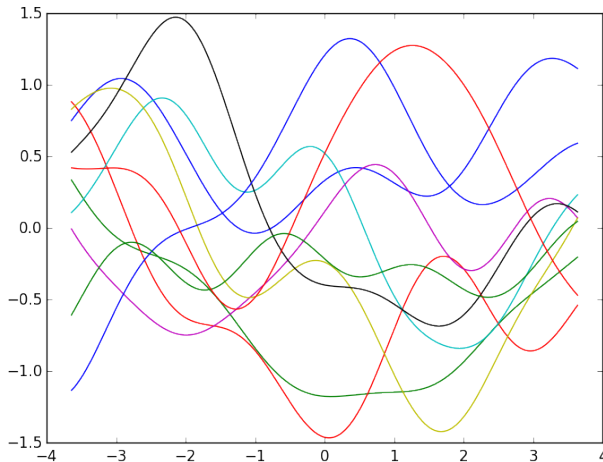
$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 e^{-\frac{1}{2\ell^2}(\mathbf{x}_i - \mathbf{x}_j)^T(\mathbf{x}_i - \mathbf{x}_j)} \quad (14)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

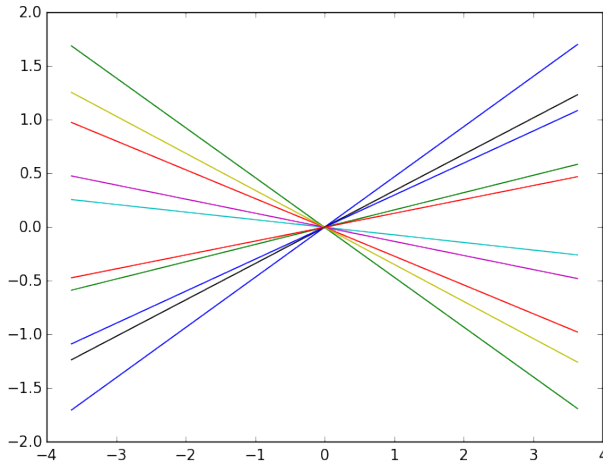


# Gaussian Processes<sup>1</sup>



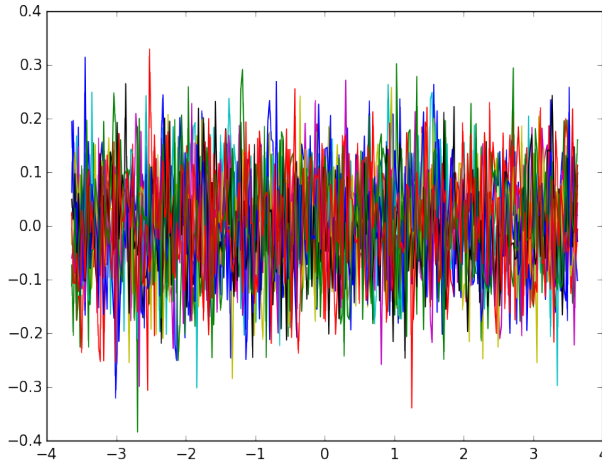
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



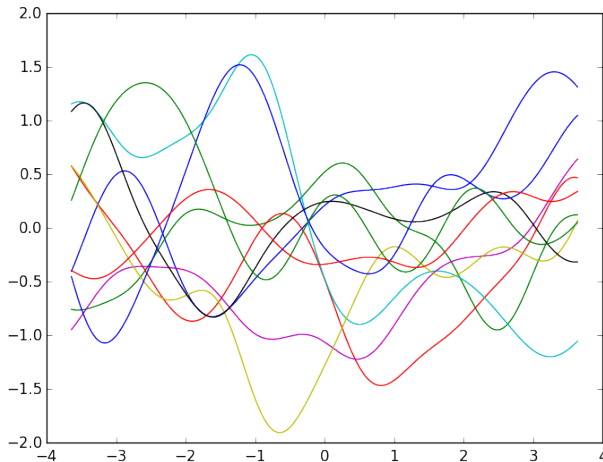
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



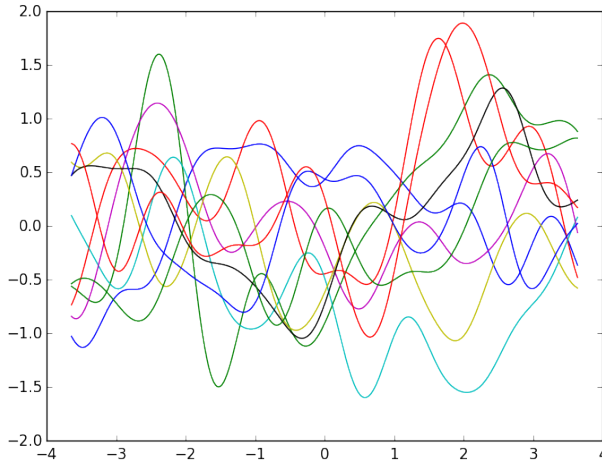
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



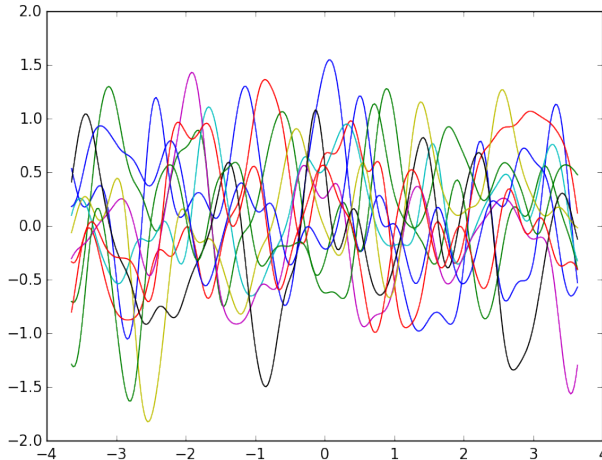
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



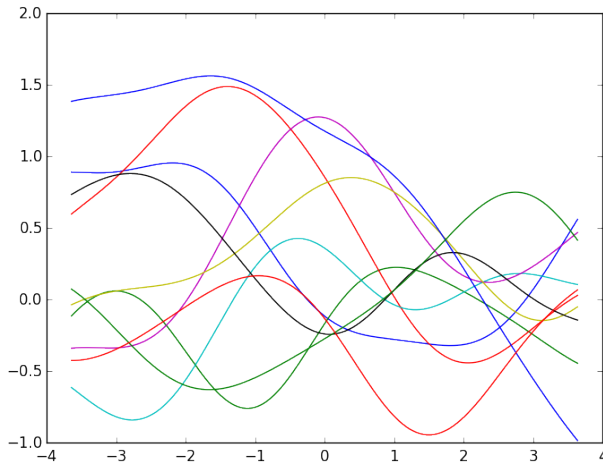
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



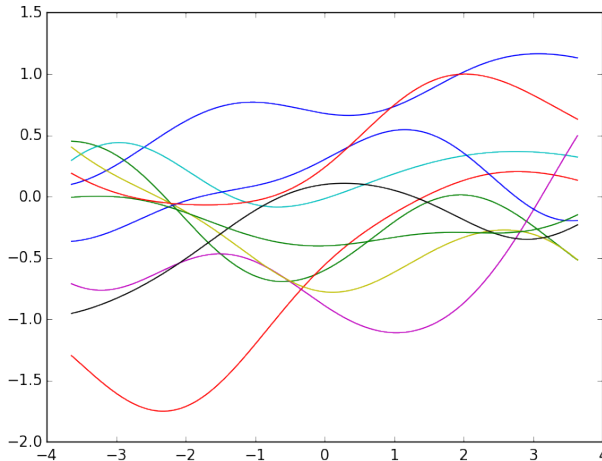
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

