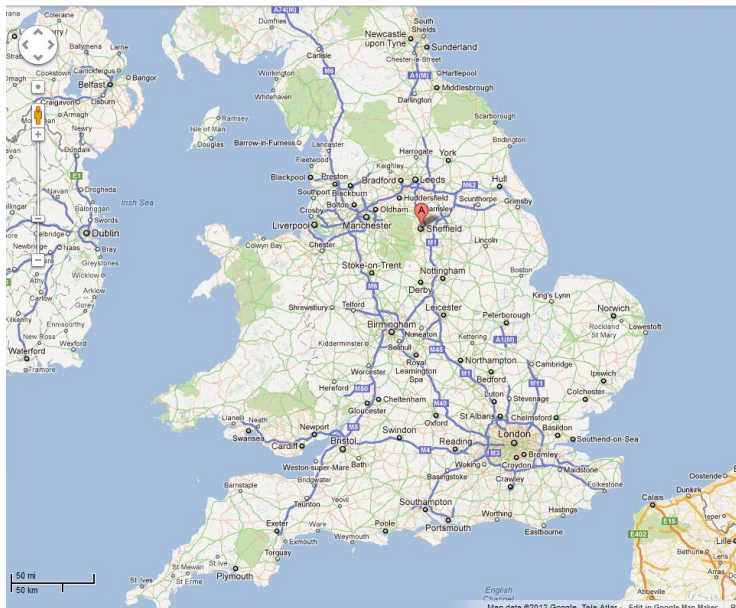# Probabilistic Models for Learning Data Representations

## Andreas Damianou

Department of Computer Science, University of Sheffield, UK

# Sheffield

# Outline

# Outline

# Probabilistic Models

"Probabilistic modelling involves the determination of a statistical model, a method for fitting that model to observed data, and a method for using the fitted model to solve the task at hand."

*D. Blei, D. Mimno*

# Treating Data as Random Variables



|   | 1 | 2 |   | 10000 |
|---|---|---|---|---|
| 1 | 0.81472 | 0.27603 | · · · | 0.58225 |
| 2 | 0.90579 | 0.6797 |   | 0.54074 |
| 3 | 0.12699 | 0.6551 |   | 0.86994 |
| 4 | 0.91338 | 0.16261 |   | 0.26478 |
| 5 | 0.63236 | 0.119 |   | 0.31807 |
|   | ⋮ |   |   | ⋮ |
| 50 | 0.75469 | 0.33712 | · · · | 0.64555 |

# Treating Data as Random Variables

# Treating Data as Random Variables



| | 1 | 2 | | 10000 |
|---|---|---|---|---|
| 1 | 0.81472 | 0.27603 | | 0.58225 |
| 2 | 0.90579 | 0.6797 | | 0.54074 |
| 3 | 0.12699 | 0.6551 | | 0.86994 |
| 4 | 0.91338 | 0.16261 | | 0.26478 |
| 5 | 0.63236 | 0.119 | | 0.31807 |
| | ⋮ | | | ⋮ |
| 50 | 0.75469 | 0.33712 | | 0.64555 |

$$\Downarrow$$
$$\mathbf{Y}$$

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

# Gaussian distribution

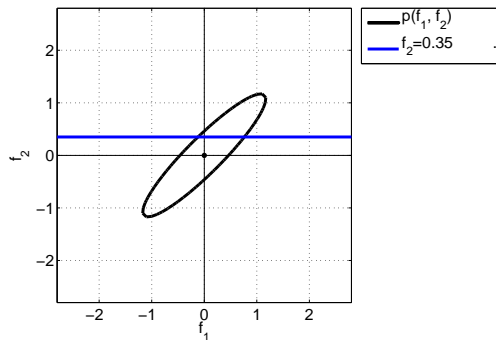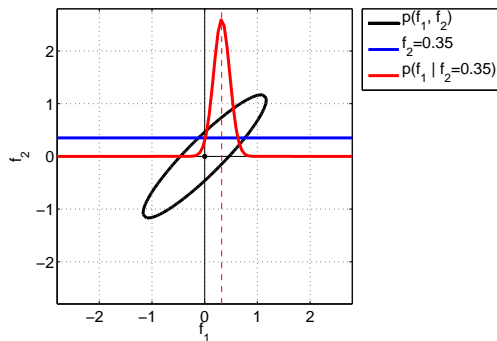Probability model: $p(f_1, f_2) \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$



Covariance between $f_1$ and $f_2$:

$$\mathbf{K} = \begin{bmatrix} 1 & 0.92 \\ 0.92 & 1 \end{bmatrix}$$

# Gaussian distribution

Probability model: $p(f_1, f_2) \sim \mathcal{N}\left(\mathbf{0}, \mathbf{K}\right)$

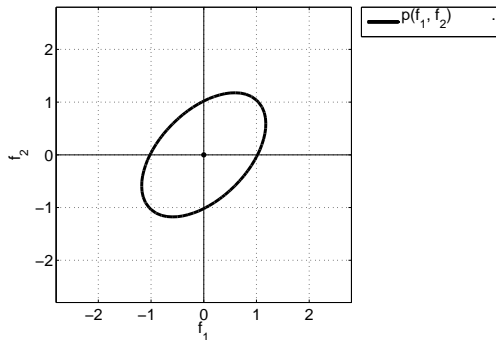

Covariance between $f_1$ and $f_2$:

$$\mathbf{K} = \begin{bmatrix} 1 & 0.92 \\ 0.92 & 1 \end{bmatrix}$$

# Gaussian distribution

Probability model: $p(f_1, f_2) \sim \mathcal{N}\left(\mathbf{0}, \mathbf{K}\right)$
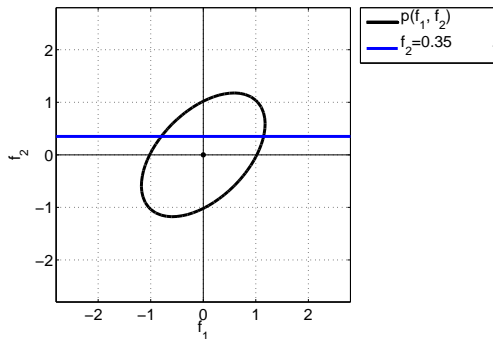


