

Modelling and consolidating complex data with Gaussian process models

Andreas Damianou

Department of Neuro- and Computer Science, University of
Sheffield, UK

ICS-FORTH, Heraklion, Greece, 10/06/2014

Outline

Part 1: Gaussian processes

 GPs for nonparametric, nonlinear regression

Introducing latent spaces: GP-LVM

Multiple views: MRD

Summary

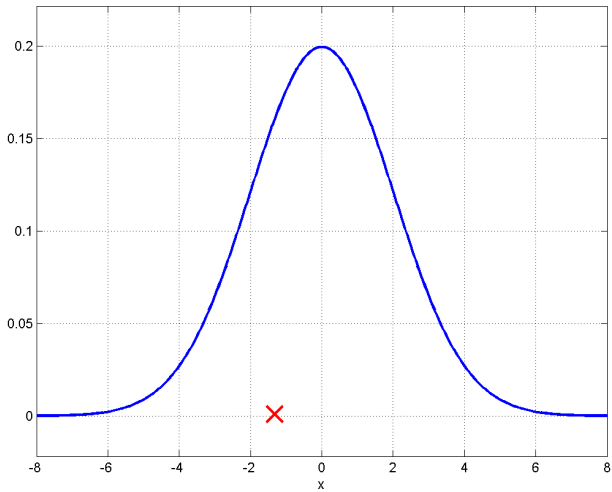
Part 2: Deep Gaussian processes

Introducing Gaussian Processes:

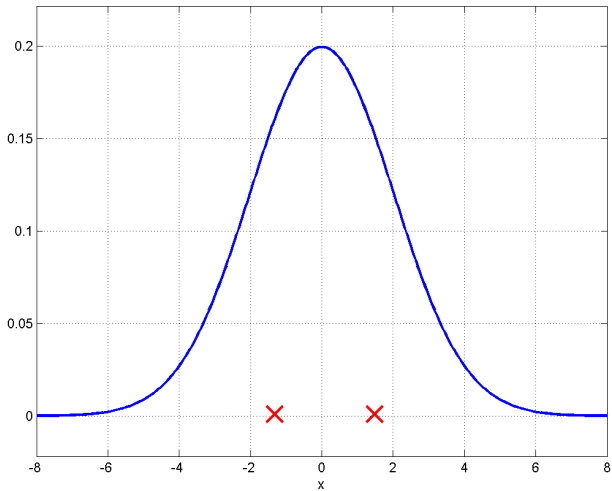
- A Gaussian **distribution** depends on a mean and a covariance vector / matrix.
- A Gaussian **process** depends on a mean and a covariance function.

Next: Demo, from Gaussian distributions to Gaussian processes.