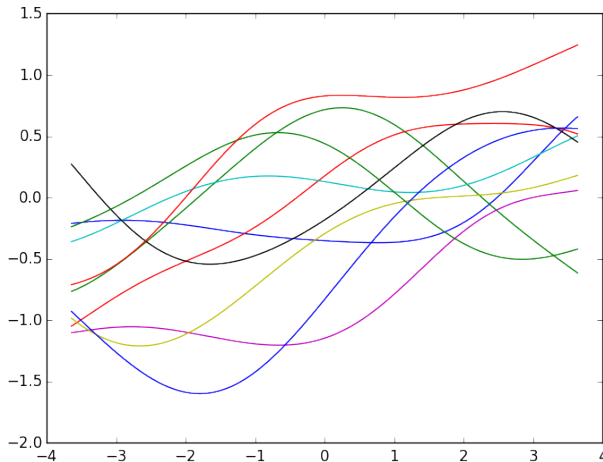
# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

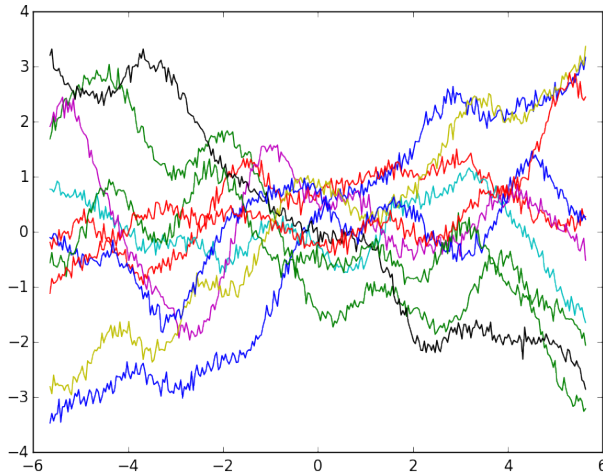


# Gaussian Processes<sup>1</sup>



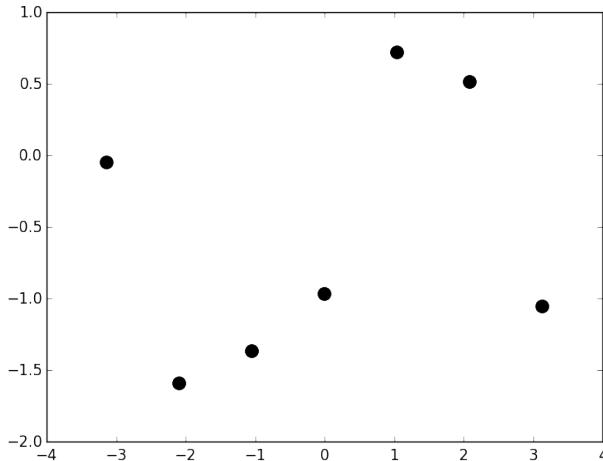
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

The (predictive) Posterior

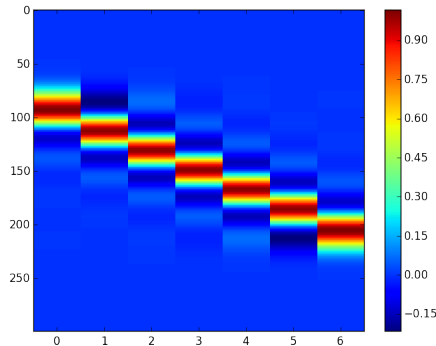
$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right) \quad (15)$$

$$\begin{aligned} p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}, \boldsymbol{\theta}) &= \mathcal{N}(k(\mathbf{x}_*, \mathbf{X})^\top K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f}, \\ &\quad k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^\top K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{x}_*)) \end{aligned} \quad (16)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

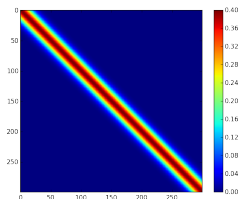


$$k(\mathbf{x}_*, \mathbf{X})^T K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f} \quad (17)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

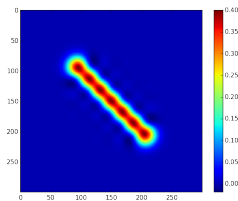
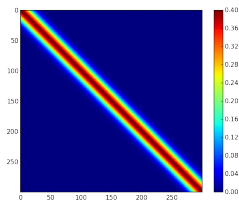


$$k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^T K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{x}_*) \quad (18)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>

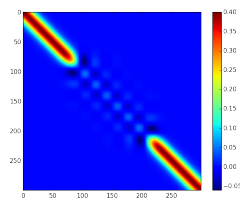
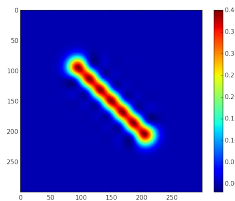
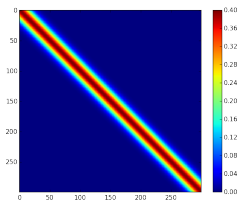


$$k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^T K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{x}_*) \quad (19)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



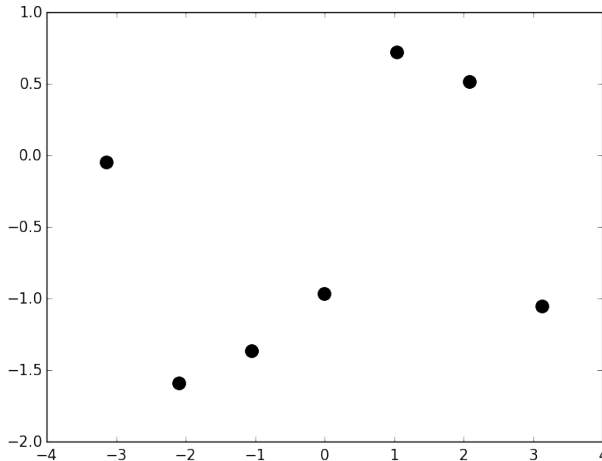
$$k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^T K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{x}_*) \quad (20)$$

---

<sup>1</sup>Bishop 2006, p. 6.4.2

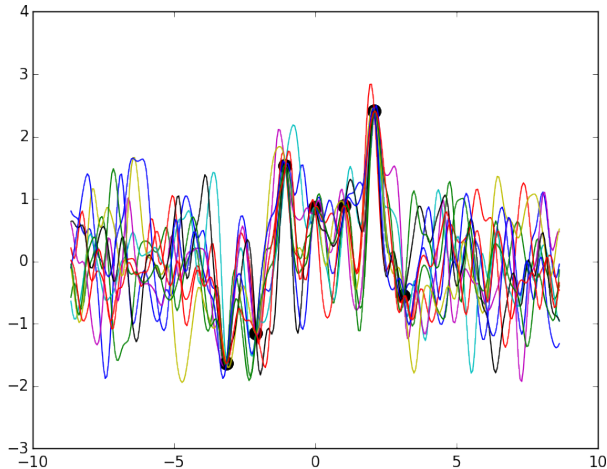


# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

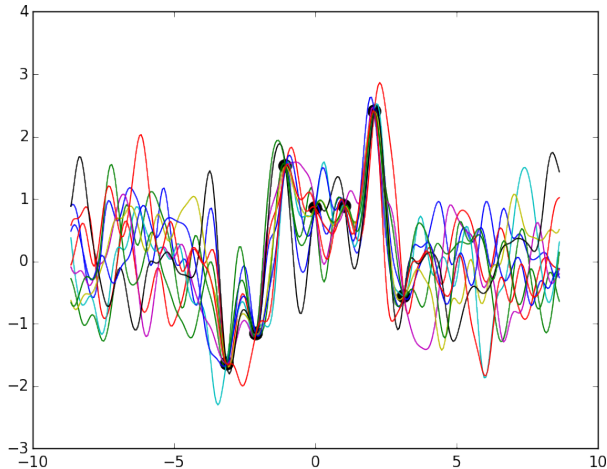
# Gaussian Processes<sup>1</sup>



---

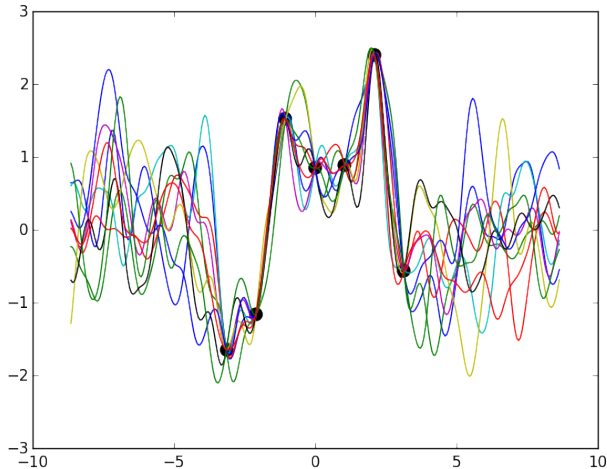
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