Covariance between $f_1$ and $f_2$:

$$\mathbf{K} = \begin{bmatrix} 1 & 0.92 \\ 0.92 & 1 \end{bmatrix}$$

# Gaussian distribution

Probability model: $p(f_1, f_2) \sim \mathcal{N}\left(\mathbf{0}, \mathbf{K}\right)$

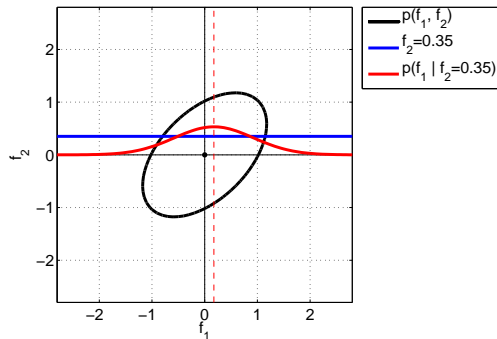

Covariance between $f_1$ and $f_2$:

$$\mathbf{K} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

# Gaussian distribution

Probability model: $p(f_1, f_2) \sim \mathcal{N}\left(\mathbf{0}, \mathbf{K}\right)$



Covariance between $f_1$ and $f_2$:

$$\mathbf{K} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

# Gaussian distribution

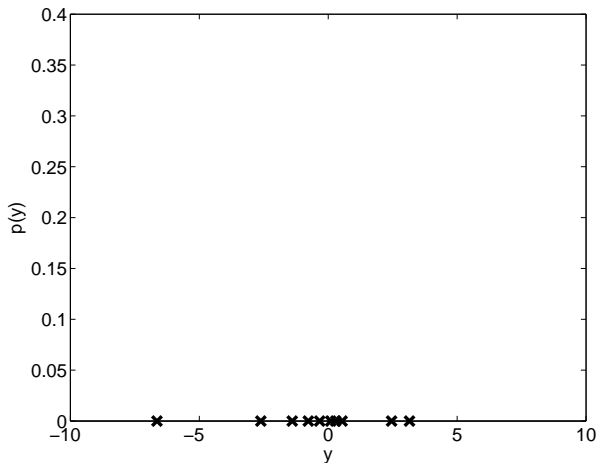Probability model: $p(f_1, f_2) \sim \mathcal{N}\left(\mathbf{0}, \mathbf{K}\right)$



Covariance between $f_1$ and $f_2$:

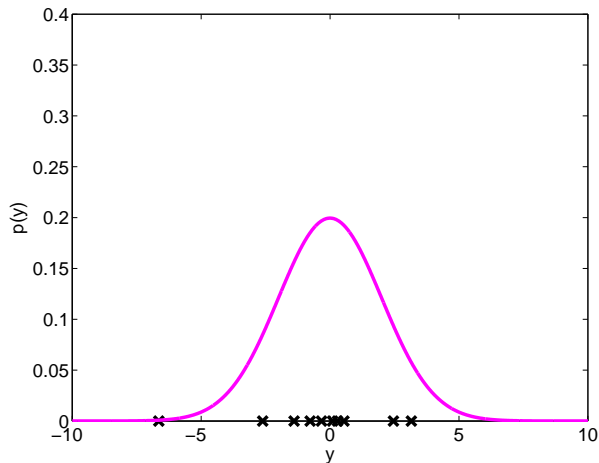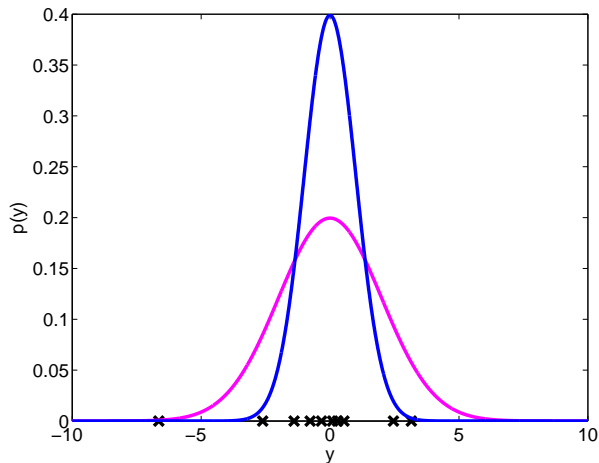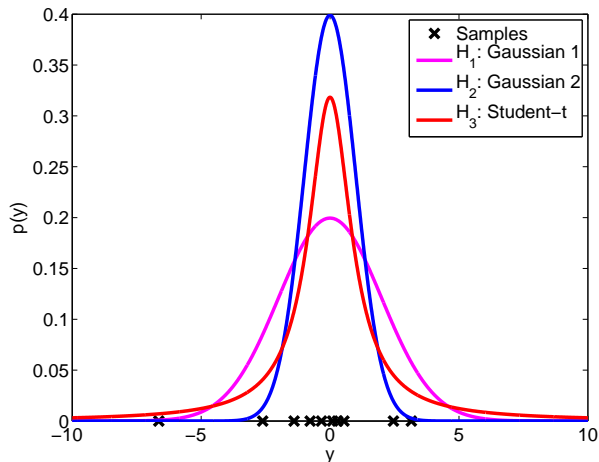$$\mathbf{K} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

# Model fitting

Which distribution (Hypothesis, $\mathcal{H}$) best *explains/fits* the data?

# Model fitting

Which distribution (Hypothesis, $\mathcal{H}$) best *explains/fits* the data?
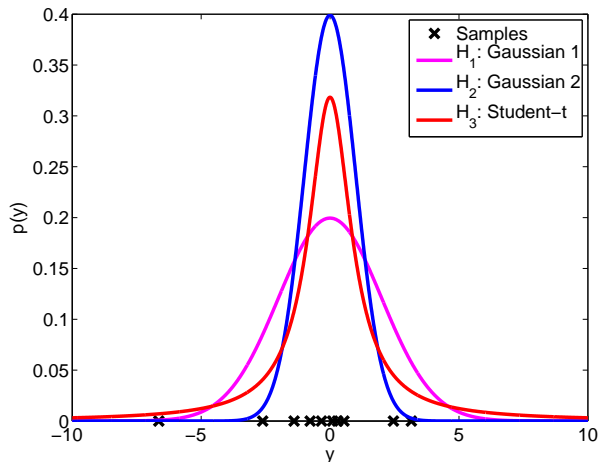
# Model fitting

Which distribution (Hypothesis, $\mathcal{H}$) best *explains/fits* the data?