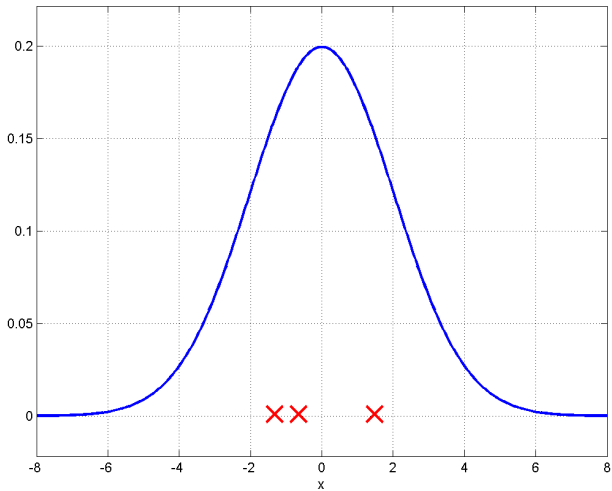
Sampling from a 1-D Gaussian



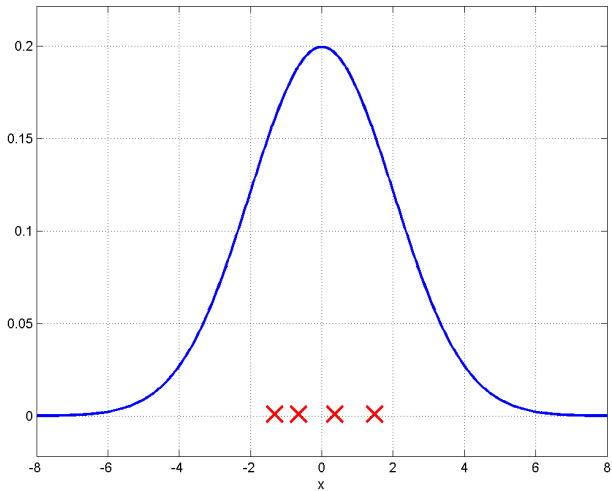
Sampling from a 1-D Gaussian



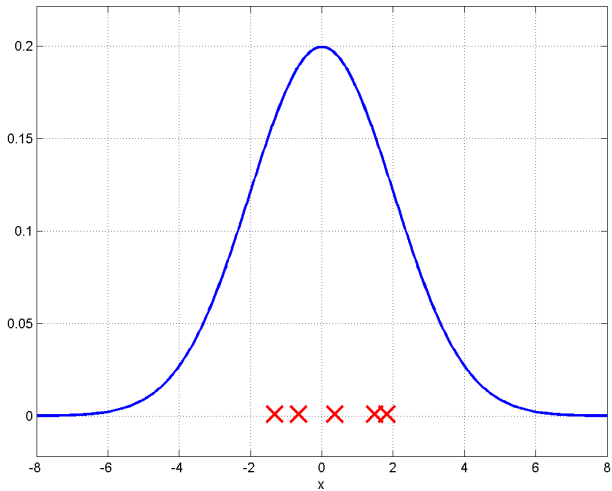
Sampling from a 1-D Gaussian



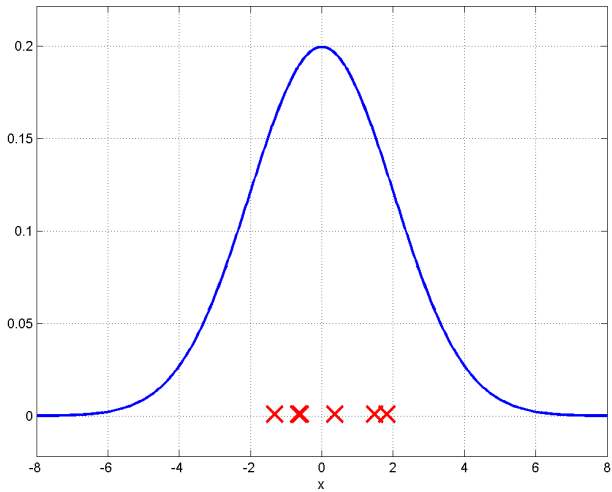
Sampling from a 1-D Gaussian



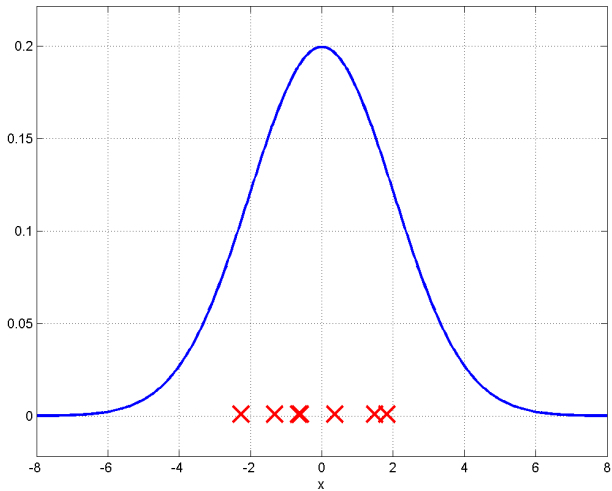
Sampling from a 1-D Gaussian



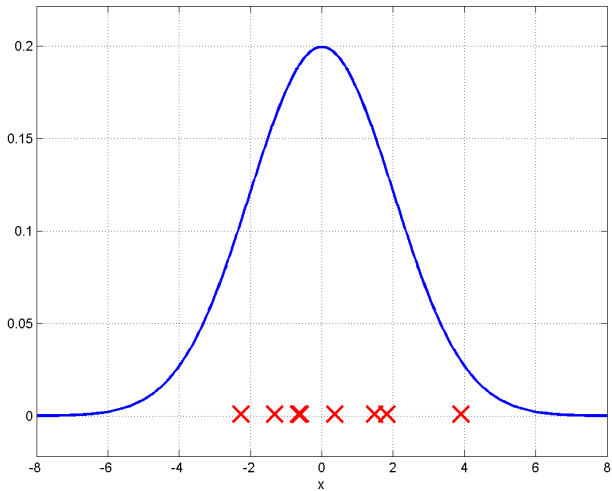
Sampling from a 1-D Gaussian



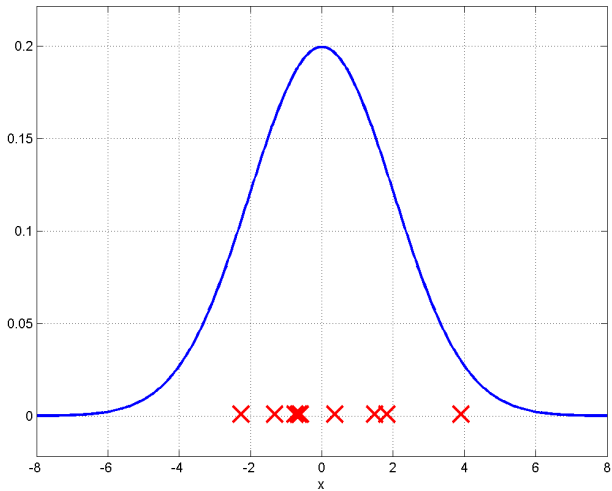
Sampling from a 1-D Gaussian



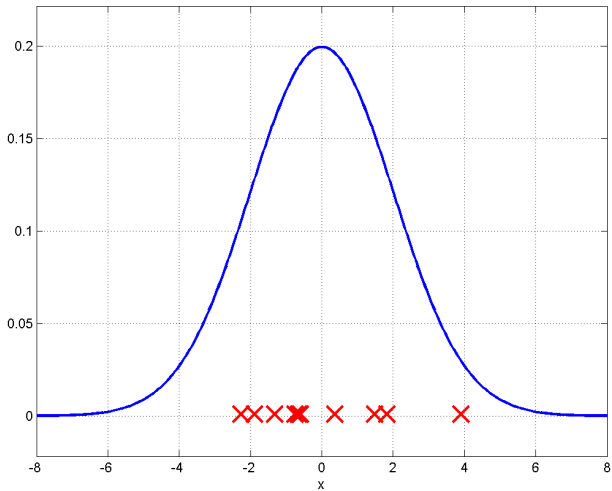
Sampling from a 1-D Gaussian



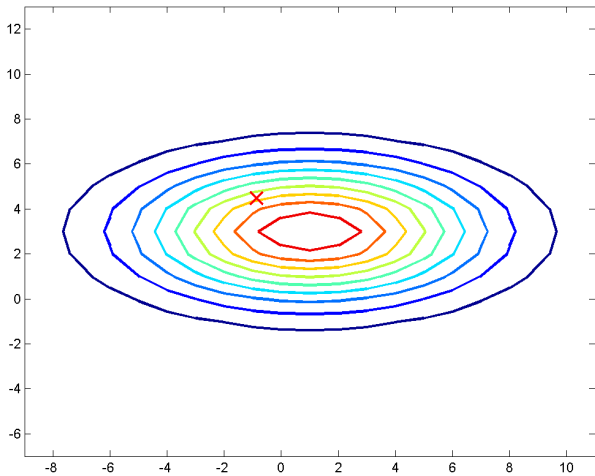
Sampling from a 1-D Gaussian



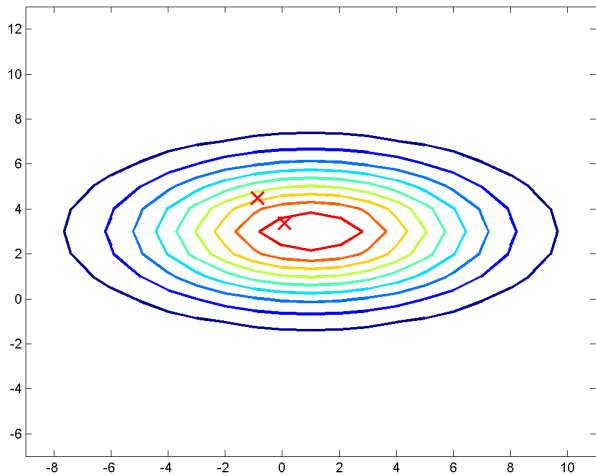
Sampling from a 1-D Gaussian



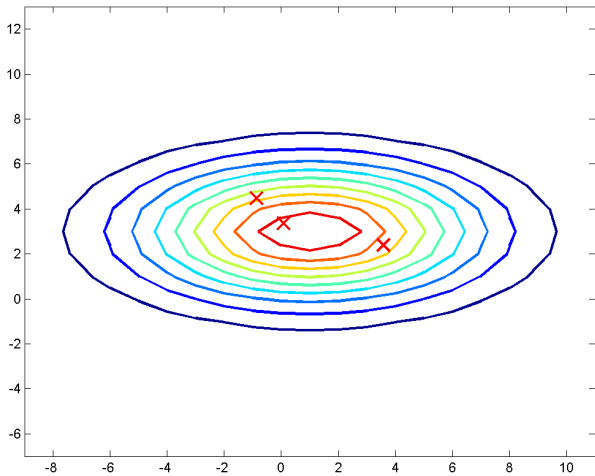
Sampling from a 2-D Gaussian



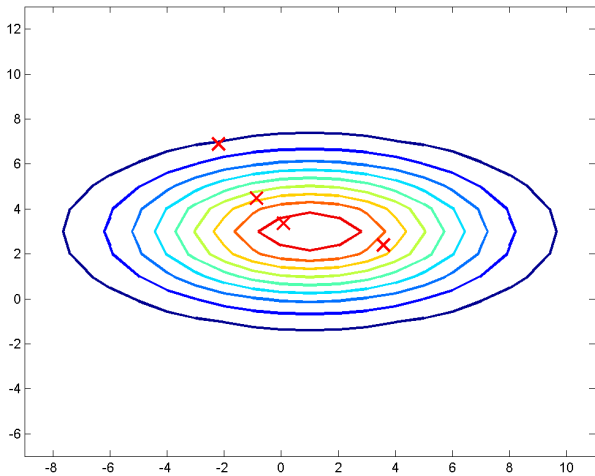
Sampling from a 2-D Gaussian



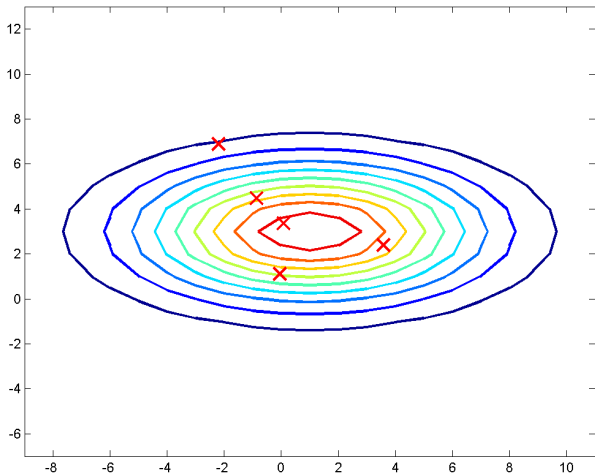
Sampling from a 2-D Gaussian



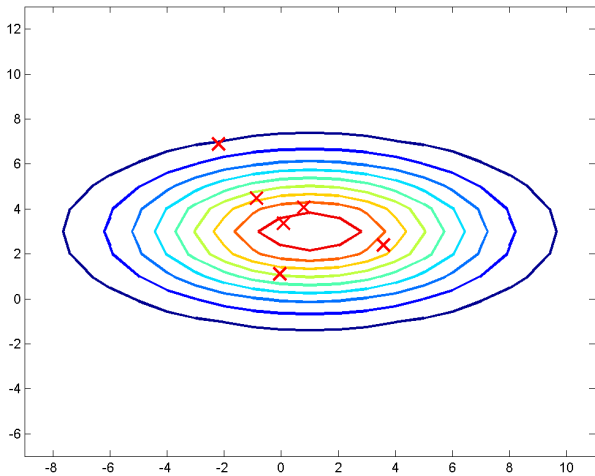
Sampling from a 2-D Gaussian



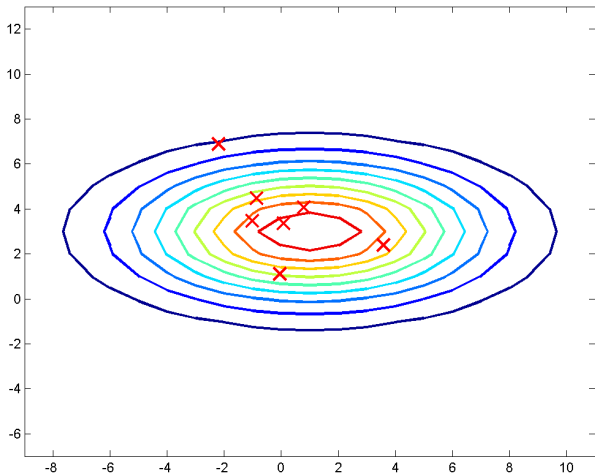
Sampling from a 2-D Gaussian



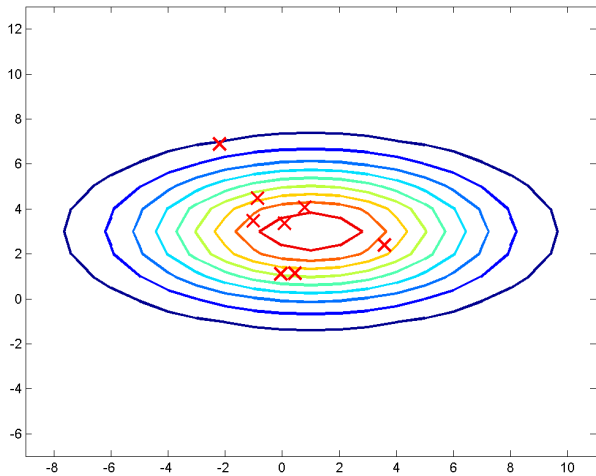