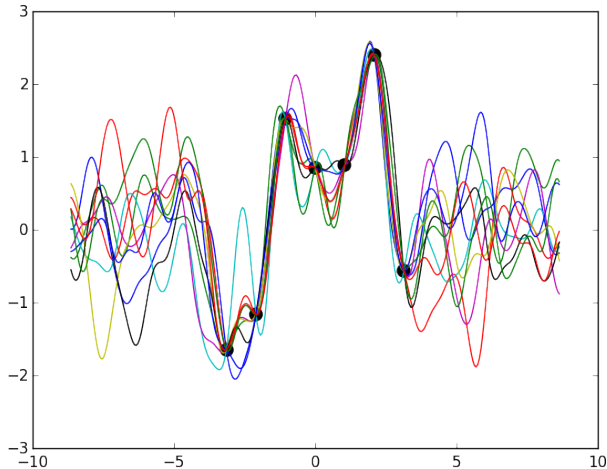
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



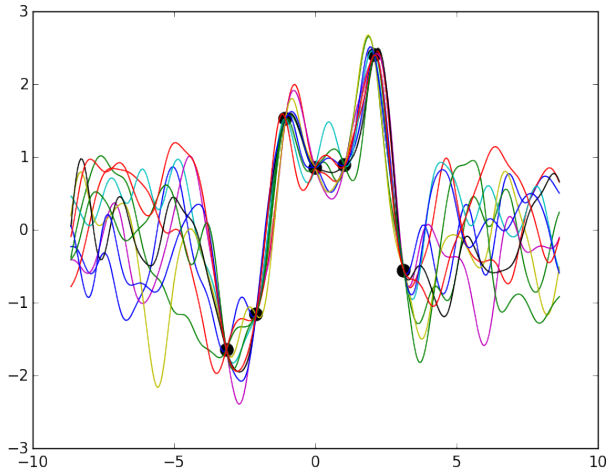
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



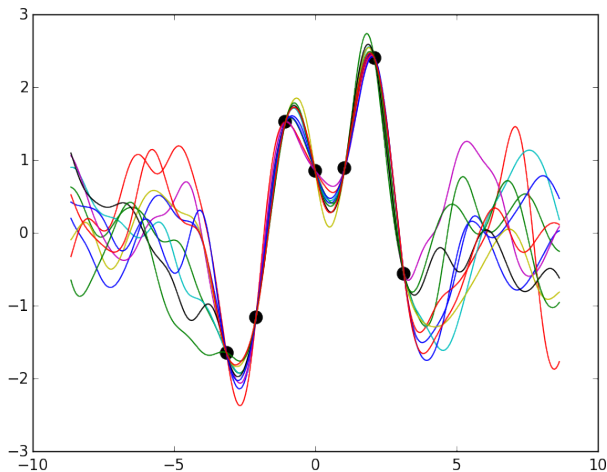
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



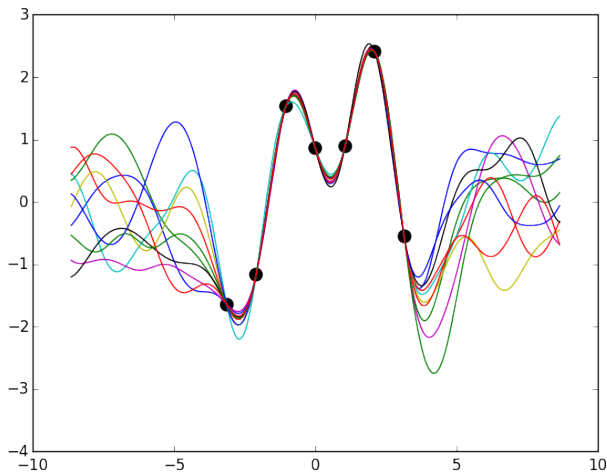
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

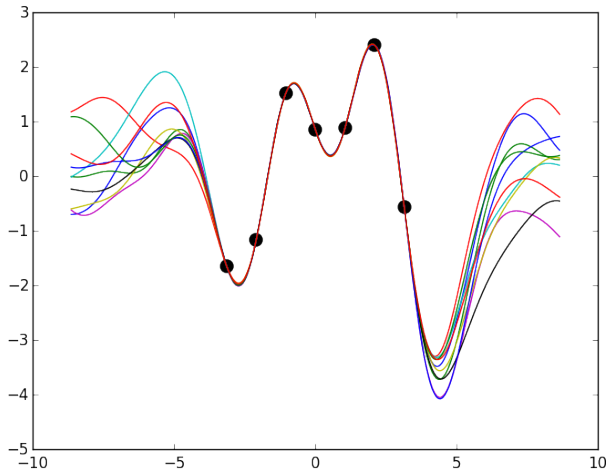
# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

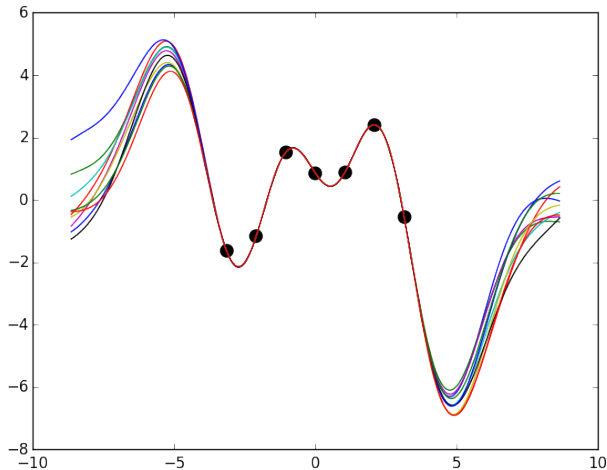


# Gaussian Processes<sup>1</sup>



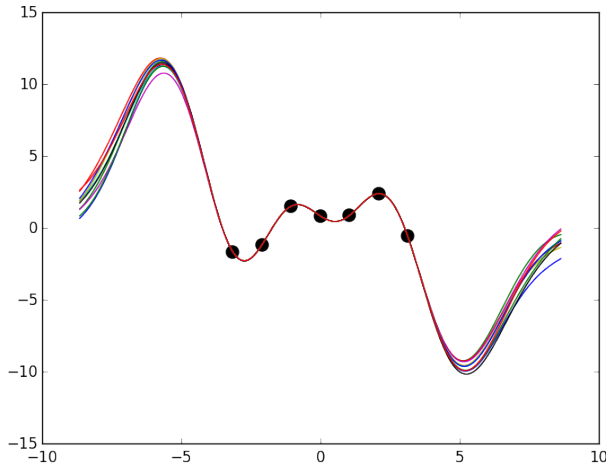
<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

# Gaussian Processes<sup>1</sup>



<sup>1</sup>Bishop 2006, p. 6.4.2

# Learning in Gaussian Processes<sup>2</sup>

## Hyper-parameters

- Prior has parameters
  - ▶ referred to as *hyper*-parameters
  - ▶ SE have lengthscale and variance
- Learning in  $\mathcal{GP}$ s implies inferring hyper-parameters from the model

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{Y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} \quad (21)$$

## Marginal Likelihood

- We are not interested in  $\mathbf{f}$  directly
- Marginalise out  $\mathbf{f}$ !