# Model fitting

Which distribution (Hypothesis, $\mathcal{H}$) best *explains/fits* the data?

# Model fitting

Which distribution (Hypothesis, $\mathcal{H}$) best *explains/fits* the data?



Model fitting can be done with *maximum likelihood*.

# Bayes' rule

Taking things one step further: assume a model (hypothesis) $\mathcal{H}$ and a distribution for its parameters, $\theta$.

► Assume a prior distribution for our parameters, $\theta$.

► Assume a likelihood for the observed data, $y$, *given* the parameters.

► Find the posterior of the parameters, given the data.

► The normaliser of the posterior is the model evidence.

► All linked through *Bayes' rule*:

$$p(\theta|y, \mathcal{H}) = \frac{p(y|\theta, \mathcal{H})p(\theta|\mathcal{H})}{p(y|\mathcal{H}) = \int_\theta p(y|\theta, \mathcal{H})}$$

# Bayes' rule

Taking things one step further: assume a model (hypothesis) $\mathcal{H}$ and a distribution for its parameters, $\theta$.

- Assume a prior distribution for our parameters, $\theta$.
- Assume a likelihood for the observed data, $y$, *given* the parameters.
- Find the posterior of the parameters, given the data.
- The normaliser of the posterior is the model evidence.
- All linked through *Bayes' rule*:

$$p(\theta|y, \mathcal{H}) = \frac{p(y|\theta, \mathcal{H})p(\theta|\mathcal{H})}{p(y|\mathcal{H}) = \int_\theta p(y|\theta, \mathcal{H})}$$

# Bayes' rule

Taking things one step further: assume a model (hypothesis) $\mathcal{H}$ and a distribution for its parameters, $\theta$.

- Assume a prior distribution for our parameters, $\theta$.
- Assume a likelihood for the observed data, $y$, *given* the parameters.
- Find the posterior of the parameters, given the data.
- The normaliser of the posterior is the model evidence.
- All linked through *Bayes' rule*:

$$p(\theta|y, \mathcal{H}) = \frac{p(y|\theta, \mathcal{H})p(\theta|\mathcal{H})}{p(y|\mathcal{H}) = \int_\theta p(y|\theta, \mathcal{H})}$$

# Bayes' rule

Taking things one step further: assume a model (hypothesis) $\mathcal{H}$ and a distribution for its parameters, $\theta$.

- Assume a prior distribution for our parameters, $\theta$.
- Assume a likelihood for the observed data, $y$, *given* the parameters.
- Find the posterior of the parameters, given the data.
- The normaliser of the posterior is the model evidence.
- All linked through *Bayes' rule*:

$$p(\theta|y, \mathcal{H}) = \frac{p(y|\theta, \mathcal{H})p(\theta|\mathcal{H})}{p(y|\mathcal{H}) = \int_\theta p(y|\theta, \mathcal{H})}$$

# Bayes' rule

Taking things one step further: assume a model (hypothesis) $\mathcal{H}$ and a distribution for its parameters, $\theta$.

- Assume a prior distribution for our parameters, $\theta$.
- Assume a likelihood for the observed data, $y$, *given* the parameters.
- Find the posterior of the parameters, given the data.
- The normaliser of the posterior is the model evidence.
- All linked through *Bayes' rule*:

$$p(\theta|y, \mathcal{H}) = \frac{p(y|\theta, \mathcal{H})p(\theta|\mathcal{H})}{p(y|\mathcal{H}) = \int_{\theta} p(y|\theta, \mathcal{H})}$$

# Occam's razor

"Everything should be made as simple as possible, but not simpler". *A. Einstein*



Fig. 1. This figure is reproduced with permission from MacKay (1991). It has also appeared in MacKay (1992) and MacKay (2003, chapter 28). The Y-axis indexes all possible data sets (under some idealized ordering). Each curve gives a probability distribution over data sets, so must enclose an area of 1. $H_1$ is a simple model focusing on data in region $C_1$. Given data is this region, $H_1$ has more evidence than a more powerful model $H_2$, which would be favored given more complex data (outside $C_1$). [Murray and Ghahramani, 2001]

# Latent Variables

- What are the *latent* features of "cuteness"?

# Another example: latent *process*

Is Beckham an expert in Newtonian & trajectory mechanics?

Is Beckham an expert in Newtonian & trajectory mechanics?



$$m\frac{\mathrm{d}^2\vec{x}(t)}{\mathrm{d}t^2} = -\nabla V(\vec{x}(t)), \quad \vec{x} = (x, y, z)$$

$$R_s = \sqrt{x^2 + y^2}$$

$$= \sqrt{\left(\frac{2v^2\cos^2\theta}{g}\left(\frac{\sin\theta}{\cos\theta} - m\right)\right)^2 + \left(m\frac{2v^2\cos^2\theta}{g}\left(\frac{\sin\theta}{\cos\theta} - m\right)\right)^2}$$

# Outline

# Introducing Gaussian Processes:

- A Gaussian distribution depends on a mean and a covariance matrix.
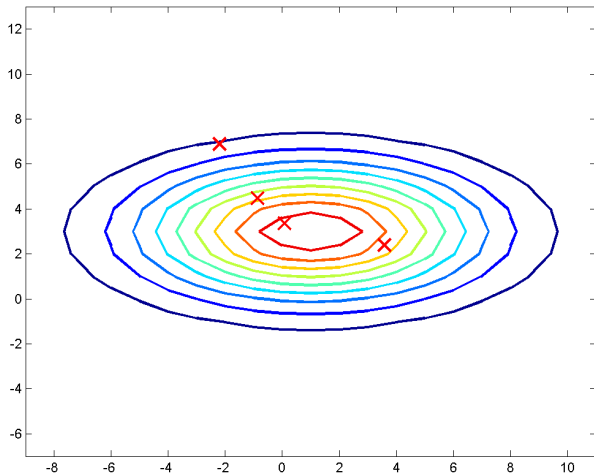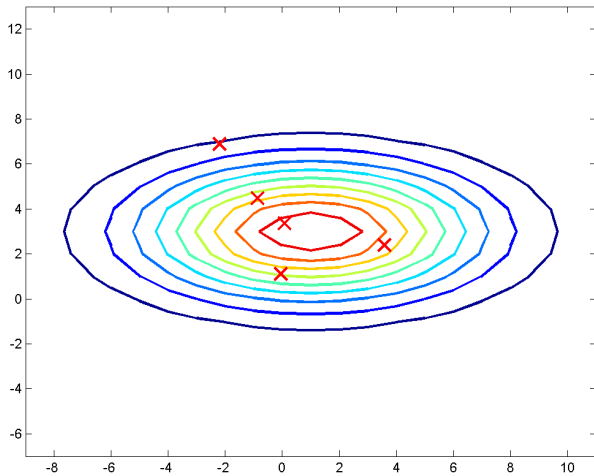- A Gaussian process depends on a mean and a covariance function.