Sampling from a 2-D Gaussian



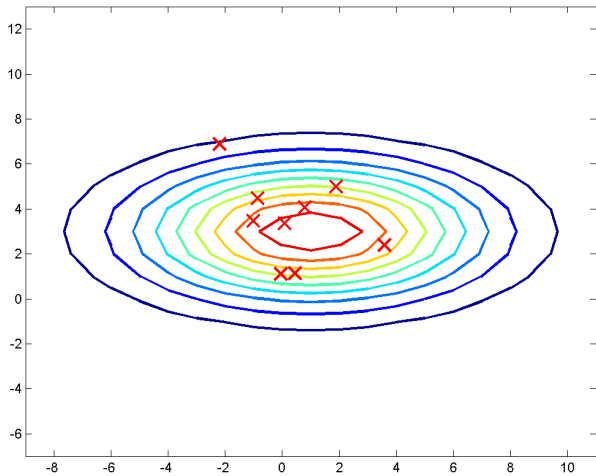
Sampling from a 2-D Gaussian



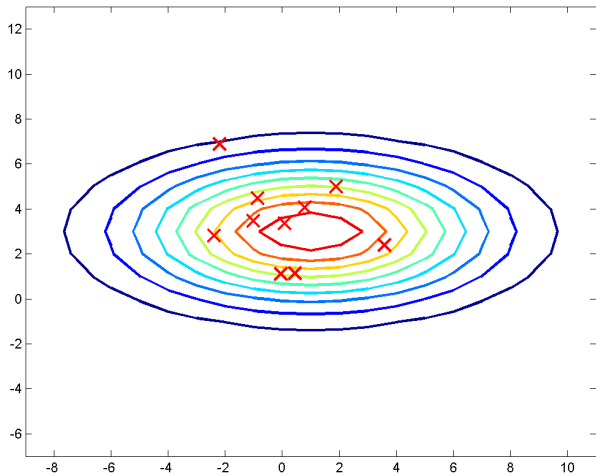
Sampling from a 2-D Gaussian



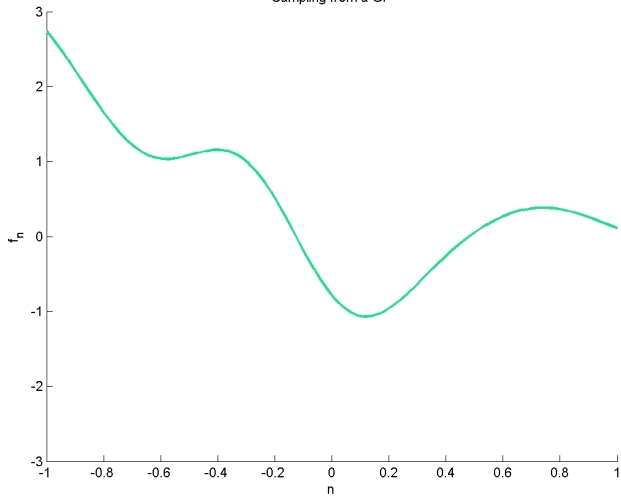
Sampling from a 2-D Gaussian



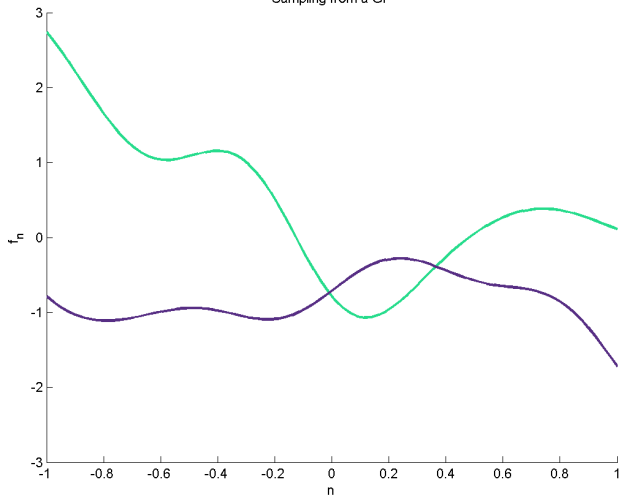
Sampling from a 2-D Gaussian



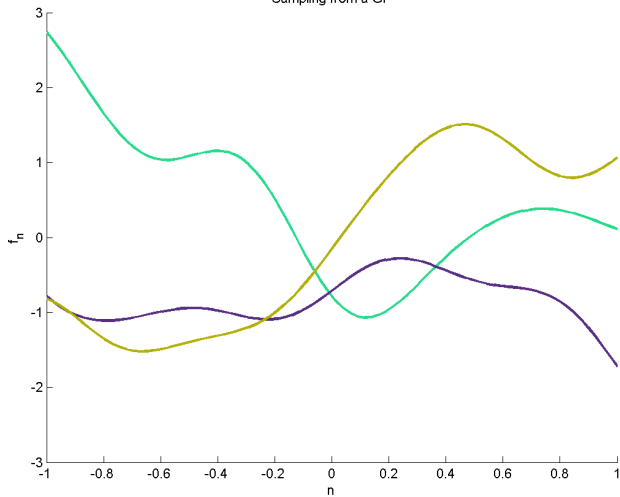
Sampling from a GP



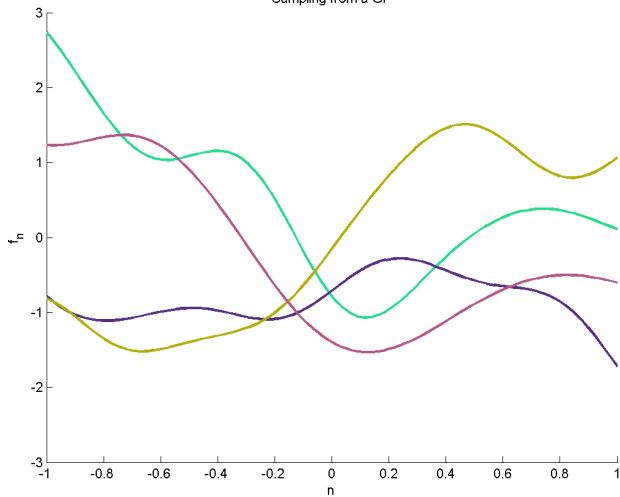
Sampling from a GP



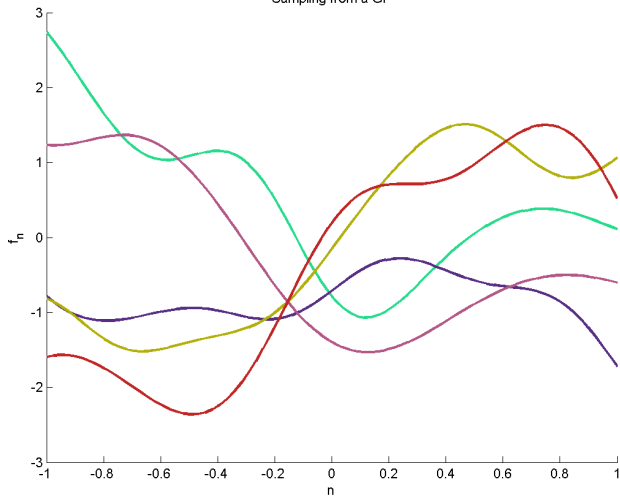
Sampling from a GP



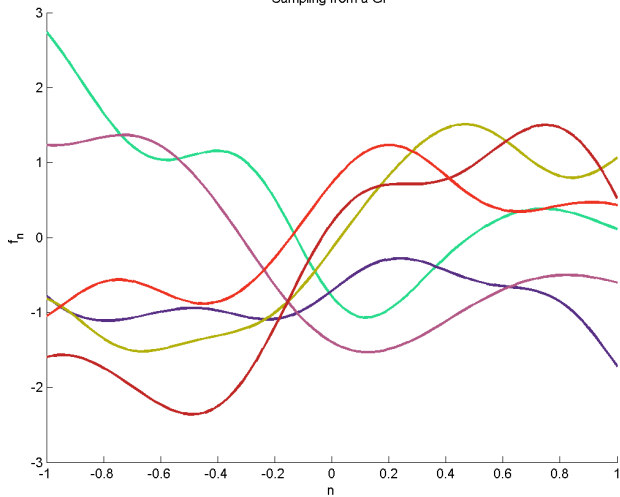
Sampling from a GP



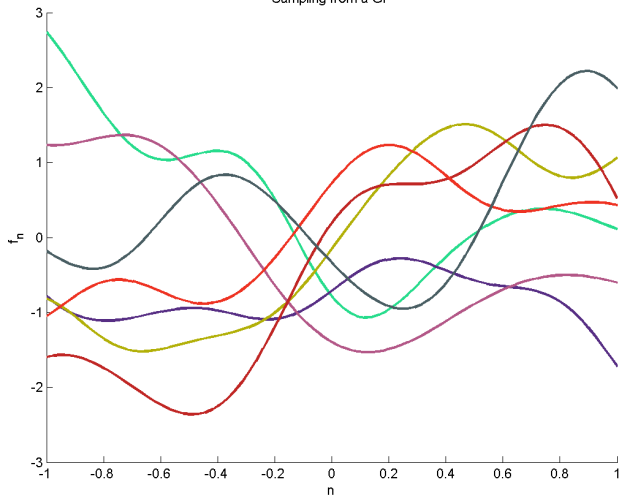
Sampling from a GP



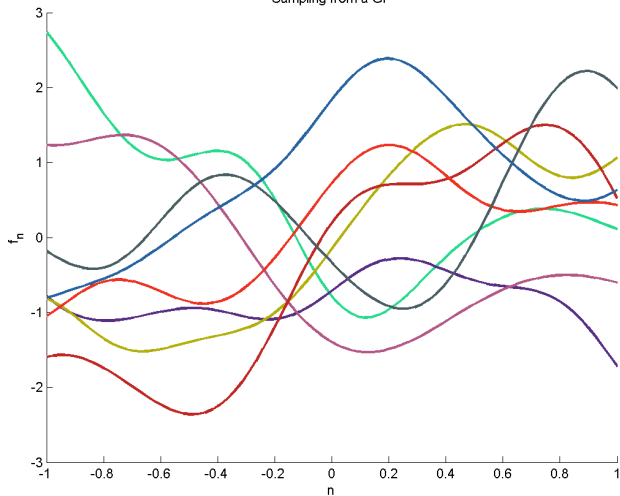
Sampling from a GP



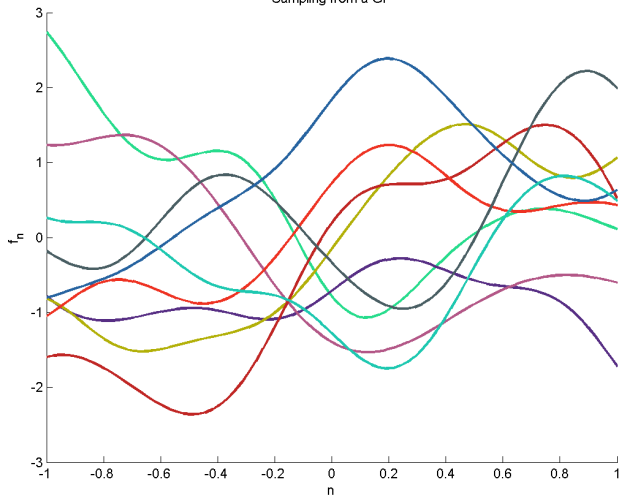
Sampling from a GP



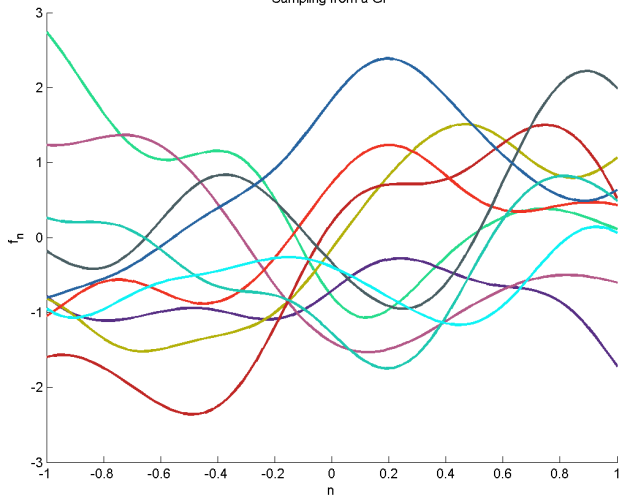
Sampling from a GP



Sampling from a GP



Sampling from a GP



Infinite model... but we *always* work with finite sets!

$$p(f_A, f_B) \sim \mathcal{N}(\mu, \mathbf{K}).$$

with:

$$\mu = \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix} \quad \text{and} \quad \mathbf{K} = \begin{bmatrix} \mathbf{K}_{AA} & \mathbf{K}_{AB} \\ \mathbf{K}_{BA} & \mathbf{K}_{BB} \end{bmatrix}$$