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{Y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} \quad (22)$$

## Marginal Likelihood

- We are not interested in  $\mathbf{f}$  directly
- Marginalise out  $\mathbf{f}$ !

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{Y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} \quad (23)$$

## Marginal Likelihood

- We are not interested in  $\mathbf{f}$  directly
- Marginalise out  $\mathbf{f}$ !

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

$$\operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} -\log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})) = \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \quad (24)$$

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi) \quad (25)$$

## Type-II Maximum Likelihood

- Can be minimised using gradient based methods
- Data-fit:  $\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$
- Complexity:  $\frac{1}{2} \log |\mathbf{K}|$

---

<sup>2</sup>Bishop 2006, p. 6.4.3



# Learning in Gaussian Processes<sup>2</sup>

$$\operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} -\log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})) = \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \quad (26)$$

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi) \quad (27)$$

## Type-II Maximum Likelihood

- Can be minimised using gradient based methods
- Data-fit:  $\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$
- Complexity:  $\frac{1}{2} \log |\mathbf{K}|$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

$$\operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}} -\log(p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})) = \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \quad (28)$$

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi) \quad (29)$$

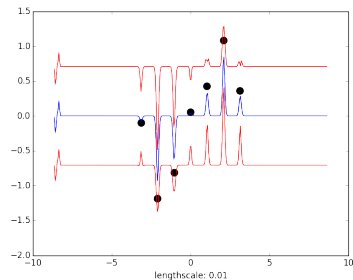
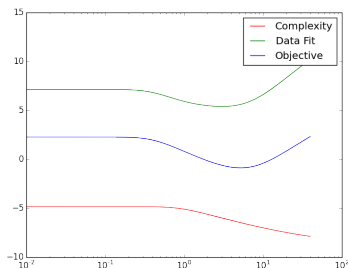
## Type-II Maximum Likelihood

- Can be minimised using gradient based methods
- Data-fit:  $\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$
- Complexity:  $\frac{1}{2} \log |\mathbf{K}|$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

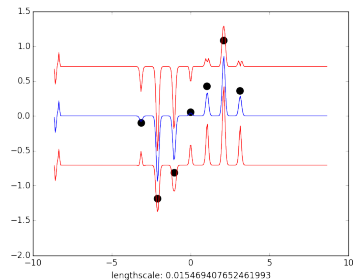
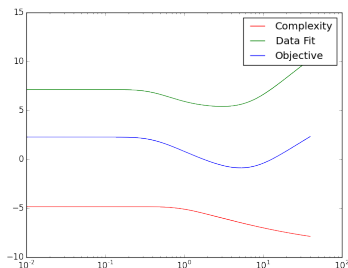


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

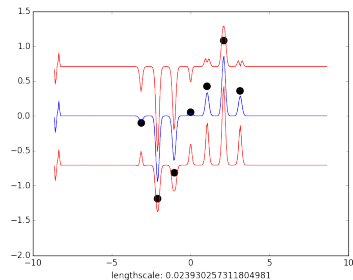
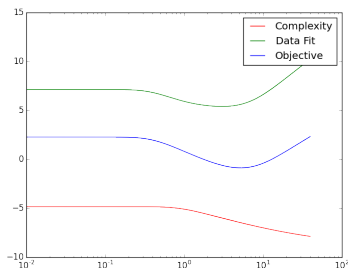


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

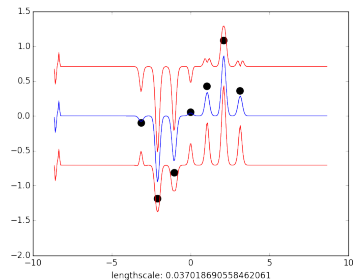
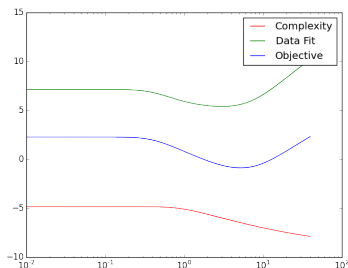


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

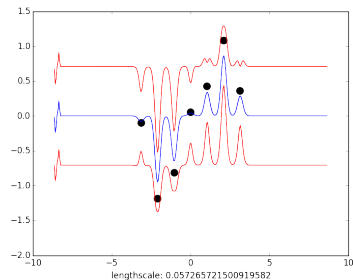
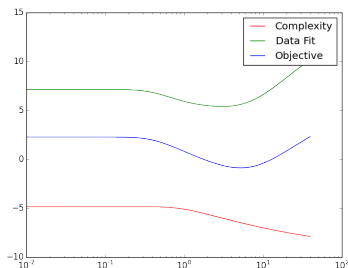


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

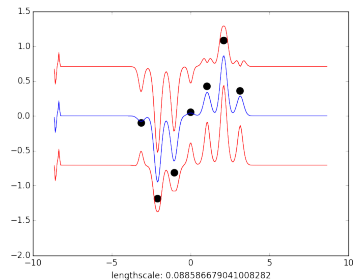
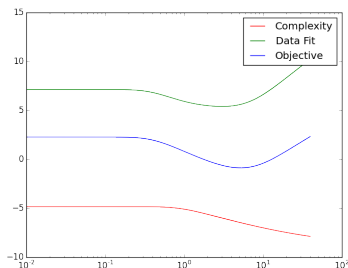


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>



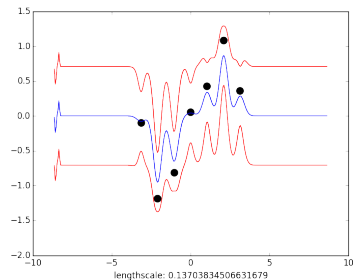
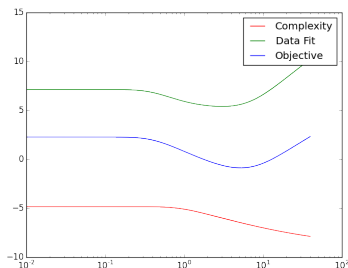
$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3



# Learning in Gaussian Processes<sup>2</sup>

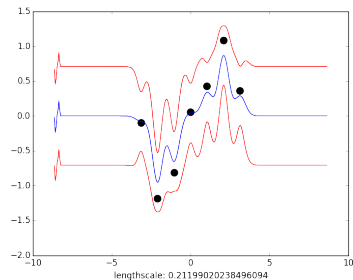
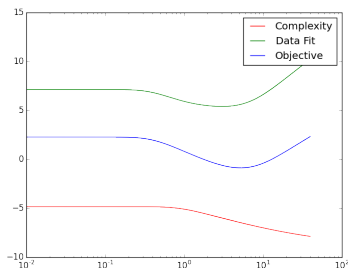


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

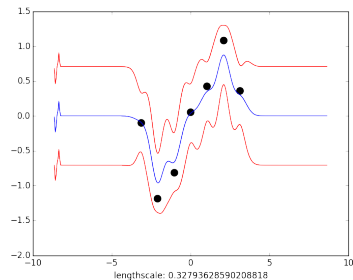
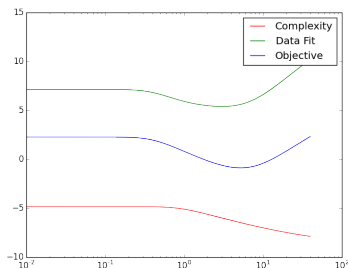