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian
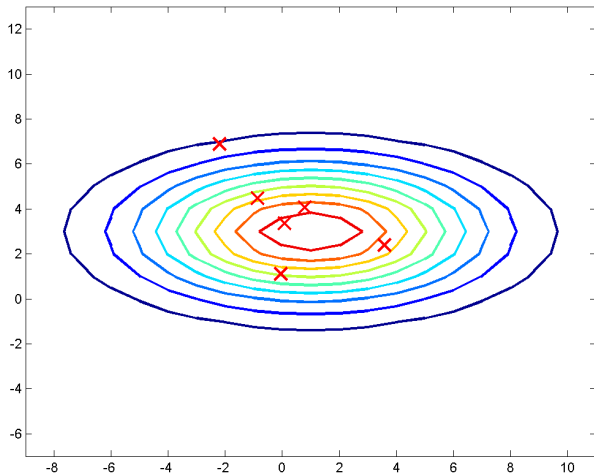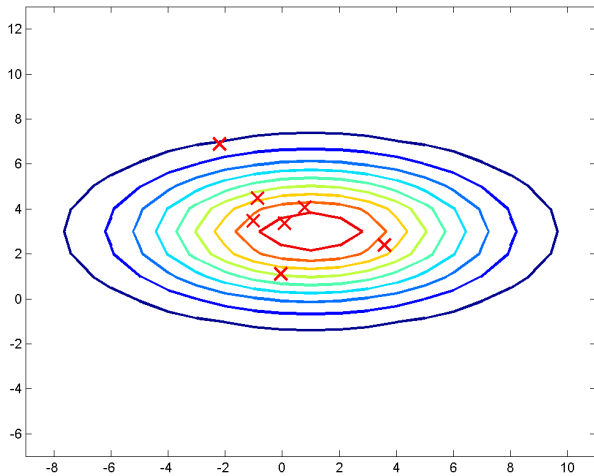
Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

Sampling from a 1-D Gaussian

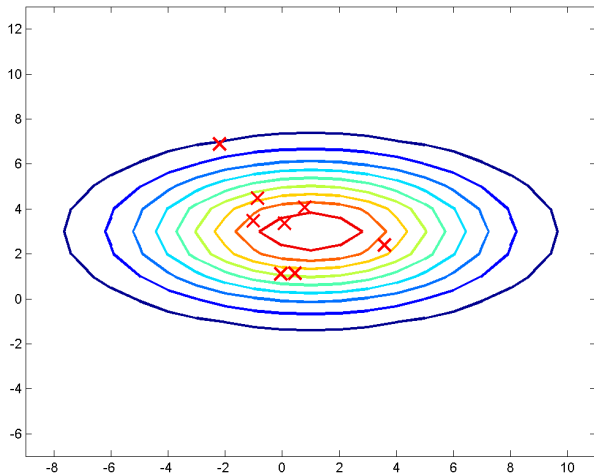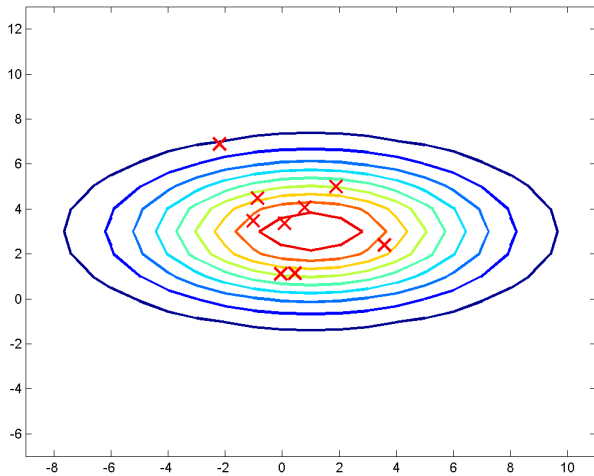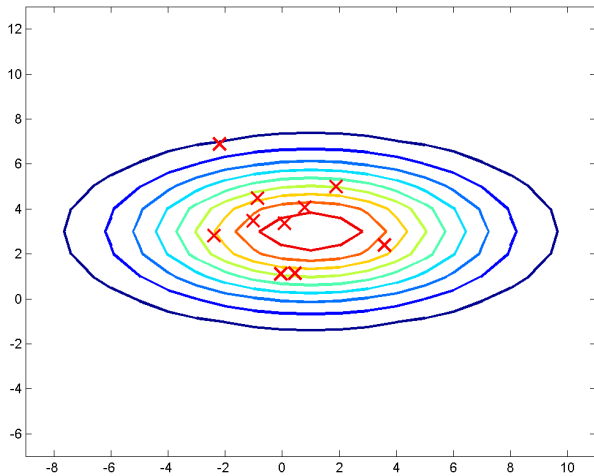Sampling from a 2-D Gaussian