Infinite model... but we *always* work with finite sets!

$p(f_A, f_B) \sim \mathcal{N}(\mu, \mathbf{K})$. then:

$$p(f_A) = \int_{f_B} p(f_A, f_B) \mathrm{d}f_B = \mathcal{N}(\mu_A, \mathbf{K}_{AA})$$

$$p(f_B) = \int_{f_A} p(f_A, f_B) \mathrm{d}f_A = \mathcal{N}(\mu_B, \mathbf{K}_{BB})$$

Infinite model... but we *always* work with finite sets!

$p(f_A, f_B) \sim \mathcal{N}(\mu, \mathbf{K})$. then:

$$p(f_A) = \int_{f_B} p(f_A, f_B) \mathrm{d}f_B = \mathcal{N}(\mu_A, \mathbf{K}_{AA})$$

$$p(f_B) = \int_{f_A} p(f_A, f_B) \mathrm{d}f_A = \mathcal{N}(\mu_B, \mathbf{K}_{BB})$$

In the GP context:

$$\begin{aligned} p(\overbrace{f_1, f_2, \dots, f_N}^{\text{training data}}) &= p(f(x_1), f(x_2), \dots, f(x_N)) \\ &= \int_{\mathbb{R} - \{X\}} p(f(\{x_i \in \mathbb{R}\})) \\ &= \mathcal{N}(\mu_X, \mathbf{K}_{XX}) \end{aligned}$$

Posterior is also Gaussian!

$p(f_A, f_B) \sim \mathcal{N}(\mu, \mathbf{K})$. then:

$$p(f_A|f_B) = \mathcal{N}(\cdots, \cdots)$$

$$p(f_B|f_A) = \mathcal{N}(\cdots, \cdots)$$

Posterior is also Gaussian!