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

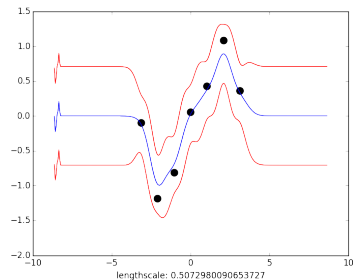
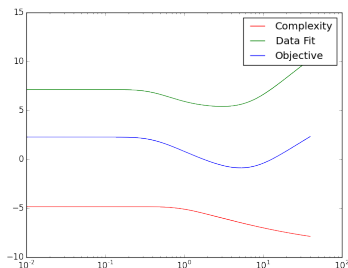


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

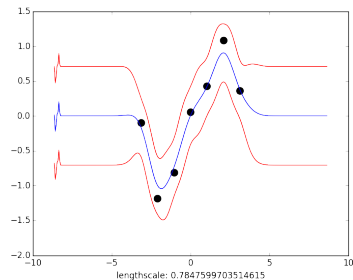
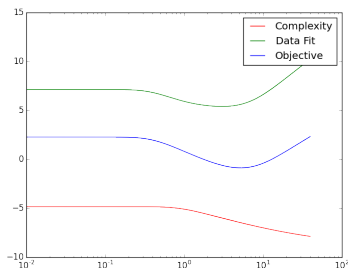


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

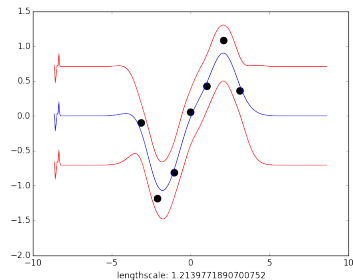
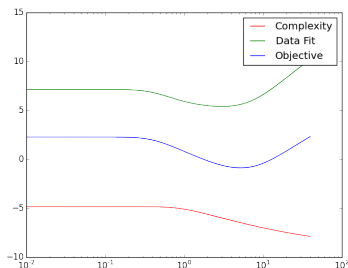


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

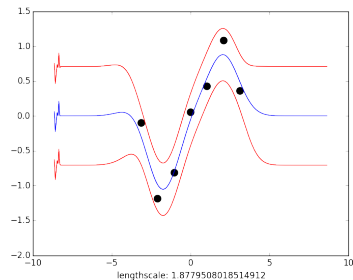
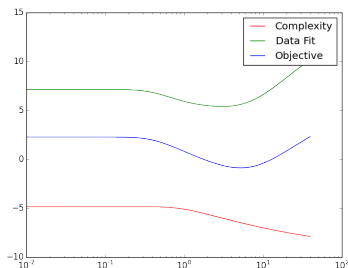


$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

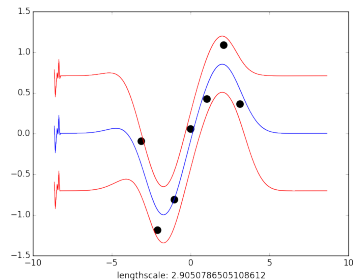
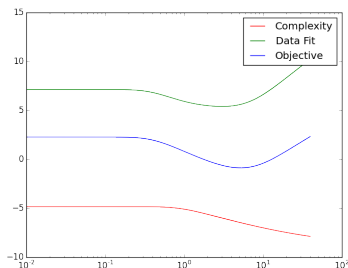


$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>



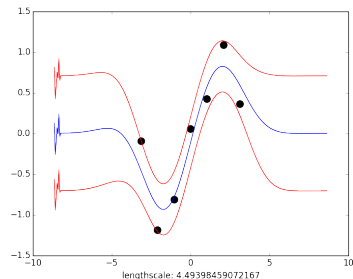
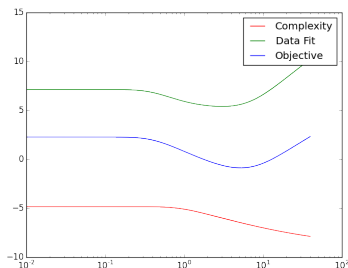
$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3



# Learning in Gaussian Processes<sup>2</sup>

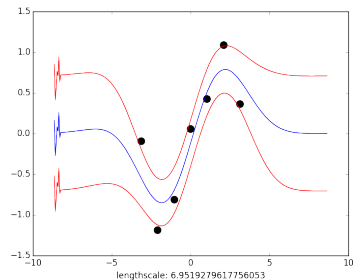
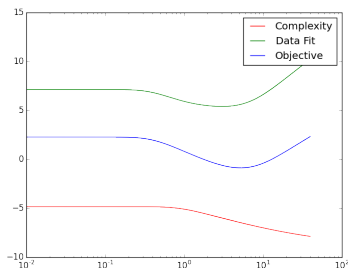


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

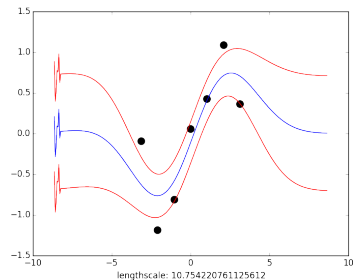
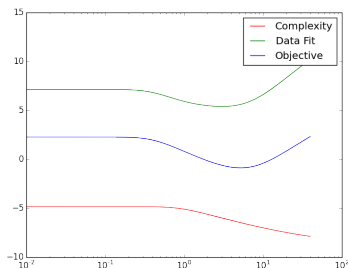


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

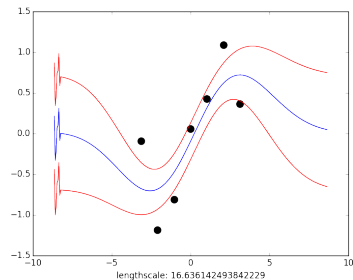
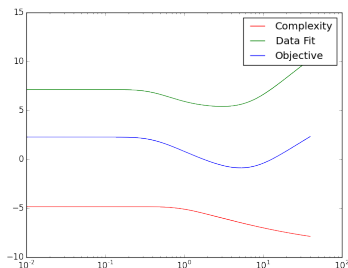


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

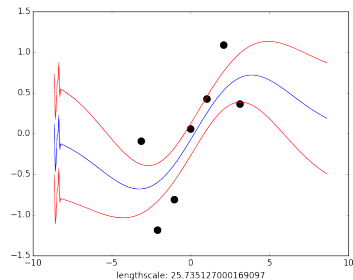
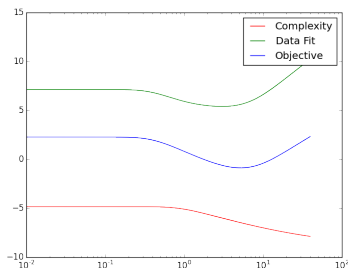


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

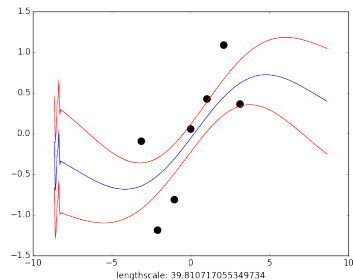
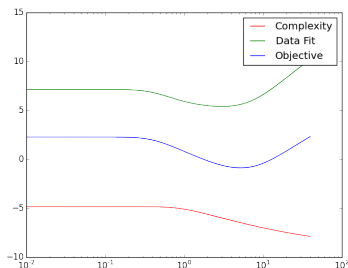


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>

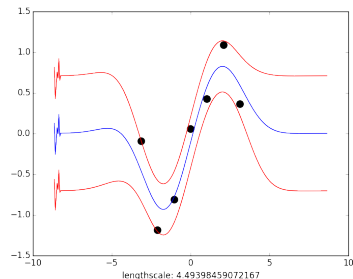
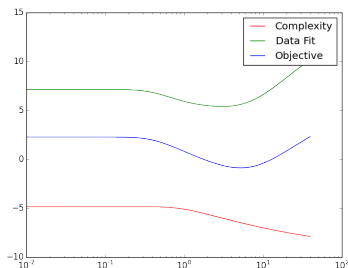


$$\mathcal{L}(\theta) = \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}| + \frac{N}{2} \log(2\pi)$$