Sampling from a 2-D Gaussian

Sampling from a 2-D Gaussian

Sampling from a 2-D Gaussian

Sampling from a 2-D Gaussian

Sampling from a 2-D Gaussian

Sampling from a 2-D Gaussian
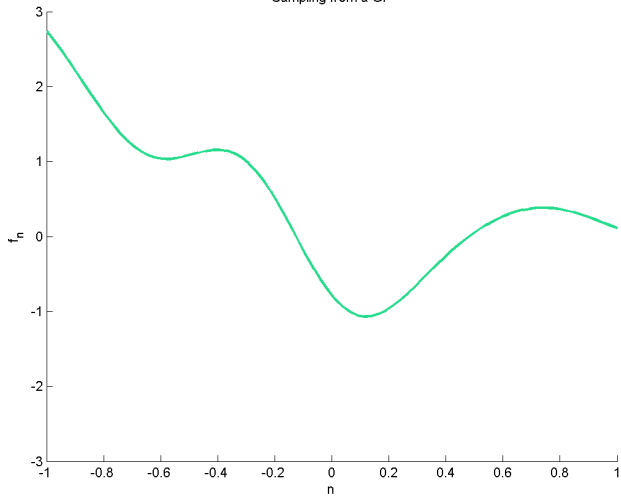
Sampling from a 2-D Gaussian
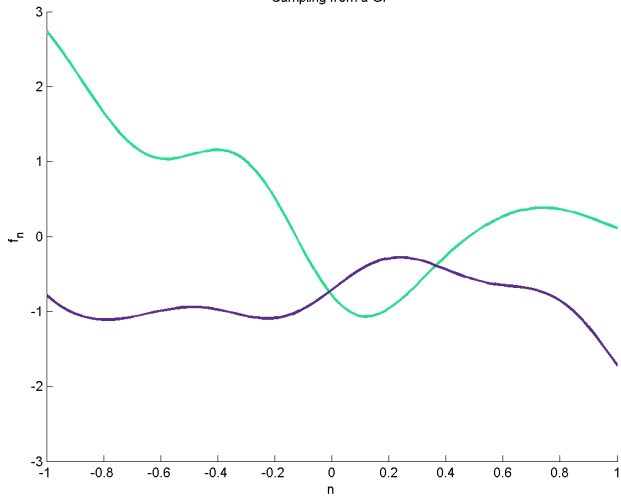
Sampling from a 2-D Gaussian
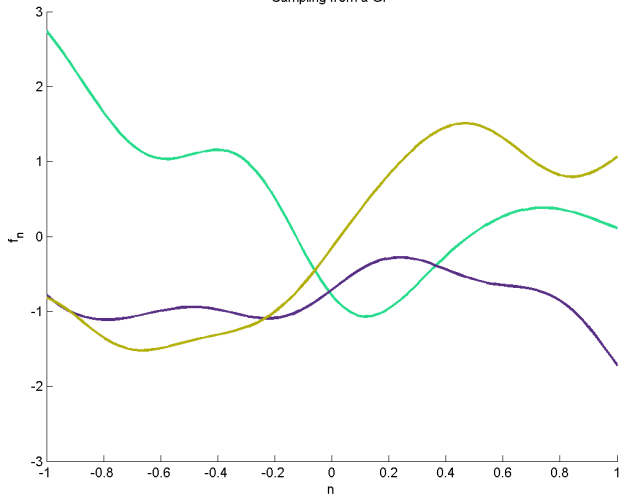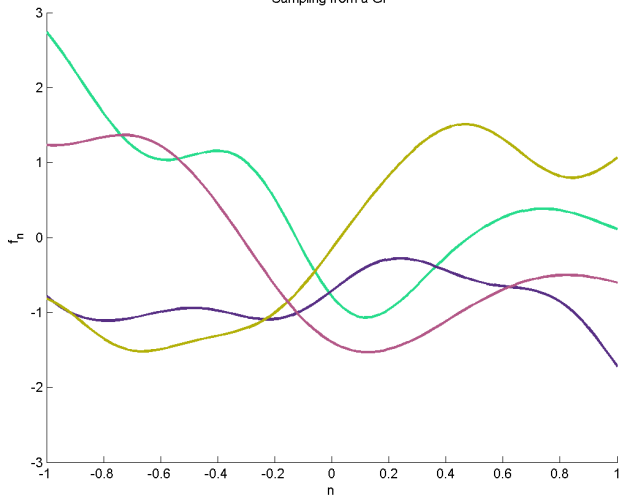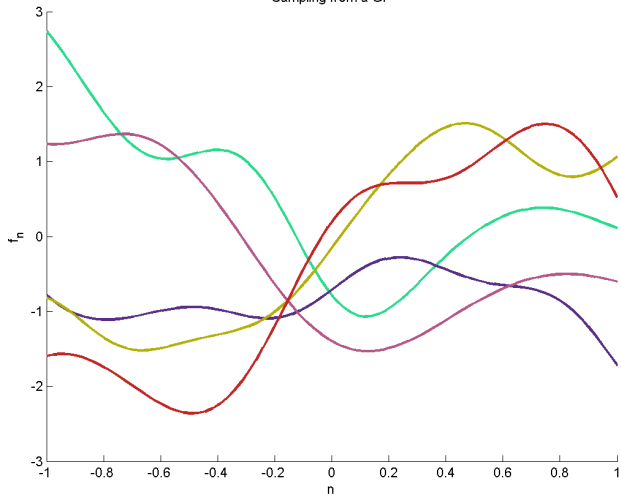
Sampling from a 2-D Gaussian