$$p(f_A, f_B) \sim \mathcal{N}(\mu, \mathbf{K}). \quad \text{then:}$$

$$p(f_A|f_B) = \mathcal{N}(\cdots, \cdots)$$

$$p(f_B|f_A) = \mathcal{N}(\cdots, \cdots)$$

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

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

Another view: from lin. regression to GPs

- Bayesian linear regression: $y = \phi(x)w + \epsilon$

$$\begin{aligned} p(y|x) &= \int_w p(y|w, x) \quad p(w) = \\ &= \int_w \mathcal{N}(\phi(x)w, \sigma^2) \mathcal{N}(0, \sigma_w^2) \end{aligned}$$

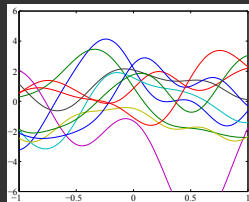
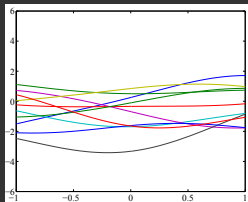
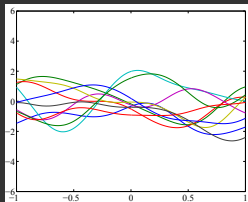
- Gaussian process: $y = f(x) + \epsilon$:

$$\begin{aligned} p(y|x) &= \int_f p(y|f, x) \quad p(f|x) = \\ &= \int_f \mathcal{N}(f, \sigma^2) \mathcal{N}(\mu(x), k(X, X)) \end{aligned}$$

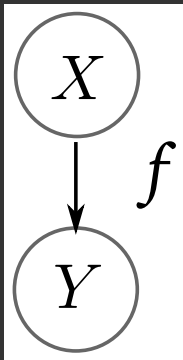
where k is any valid *covariance function*.

Covariance samples and hyperparameters

- The hyperparameters of the cov. function define the properties (and NOT an explicit form) of the sampled functions

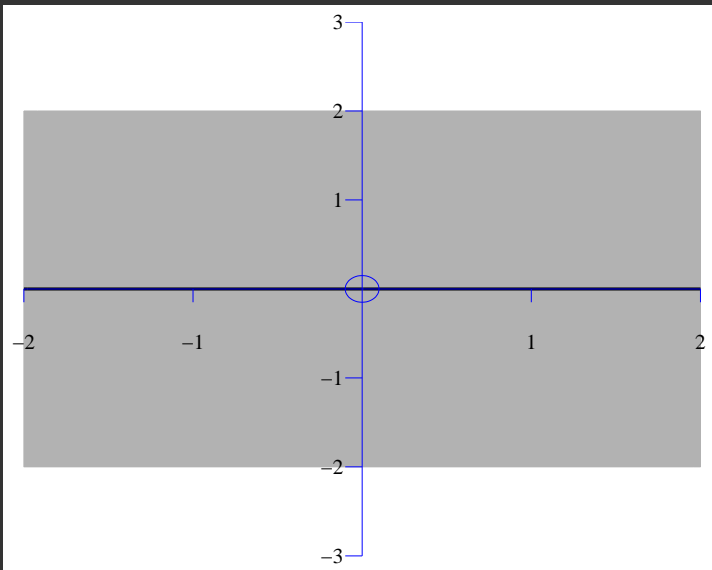


Formally...