---

<sup>2</sup>Bishop 2006, p. 6.4.3

# Learning in Gaussian Processes<sup>2</sup>



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

---

<sup>2</sup>Bishop 2006, p. 6.4.3

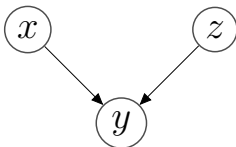
Introduction

Recap

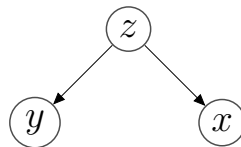
Representation Learning



# Graphical Models<sup>3</sup>



$$p(x, y, z) = p(y|x, z)p(x)p(z)$$



$$p(x, y, z) = p(y|z)p(x|z)p(z)$$

$$p(\{x_i\}_{i=1}^N) = \prod_{i=1}^N p(x_i|\text{pa}_i) \quad (30)$$

---

<sup>3</sup>Bishop 2006, pp. 8.0, 8.1.

# Latent Variable Models<sup>4</sup>

## Machine Learning

- What is our task?
- $p(y)$
- Unobservables
  - ▶ Latent variables
- Explaining away



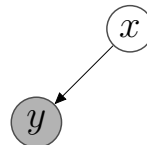
---

<sup>4</sup>Bishop 2006, p. 364.

# Latent Variable Models<sup>4</sup>

## Machine Learning

- What is our task?
- $p(y)$
- Unobservables
  - ▶ Latent variables
- Explaining away



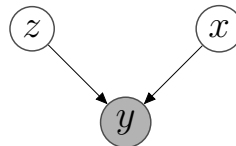
---

<sup>4</sup>Bishop 2006, p. 364.

# Latent Variable Models<sup>4</sup>

## Machine Learning

- What is our task?
- $p(y)$
- Unobservables
  - ▶ Latent variables
- Explaining away



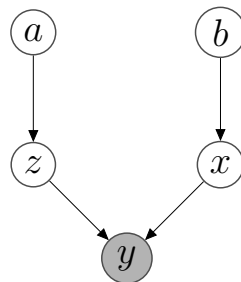
---

<sup>4</sup>Bishop 2006, p. 364.

# Latent Variable Models<sup>4</sup>

## Machine Learning

- What is our task?
- $p(y)$
- Unobservables
  - ▶ Latent variables
- Explaining away



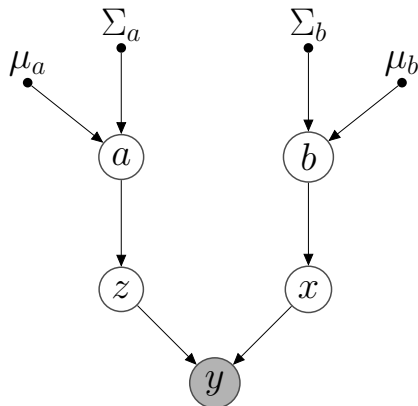
---

<sup>4</sup>Bishop 2006, p. 364.

# Latent Variable Models<sup>4</sup>

## Machine Learning

- What is our task?
- $p(y)$
- Unobservables
  - ▶ Latent variables
- Explaining away



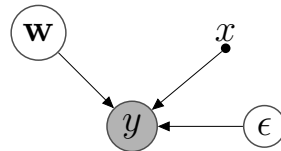
---

<sup>4</sup>Bishop 2006, p. 364.

# Latent Variable Models<sup>4</sup>

## Machine Learning

- What is our task?
- $p(y)$
- Unobservables
  - ▶ Latent variables
- Explaining away



$$\mathbf{y}_i = \mathbf{w}\mathbf{x}_i + \epsilon \quad (31)$$

---

<sup>4</sup>Bishop 2006, p. 364.

# Latent Variable Models<sup>4</sup>

## Latent Variables<sup>a</sup>

---

<sup>a</sup>Bishop 2006, p. 366.

*“The primary role of the latent variable is to allow a complicated distribution over the observed variables be represented in terms of a model constructed from a simpler (typically exponential family) conditional distribution.”*

---

<sup>4</sup>Bishop 2006, p. 364.



# Latent Variable Models<sup>4</sup>



---

<sup>4</sup>Bishop 2006, p. 364.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.



# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.



# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.



# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>



---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Generative Models<sup>5</sup>

- We want to model this data, i.e.  $p(y)$
- $\mathbf{y}_i \in \mathbb{R}^{1400 \times 931 \times 3 = 3910200}$
- Latent variable representation
- Conditional distribution parametrised by latent variable



---

<sup>5</sup>Bishop 2006, p. 8.1.2.



# Generative Models<sup>5</sup>

```
1 t = np.sort(5*np.random.randn(30))
2 for i in range(0,len(t)):
3     pos[0] = int(radius*np.sin(t[i])+offset[0]/2)
4     pos[1] = int(radius*np.cos(t[i])+offset[1]/2)
5     background_draw.paste(stella,pos,stella)
6     background_draw.save('lvm_'+str(i)+'.png')
7     background_draw = copy.deepcopy(background)
```

$$p(\mathbf{Y}|\mathbf{X}) \quad (32)$$

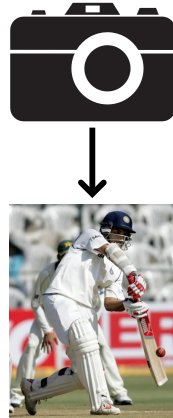
---

<sup>5</sup>Bishop 2006, p. 8.1.2.

# Sensory Data

## What we are doing

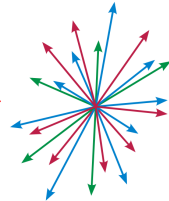
- Sensory representation
  - ▶ Capturing process
  - ▶ Pixels, Waveforms
- Degrees of freedom and dimensionality



# Sensory Data

## What we are doing

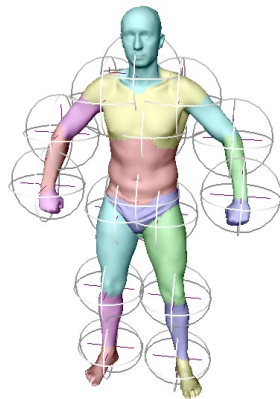
- Sensory representation
  - ▶ Capturing process
  - ▶ Pixels, Waveforms
- Degrees of freedom and dimensionality



# Sensory Data

## What we are doing

- Sensory representation
  - ▶ Capturing process
  - ▶ Pixels, Waveforms
- Degrees of freedom and dimensionality



## Outline

- Re-visit PCA
- PCA as a Latent Variable Model
- Factor Analysis
- Example of Intractability



# Re-visit: Principal Component Analysis

- Given data  $\mathbf{Y}$  project to directions of maximum variance
- Provides no uncertainty

$$\operatorname{argmax}_{\mathbf{v}} \sigma(\mathbf{Y}\mathbf{v}, \mathbf{Y}\mathbf{v}) \quad (33)$$

$$\mathbf{v}^T \mathbf{Y}^T \mathbf{Y} \mathbf{v} \quad (34)$$

$$\text{subject to: } \mathbf{v}^T \mathbf{v} = 1 \quad (35)$$

# Latent Variable Models<sup>6</sup>

$$p(\mathbf{Y}) \tag{36}$$

- We have observed some data  $\mathbf{Y}$
- Lets assume that  $\mathbf{Y} \in \mathbb{R}^{N \times d}$  have been generated from  $\mathbf{X} \in \mathbb{R}^{N \times q}$
- $\mathbf{X}$  - latent variable
- $f$  - generative mapping

---

<sup>6</sup>Bishop 2006, p. 8.1.2.