Sampling from a GP

Sampling from a GP

Sampling from a GP

Sampling from a GP

Sampling from a GP

Sampling from a GP

Sampling from a GP

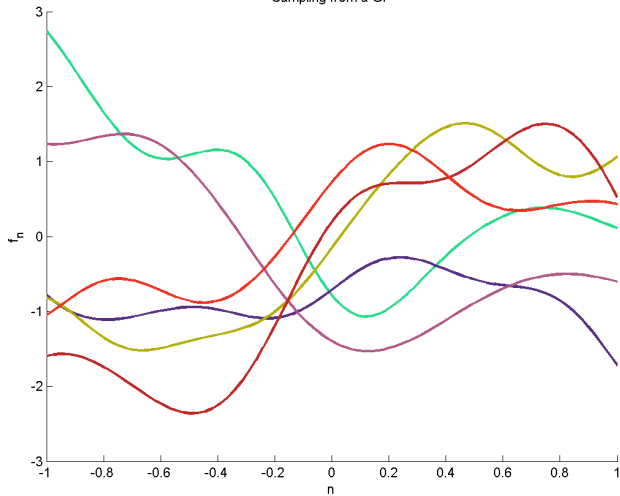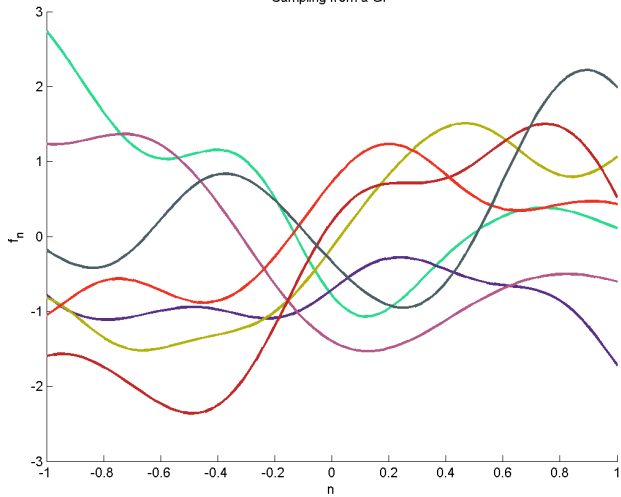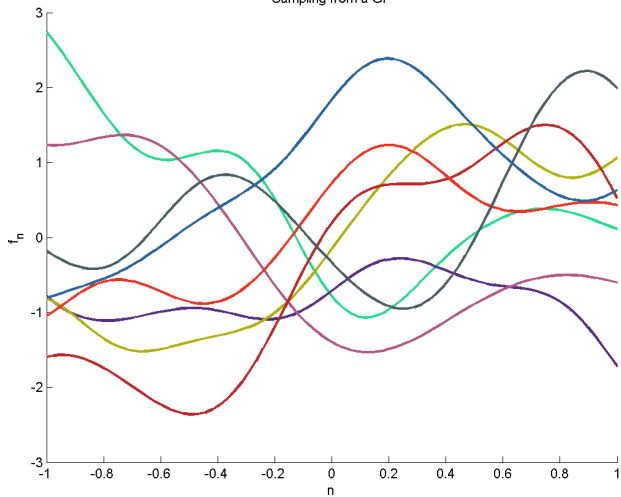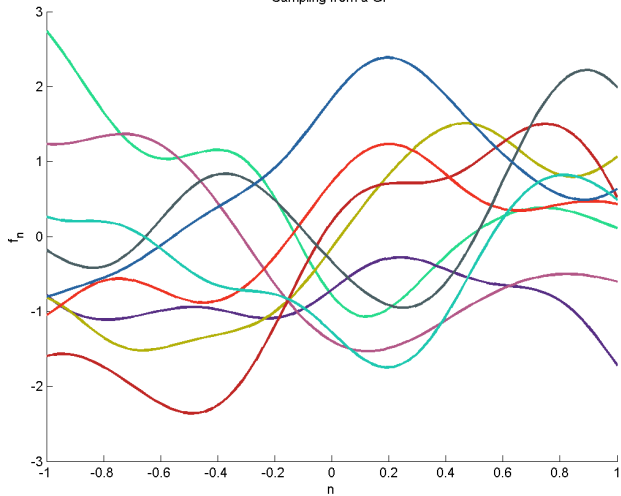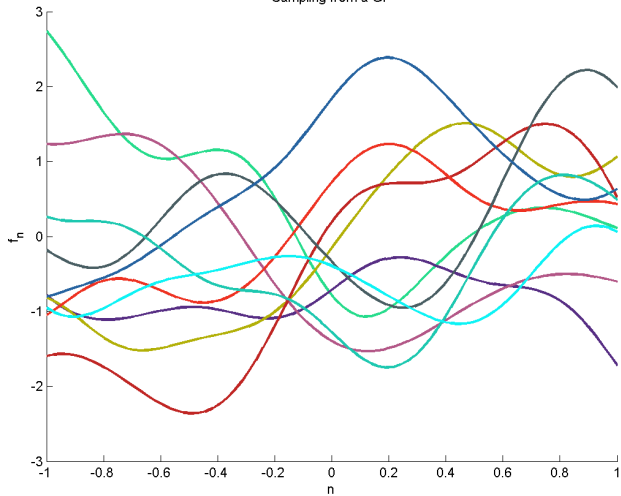Sampling from a GP

Sampling from a GP

Sampling from a GP

Let's start with a multivariate Gaussian:

$$p(\underbrace{f_1, f_2, \cdots, f_s}_{\mathbf{f}_A}, \underbrace{f_{s+1}, f_{s+2}, \cdots, f_N}_{\mathbf{f}_B}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$

with:

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix} \text{ and } \mathbf{K} = \begin{bmatrix} \mathbf{K}_{AA} & \mathbf{K}_{AB} \\ \mathbf{K}_{BA} & \mathbf{K}_{BB} \end{bmatrix}$$

Marginalisation property:

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$

$$p(\mathbf{f}_A) = \int_{\mathbf{f}_B} p(\mathbf{f}_A, \mathbf{f}_B) d\mathbf{f}_B = \mathcal{N}(\boldsymbol{\mu}_A, \mathbf{K}_{AA})$$

Let's start with a multivariate Gaussian:

$$p(\underbrace{f_1, f_2, \cdots, f_s}_{\mathbf{f}_A}, \underbrace{f_{s+1}, f_{s+2}, \cdots, f_N}_{\mathbf{f}_B}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$

with:

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix} \text{ and } \mathbf{K} = \begin{bmatrix} \mathbf{K}_{AA} & \mathbf{K}_{AB} \\ \mathbf{K}_{BA} & \mathbf{K}_{BB} \end{bmatrix}$$

Marginalisation property:

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$
$$p(\mathbf{f}_A) = \int_{\mathbf{f}_B} p(\mathbf{f}_A, \mathbf{f}_B) d\mathbf{f}_B = \mathcal{N}(\boldsymbol{\mu}_A, \mathbf{K}_{AA})$$

In the GP context:

$$\boldsymbol{\mu}_{\infty} = \begin{bmatrix} \boldsymbol{\mu}_{\mathbf{x}} \\ \cdots \\ \cdots \end{bmatrix} \quad \text{and} \quad \mathbf{K}_{\infty} = \begin{bmatrix} \mathbf{K}_{\mathbf{xx}} & \cdots \\ \cdots & \cdots \end{bmatrix}$$

## Posterior is also Gaussian!

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$
$$p(\mathbf{f}_A | \mathbf{f}_B) = \mathcal{N}(\cdots, \cdots)$$
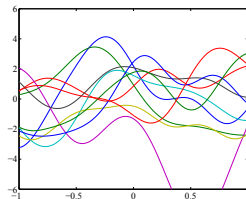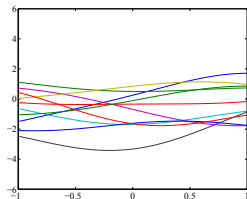
In the GP context this can be used for inter/extrapolation:

$$p(f_* | f_1, \cdots, f_N) = p(f(x_*) | f(x_1), \cdots, f(x_N)) \sim \mathcal{N}$$

But where is $\mathbf{K}_{\cdot\cdot}$ coming from in GPs?