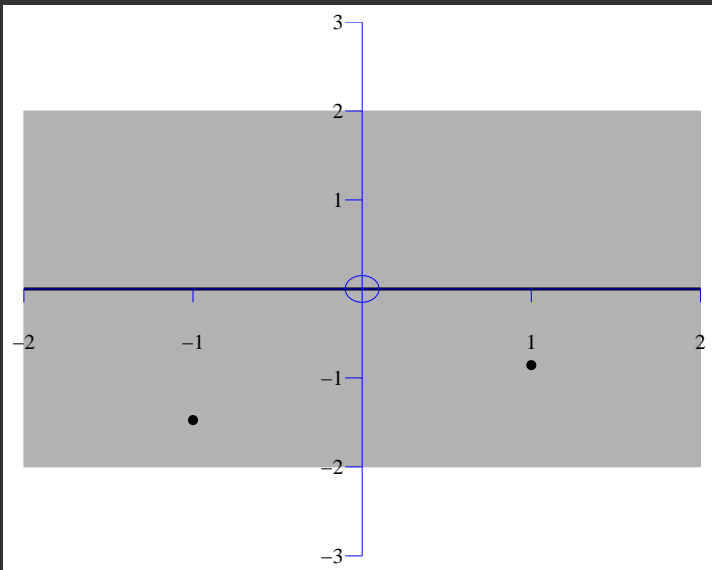


- We write: $f \sim \mathcal{GP}(0, k(x, x))$
- Optimize w.r.t the parameters of k