# Latent Variable Models<sup>6</sup>

$$p(\mathbf{Y}|f, \mathbf{X}) \quad (37)$$

$$\mathbf{f} : \mathbf{X} \rightarrow \mathbf{Y} \quad (38)$$

- We have observed some data  $\mathbf{Y}$
- Lets assume that  $\mathbf{Y} \in \mathbb{R}^{N \times d}$  have been generated from  $\mathbf{X} \in \mathbb{R}^{N \times q}$
- $\mathbf{X}$  - latent variable
- $f$  - generative mapping

---

<sup>6</sup>Bishop 2006, p. 8.1.2.



# Latent Variable Models<sup>6</sup>

$$p(\mathbf{Y}|f, \mathbf{X}) \quad (39)$$

$$\mathbf{f} : \mathbf{X} \rightarrow \mathbf{Y} \quad (40)$$

- We have observed some data  $\mathbf{Y}$
- Lets assume that  $\mathbf{Y} \in \mathbb{R}^{N \times d}$  have been generated from  $\mathbf{X} \in \mathbb{R}^{N \times q}$
- $\mathbf{X}$  - latent variable
- $f$  - generative mapping

---

<sup>6</sup>Bishop 2006, p. 8.1.2.

# Latent Variable Models<sup>6</sup>

$$p(\mathbf{Y}|f, \mathbf{X}) \quad (41)$$

$$\mathbf{f} : \mathbf{X} \rightarrow \mathbf{Y} \quad (42)$$

- We have observed some data  $\mathbf{Y}$
- Lets assume that  $\mathbf{Y} \in \mathbb{R}^{N \times d}$  have been generated from  $\mathbf{X} \in \mathbb{R}^{N \times q}$
- $\mathbf{X}$  - latent variable
- $f$  - generative mapping

---

<sup>6</sup>Bishop 2006, p. 8.1.2.

# Linear Latent Variable Models<sup>7</sup>

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}) = \prod_i^N p(\mathbf{y}_i|\mathbf{W}, \mathbf{x}_i) \quad (43)$$

## Regression

- Regression without inputs?
- Solve the task: Given some data
  - ▶ a representation of this data
  - ▶ and a mapping that have generated the

---

<sup>7</sup>Bishop 2006, p. 12.2.0.

# WTF?

## The strength of Priors

- Encodes prior belief
- This can also be seen as a preference
  - ▶ Given several perfectly valid solutions which one do i prefer
  - ▶ Regularises solution space
- Latent variable models what do we prefer?

# Factor Analysis<sup>8</sup>

$$\mathbf{y}_i = \mathbf{W}\mathbf{x}_i + \epsilon \quad (44)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \Psi) \quad (45)$$

- Assume the generating mapping to be linear
- Assume  $\Psi$  diagonal
- For regression we assumed that we knew the inputs  $\mathbf{X}$
- Now we do not

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

$$\mathbf{y}_i = \mathbf{W}\mathbf{x}_i + \epsilon \quad (46)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \Psi) \quad (47)$$

- Assume the generating mapping to be linear
- Assume  $\Psi$  diagonal
- For regression we assumed that we knew the inputs  $\mathbf{X}$
- Now we do not

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

$$\mathbf{y}_i = \mathbf{W}\mathbf{x}_i + \epsilon \quad (48)$$

$$p(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{W}\mathbf{x}_i, \boldsymbol{\Psi}) \quad (49)$$

$$p(\mathbf{x}_i) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (50)$$

- Assume the generating mapping to be linear
- Assume  $\boldsymbol{\Psi}$  diagonal
- For regression we assumed that we knew the inputs  $\mathbf{X}$
- Now we do not  $\Rightarrow$  specify a prior

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

$$\begin{aligned} p(\mathbf{y}_i | \boldsymbol{\theta}) &= \int p(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}) p(\mathbf{x}_i) d\mathbf{x}_i = \int \mathcal{N}(\mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \\ &= \mathcal{N}(\mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T) \end{aligned}$$

- $\mathbf{X}$  and  $\mathbf{W}$  are related
- Integrate out  $\mathbf{X}$ 
  - ▶ pick  $\boldsymbol{\mu}_0 = 0, \boldsymbol{\Sigma}_0 = \mathbf{I}$
- Low dimensional density model of  $\mathbf{Y}$ 
  - ▶ rank of  $\mathbf{W}\mathbf{W}^T$  dimensionality of  $\mathbf{X}$

<sup>8</sup>Bishop 2006, p. 12.2.4.



# Factor Analysis<sup>8</sup>

$$\begin{aligned} p(\mathbf{y}_i | \boldsymbol{\theta}) &= \int p(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}) p(\mathbf{x}_i) d\mathbf{x}_i = \int \mathcal{N}(\mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \\ &= \mathcal{N}(\mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T) \\ &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T) \end{aligned} \tag{51}$$

- $\mathbf{X}$  and  $\mathbf{W}$  are related
- Integrate out  $\mathbf{X}$ 
  - ▶ pick  $\boldsymbol{\mu}_0 = 0, \boldsymbol{\Sigma}_0 = \mathbf{I}$
- Low dimensional density model of  $\mathbf{Y}$ 
  - ▶ rank of  $\mathbf{W}\mathbf{W}^T$  dimensionality of  $\mathbf{X}$

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

$$\begin{aligned} p(\mathbf{y}_i | \boldsymbol{\theta}) &= \int p(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}) p(\mathbf{x}_i) d\mathbf{x}_i = \int \mathcal{N}(\mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \\ &= \mathcal{N}(\mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T) \\ &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T) \end{aligned} \tag{52}$$

- $\mathbf{X}$  and  $\mathbf{W}$  are related
- Integrate out  $\mathbf{X}$ 
  - ▶ pick  $\boldsymbol{\mu}_0 = 0, \boldsymbol{\Sigma}_0 = \mathbf{I}$
- Low dimensional density model of  $\mathbf{Y}$ 
  - ▶ rank of  $\mathbf{W}\mathbf{W}^T$  dimensionality of  $\mathbf{X}$

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

$$\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R} \quad (53)$$

$$p(\mathbf{y}_i|\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T) \quad (54)$$

$$= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T) \quad (55)$$

$$(56)$$

## Identifiability

- The marginal likelihood is invariant to a rotation
  - ▶ no unique solution
  - ▶ model is the same but interpretation tricky

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

$$\mathbf{W}_{ML} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{Y}|\boldsymbol{\theta}) \quad (57)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (58)$$

$$\mathbf{W}_{ML} = \mathbf{U}_q(\Lambda - \sigma^2 \mathbf{I})^{\frac{1}{2}} \quad (59)$$

$$\mathbf{S} = \mathbf{U}\Lambda\mathbf{U}^T \quad (60)$$

## Probabilistic PCA

- Dimensions of  $\mathbf{Y}$  independent given  $\mathbf{X}$ 
  - ▶  $\mathbf{W}$  orthogonal matrix  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

## Summary

- Factor Analysis is a linear continuous latent variable model
- Solution not unique
- PCA is Factor Analysis with two assumptions
  - ▶ factor loadings orthogonal  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$
  - ▶ noise free case  $\epsilon = \lim_{\sigma^2 \rightarrow 0} \sigma^2 \mathbf{I}$
- PCA is incredibly useful but its important to know what you are assuming, the probabilistic formulation allows you to do just that

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

# Factor Analysis<sup>8</sup>

## Summary

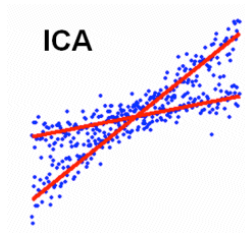
- Factor Analysis is a linear continuous latent variable model
- Solution not unique
- PCA is Factor Analysis with two assumptions
  - ▶ factor loadings orthogonal  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$
  - ▶ noise free case  $\epsilon = \lim_{\sigma^2 \rightarrow 0} \sigma^2 \mathbf{I}$
- PCA is incredibly useful but its important to know what you are assuming, the probabilistic formulation allows you to do just that

---

<sup>8</sup>Bishop 2006, p. 12.2.4.