# Posterior is also Gaussian!

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$
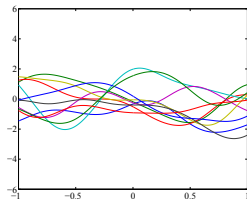$$p(\mathbf{f}_A | \mathbf{f}_B) = \mathcal{N}(\cdots, \cdots)$$

In the GP context this can be used for inter/extrapolation:

$$p(f_* | f_1, \cdots, f_N) = p(f(x_*) | f(x_1), \cdots, f(x_N)) \sim \mathcal{N}$$

But where is $\mathbf{K}_{..}$ coming from in GPs?

# Posterior is also Gaussian!

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$
$$p(\mathbf{f}_A | \mathbf{f}_B) = \mathcal{N}(\cdots, \cdots)$$

In the GP context this can be used for inter/extrapolation:

$$p(f_* | f_1, \cdots, f_N) = p(f(x_*) | f(x_1), \cdots, f(x_N)) \sim \mathcal{N}$$

But where is $\mathbf{K}_{\cdot\cdot}$ coming from in GPs?

# Covariance samples and hyperparameters

- $k(x, x') = \alpha \exp\left(-\frac{\gamma}{2}(x - x')^\top (x - x')\right)$
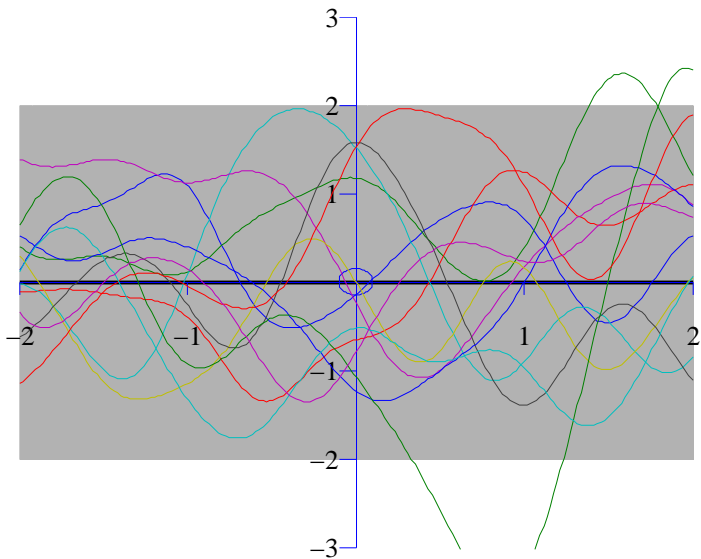- The hyperparameters of the cov. function define the properties (and NOT an explicit form) of the sampled functions

- So far we assumed: $\mathbf{f} = f(\mathbf{X})$
- Assuming that we only observe noisy versions $\mathbf{y}$ of the true outputs $\mathbf{f}$:

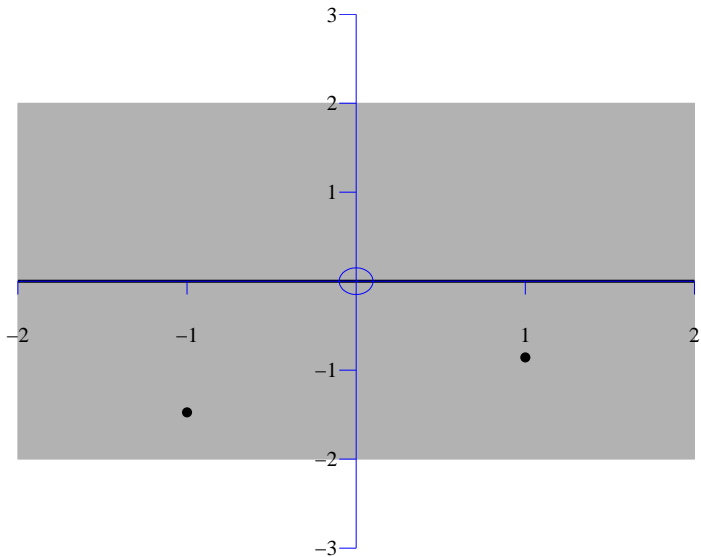$$\mathbf{y} = f(\mathbf{X}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

# Fitting the data

# Fitting the data

# Application to Disease modelling

Ricardo Andrade Pacheco.
http://ric70x7.github.io/research.html

▶ If $\mathbf{X}$ is unobserved, treat it as a parameter and optimize over it.
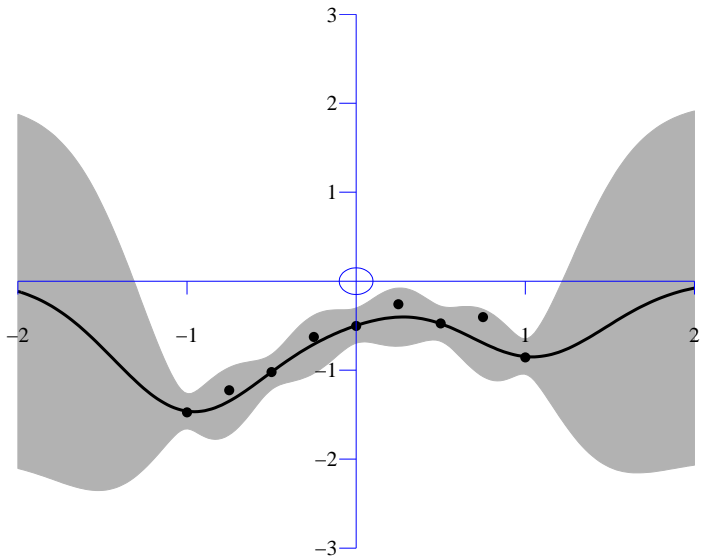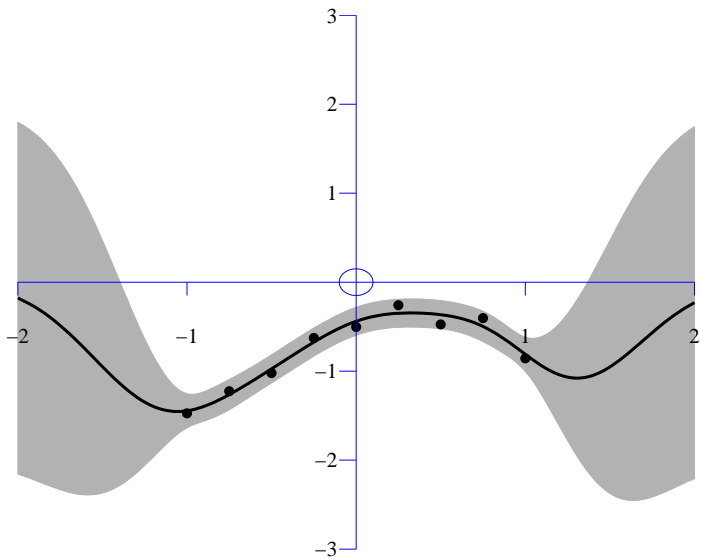
# Manifold Relevance Determination