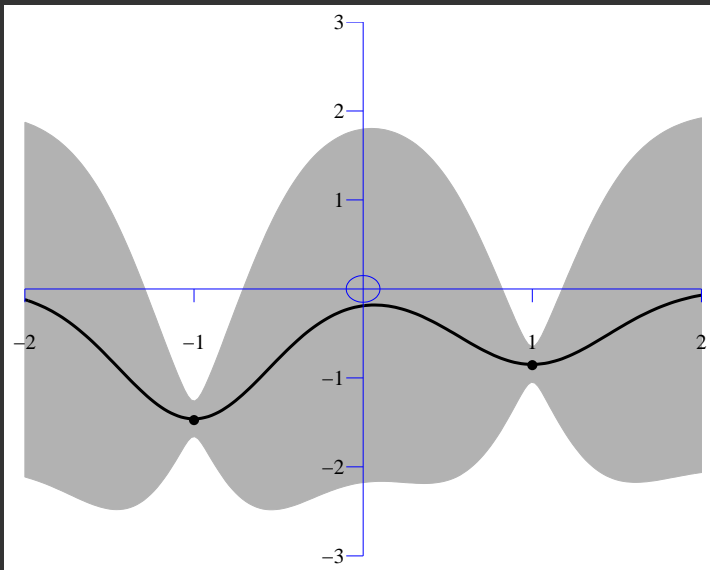
Fitting the data



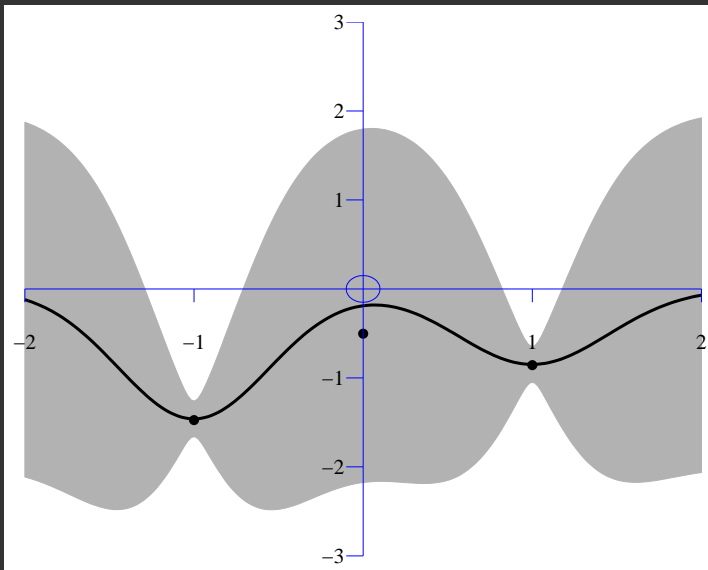
Fitting the data



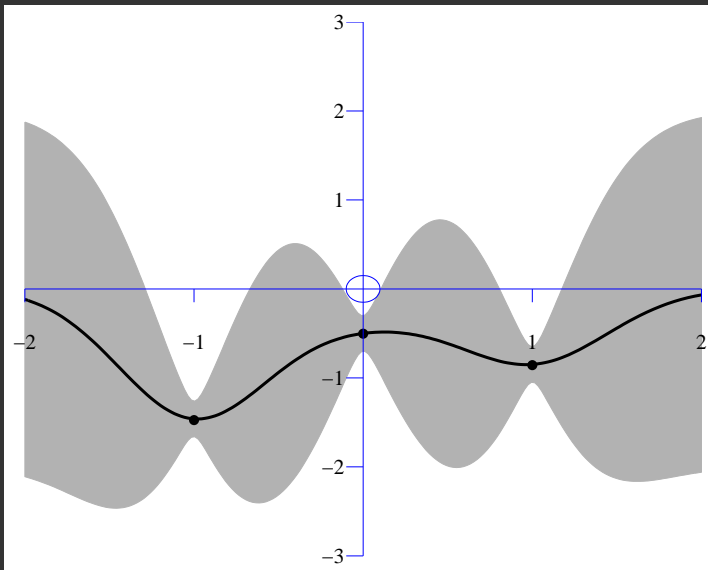
Fitting the data



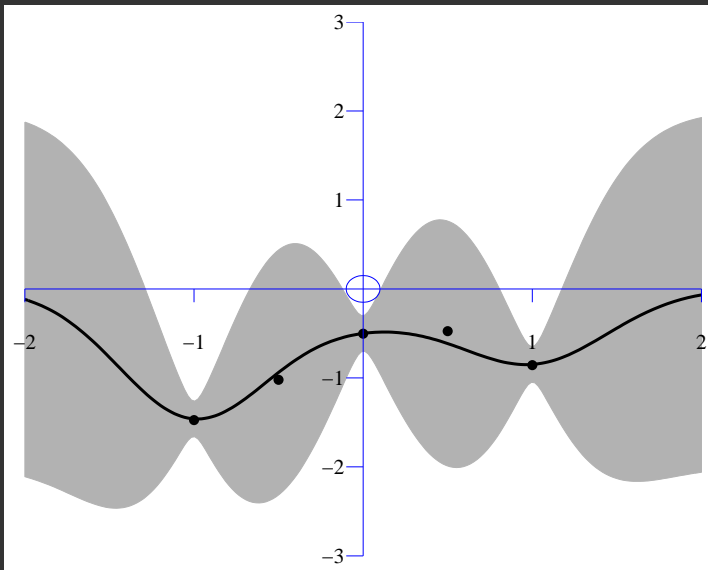
Fitting the data



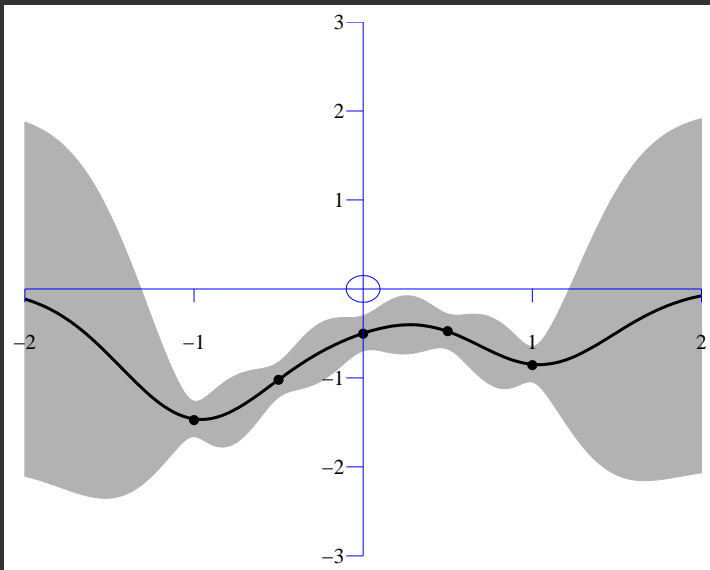
Fitting the data



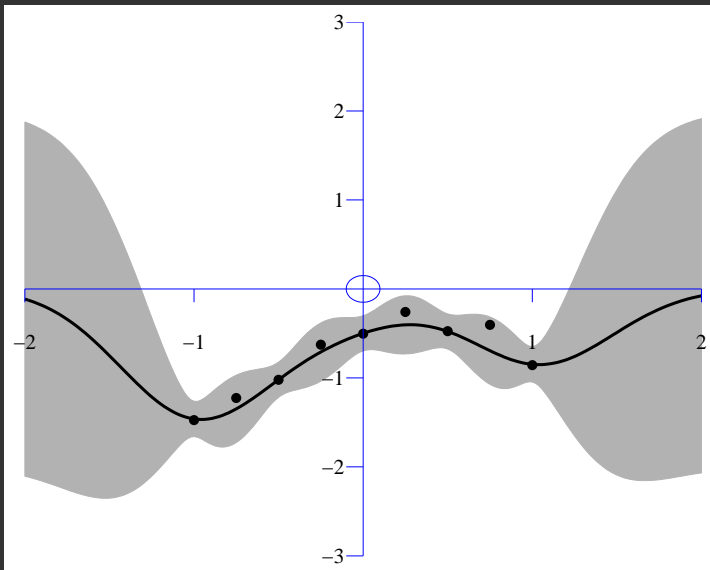
Fitting the data



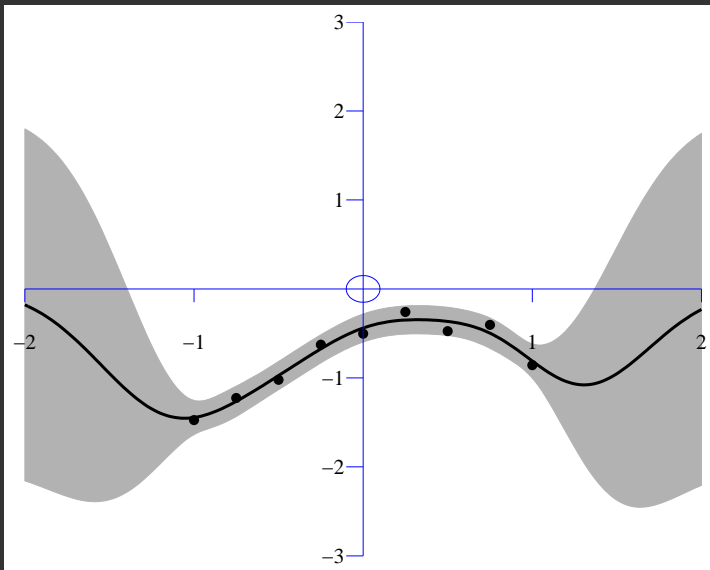
Fitting the data



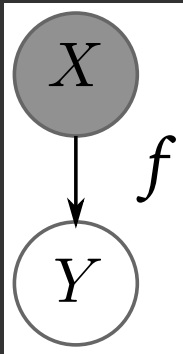
Fitting the data



Fitting the data

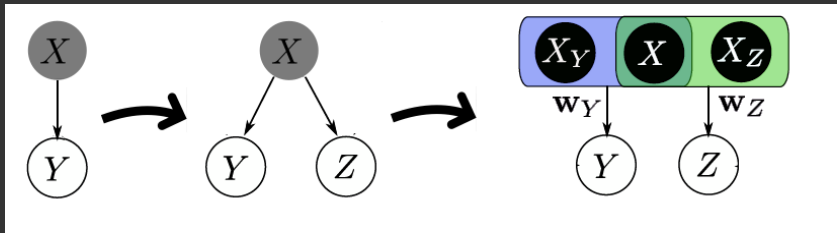


Unsupervised learning: GP-LVM

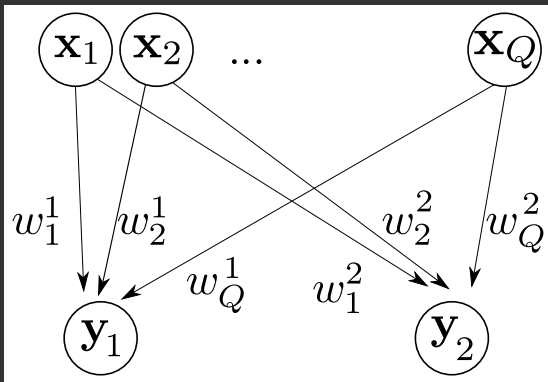


- If X is unobserved, treat it as a parameter and optimize over it.
- X is called the *latent space* assumed to have generated the (noisy) data.
- GP-LVM is interpreted as non-linear PPCA.

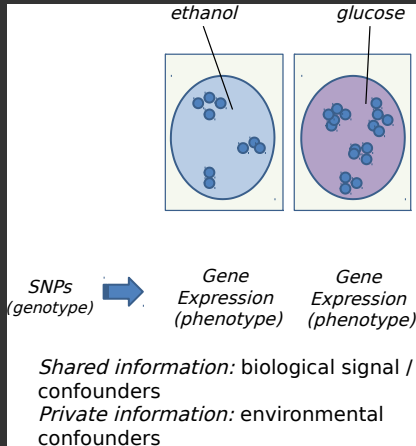
Manifold Relevance Determination



- 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.