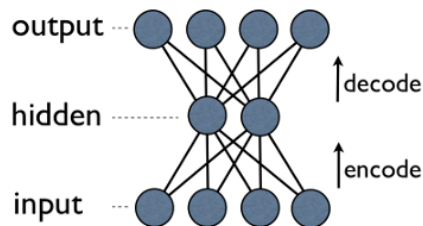
# Other Factor Analysis Models

- Independent Component Analysis
  - ▶  $p(\mathbf{X}) = \prod_i^q p(x_i)$
  - ▶ Cocktail party problem
- Auto-associative models
  - ▶  $p(\mathbf{X}|\mathbf{Y}) = \prod_i^N p(\mathbf{x}_i|\mathbf{y}_i)$
  - ▶ strange
- Lots and lots of different models differing by prior



# Other Factor Analysis Models

- Independent Component Analysis
  - ▶  $p(\mathbf{X}) = \prod_i^q p(x_i)$
  - ▶ Cocktail party problem
- Auto-associative models
  - ▶  $p(\mathbf{X}|\mathbf{Y}) = \prod_i^N p(\mathbf{x}_i|\mathbf{y}_i)$
  - ▶ strange
- Lots and lots of different models differing by prior





# Other Factor Analysis Models

- Independent Component Analysis

- ▶  $p(\mathbf{X}) = \prod_i^q p(x_i)$
- ▶ Cocktail party problem

- Auto-associative models

- ▶  $p(\mathbf{X}|\mathbf{Y}) = \prod_i^N p(\mathbf{x}_i|\mathbf{y}_i)$
- ▶ strange

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

- Lots and lots of different models differing by prior

## Assignment

You should now be able to do Task 2.3 and 2.4 in the assignment

# Gaussian Process Latent Variable Models

## History repeats itself

- In PPCA we assumed no uncertainty in the form of mapping
- We can use  $\mathcal{GP}$ s over mapping
- Gaussian Process Latent Variable Model [Lawrence 2005]

# Gaussian Process Latent Variable Models

## History repeats itself

- In PPCA we assumed no uncertainty in the form of mapping
- We can use  $\mathcal{GP}$ s over mapping
- Gaussian Process Latent Variable Model [Lawrence 2005]

# Gaussian Process Latent Variable Models

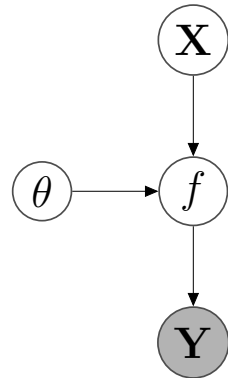
$$p(\mathbf{Y}|\mathbf{f}, \mathbf{X}, \theta) \tag{61}$$

- In PPCA we marginalised out  $\mathbf{X}$  and optimised for  $\mathbf{W}$
- Not possible for a general  $\mathcal{GP}$

# Gaussian Process Latent Variable Models

## GP-LVM

- General co-variance function (Ex. SE)
- $\mathbf{X}$  appears non-linearly in relation to  $\mathbf{Y}$
- Marginalisation of  $\mathbf{X}$  intractable



# Gaussian Process Latent Variable Models

$$\operatorname{argmax}_{\mathbf{X}, \theta} p(\mathbf{Y} | \mathbf{X}, \theta) p(\mathbf{X}) \quad (62)$$

$$p(\mathbf{Y} | \mathbf{X}, \theta) = \int p(\mathbf{Y} | \mathbf{f}) p(\mathbf{f} | \mathbf{X}, \theta) d\mathbf{f} \quad (63)$$

$$p(\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (64)$$

- **GP**-prior sufficiently regularises objective
- Need to set dimensionality of  $\mathbf{X}$

# Demo



# Bayesian PCA<sup>9</sup>

- PPCA We have no prior on  $\mathbf{W}$
- GP-LVM We have no prior on  $\mathbf{X}$
- Likelihood always increases with number of free parameters
  - ▶ larger dimensionality always better
- Cross-validation expensive
- Should learn distributions rather than point estimates



---

<sup>9</sup>Bishop 2006, p. 12.2.3

# Bayesian PCA<sup>9</sup>

- PPCA We have no prior on  $\mathbf{W}$
- GP-LVM We have no prior on  $\mathbf{X}$
- Likelihood always increases with number of free parameters
  - ▶ larger dimensionality always better
- Cross-validation expensive
- Should learn distributions rather than point estimates



---

<sup>9</sup>Bishop 2006, p. 12.2.3

# Bayesian PCA<sup>9</sup>

- PPCA We have no prior on  $\mathbf{W}$
- GP-LVM We have no prior on  $\mathbf{X}$
- Likelihood always increases with number of free parameters
  - ▶ larger dimensionality always better
- Cross-validation expensive
- Should learn distributions rather than point estimates



---

<sup>9</sup>Bishop 2006, p. 12.2.3

# Bayesian PCA<sup>9</sup>

- PPCA We have no prior on  $\mathbf{W}$
- GP-LVM We have no prior on  $\mathbf{X}$
- Likelihood always increases with number of free parameters
  - ▶ larger dimensionality always better
- Cross-validation expensive
- Should learn distributions rather than point estimates



---

<sup>9</sup>Bishop 2006, p. 12.2.3

# Bayesian PCA<sup>9</sup>

$$\begin{aligned} p(\mathbf{Y}) &= \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{X})p(\mathbf{W})d\mathbf{X}d\mathbf{W} \\ &= \int p(\mathbf{Y}|\mathbf{W})p(\mathbf{W})d\mathbf{W} = \int \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T)p(\mathbf{W})d\mathbf{W} \\ &\propto \int \exp\left(-\frac{1}{2}\text{tr}\left((\mathbf{Y} - \boldsymbol{\mu})^T(\boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T)^{-1}(\mathbf{Y} - \boldsymbol{\mu})\right)\right)p(\mathbf{W})d\mathbf{W} \end{aligned} \tag{65}$$

---

<sup>9</sup>Bishop 2006, p. 12.2.3

# Bayesian PCA<sup>9</sup>

## Inference

- Bayesian inference maximise  $p(\mathbf{Y})$ 
  - ▶ requires integrating out  $\mathbf{W}$  and  $\mathbf{X}$
  - ▶ intractable in PCA (and most models)
- Elementary functions **not** closed under integration

$$\int e^{e^{ax}} dx \quad (66)$$

- Solution through approximate inference (tomorrow & Friday & Hedvig)

---

<sup>9</sup>Bishop 2006, p. 12.2.3

## Summary

- Data often represented based on what we can measure
- Implicit representation & degrees of freedom of data
- Generative models
- Priors as preference
- Probabilistic PCA
- Bayesian inference

# Next Time

## Lecture 4

- November 12th 13-15 V1
- Summary of my part of the course
- Approximative Inference
  - ▶ Variational Bayes
- Complete assignment Task 2.3 and 2.4





# Next Time


## Lecture 4


- November 12th 13-15 V1
- Summary of my part of the course
- Approximative Inference
  - ▶ Variational Bayes
- Complete assignment Task 2.3 and 2.4



**e.o.f.**

# References I

 **Christopher M Bishop.** *Pattern recognition and machine learning*. 2006. URL: <http://www.library.wisc.edu/selectedtocs/bg0137.pdf>.

 **Neil D Lawrence.** “Probabilistic non-linear principal component analysis with Gaussian process latent variable models”. In: *The Journal of Machine Learning Research* 6 (2005), pp. 1783–1816. URL: <http://dl.acm.org/citation.cfm?id=1194904>.