- ▶ Observations come into two different *views*: $Y$ and $Z$.
- ▶ The latent space is segmented into parts private to $Y$, private to $Z$ and shared between $Y$ and $Z$.
- ▶ Used for data consolidation and discovering commonalities.

# Consolidating complementary experimental data



*ethanol*        *glucose*

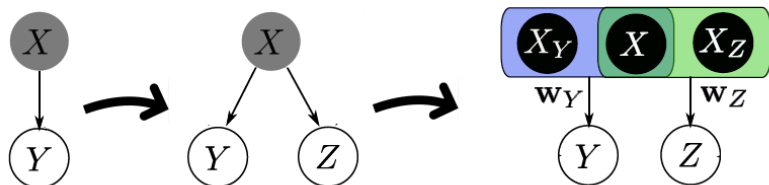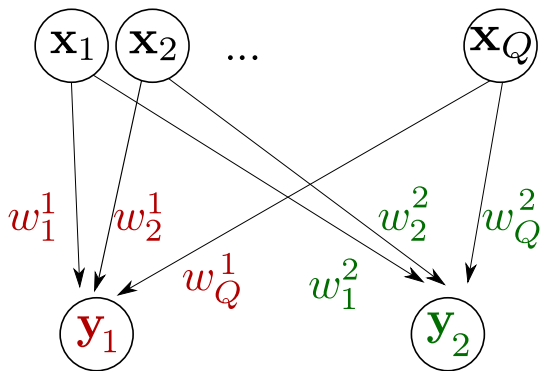SNPs (genotype) → Gene Expression (phenotype)    Gene Expression (phenotype)

*Shared information:* biological signal / confounders
*Private information:* environmental confounders

Confounders: Statistical relationships that do not reflect the true causality in the data

# Application to Health Modelling

Research agenda of Prof. Neil Lawrence's group:

# Example: faces

▸ https://youtu.be/rIPX3ClOhKY
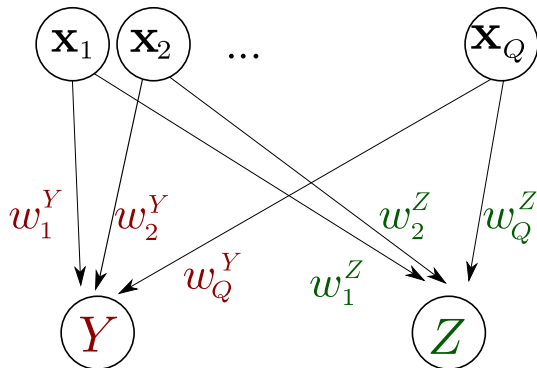
# Example: robotics

# Summary

- ..

# Thanks

Thanks to Neil Lawrence, James Hensman, Michalis Titsias, Carl Henrik Ek.

References:

- N. D. Lawrence (2006) "The Gaussian process latent variable model" Technical Report no CS-06-03, The University of Sheffield, Department of Computer Science

- N. D. Lawrence (2006) "Probabilistic dimensional reduction with the Gaussian process latent variable model" (talk)

- C. E. Rasmussen (2008), "Learning with Gaussian Processes", Max Planck Institute for Biological Cybernetics, Published: Feb. 5, 2008 (Videolectures.net)

- Carl Edward Rasmussen and Christopher K. I. Williams. Gaussian Processes for Machine Learning. MIT Press, Cambridge, MA, 2006. ISBN 026218253X.

- M. K. Titsias (2009), "Variational learning of inducing variables in sparse Gaussian processes", AISTATS 2009

- A. C. Damianou, M. K. Titsias and N. D. Lawrence (2011), "Variational Gaussian process dynamical systems", NIPS 2011

- A. C. Damianou, C. H. Ek, M. K. Titsias and N. D. Lawrence (2012), "Manifold Relevance Determination", ICML 2012

- A. C. Damianou and N. D. Lawrence (2013), "Deep Gaussian processes", AISTATS 2013

- J. Hensman (2013), "Gaussian processes for Big Data", UAI 2013

BACKUP SLIDES

# Dimensionality reduction: Linear vs non-linear



$\mathbf{y}_i = f(\mathbf{x}_i)$

*f is linear*

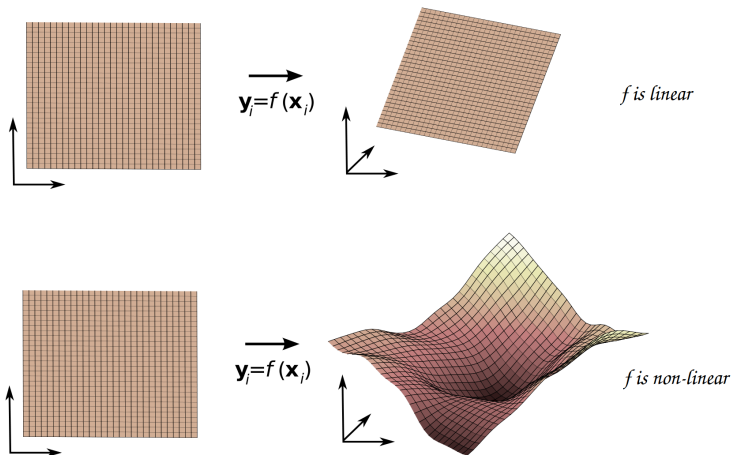$\mathbf{y}_i = f(\mathbf{x}_i)$

*f is non-linear*

*Image from: "Dimensionality Reduction the Probabilistic Way", N. Lawrence, ICML tutorial 2008*