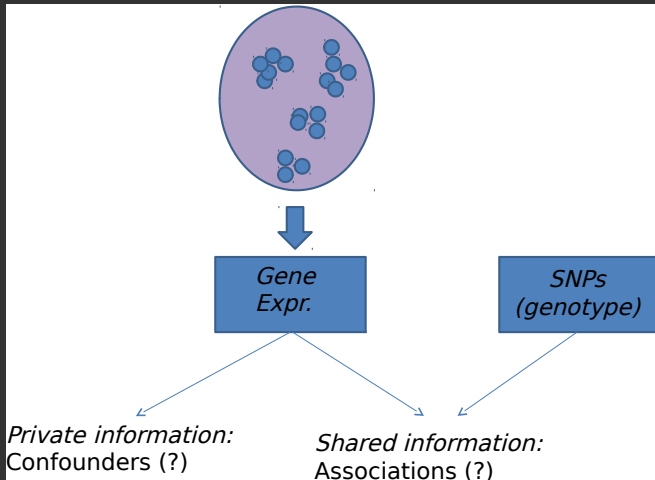


Consolidating complementary experimental data



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

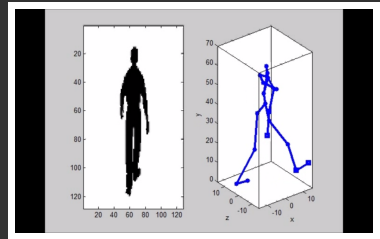
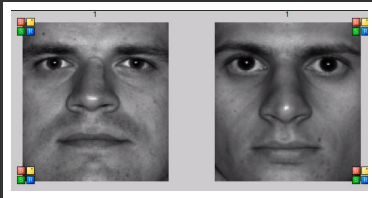
Discovering commonalities in heterogeneous data



Example 3

Motion capture / silhouette

Yale faces



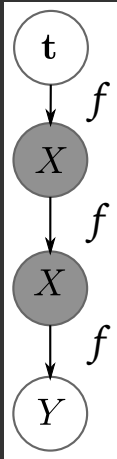
Summary

- Observed data are noisy. Assuming a **latent space** representation helps in modelling/analysis.
- The emerging structure of the latent space helps in **data understanding**.
- All of our **assumptions** can be naturally taken into account in the latent space.
- We can obtain **temporal, multi-view, deep models**.

Thanks

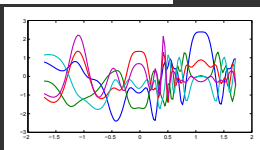
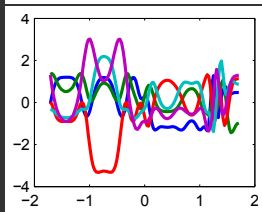
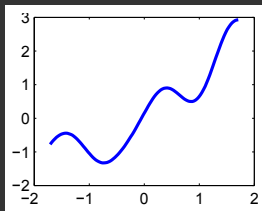
Thanks to Neil Lawrence, James Hensman, Michalis Titsias, Carl Henrik Ek.

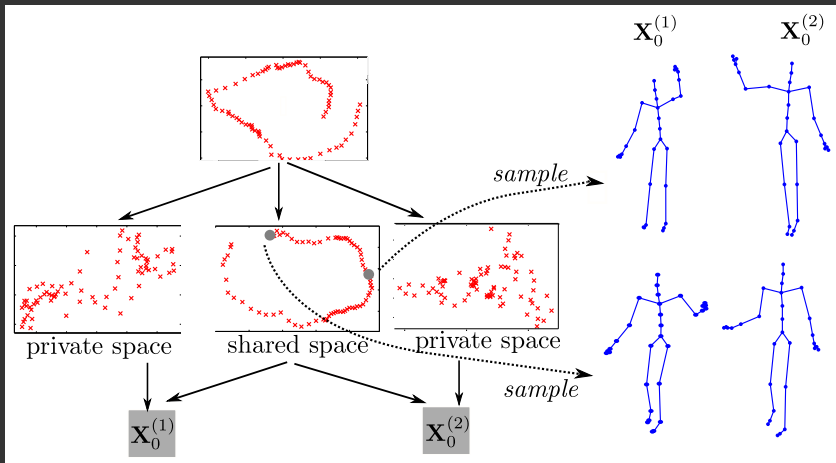
Deep Gaussian processes



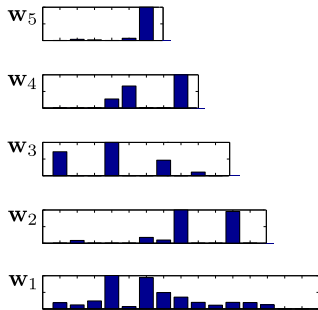
- Construct latent spaces hierarchically:
$$f = f_1(f_2(f_3 \cdots (t)))$$
- Supervised / unsupervised
- A deep GP is NOT a GP! Can learn much more complicated functions!

Sampling from a deep GP





Optimised
weights



X_5

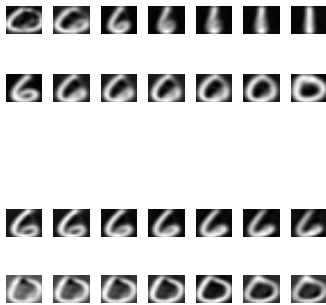
X_4

X_3

X_2

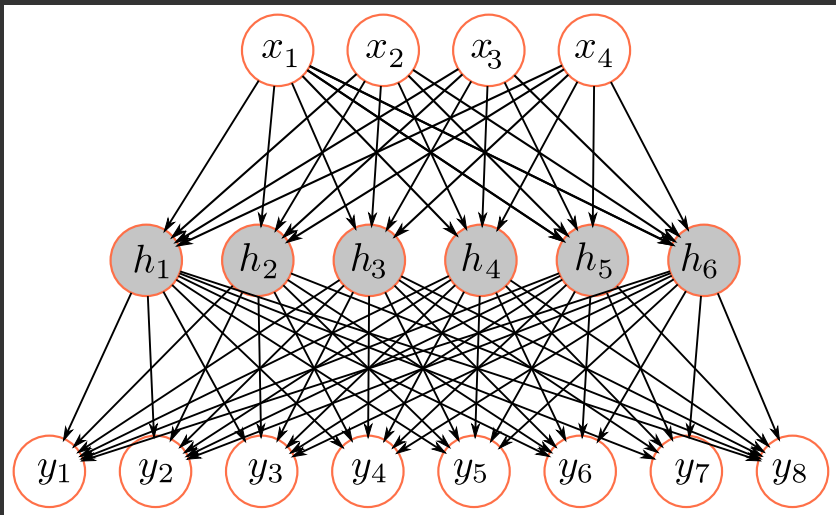
X_1

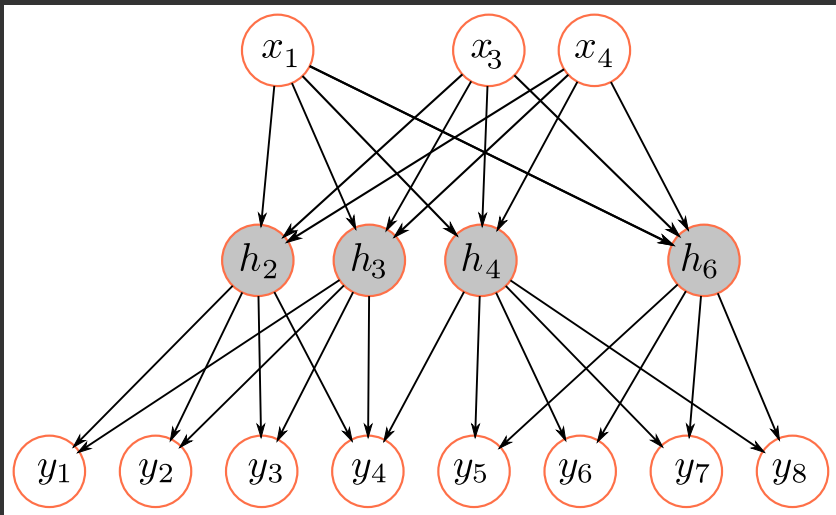
Outputs obtained
after sampling
from (certain nodes)
of layers 5,4,2,1



Generic
feature
encoding

Local
feature
encoding





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