

# Gaussian Processes

Instructor: Hemanth Venkateswara  
Computer Science & Engineering

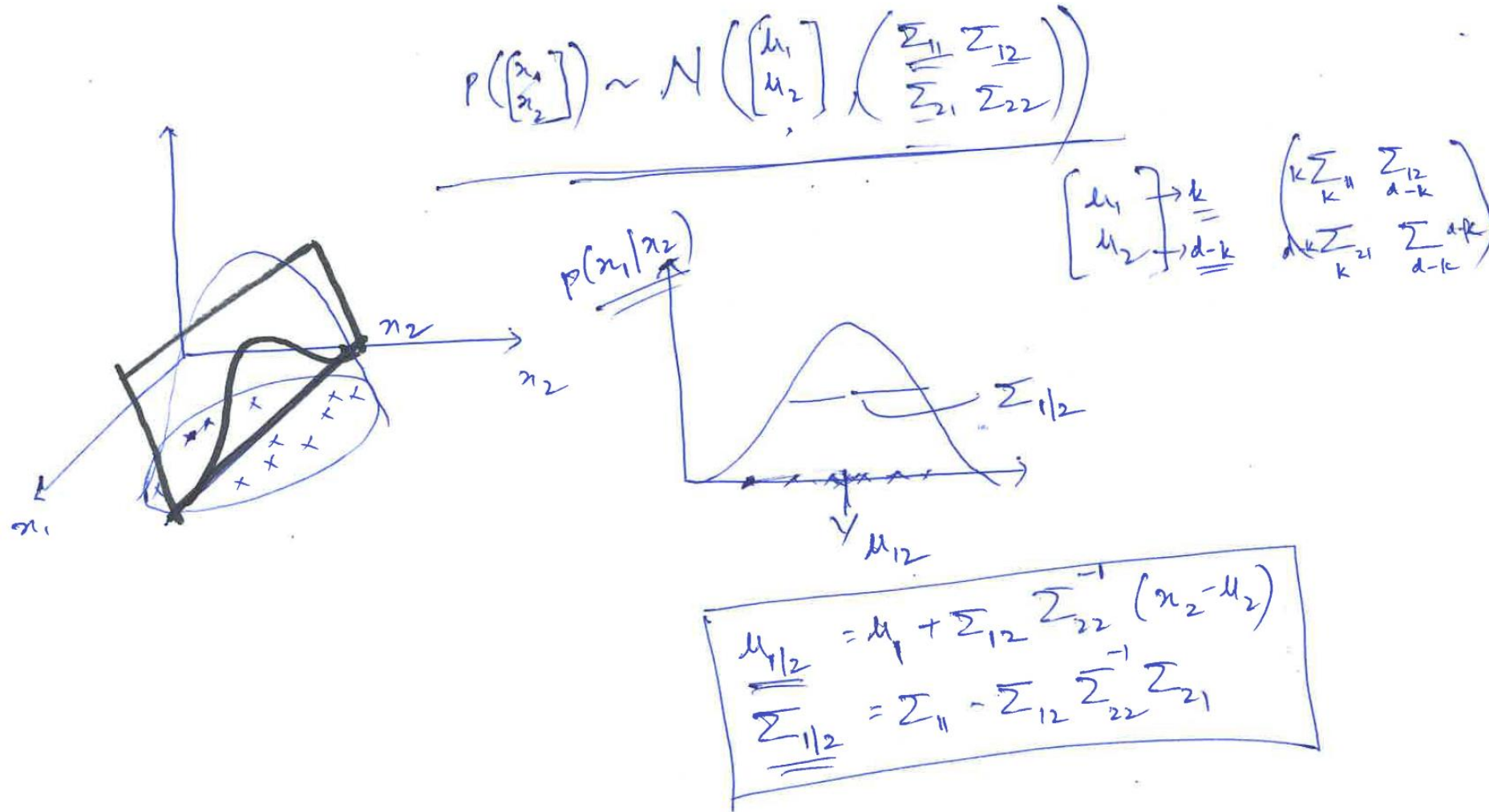


# Overview

## Introduction to Gaussian Processes

- Gaussian Conditionals
- Sampling from a Gaussian Distribution
- Gaussian Processes for Regression

# Gaussian Conditionals



# Gaussian Conditionals

**Theorem 4.3.1** (Marginals and conditionals of an MVN). *Suppose  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  is jointly Gaussian with parameters*

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{pmatrix} \quad (4.67)$$

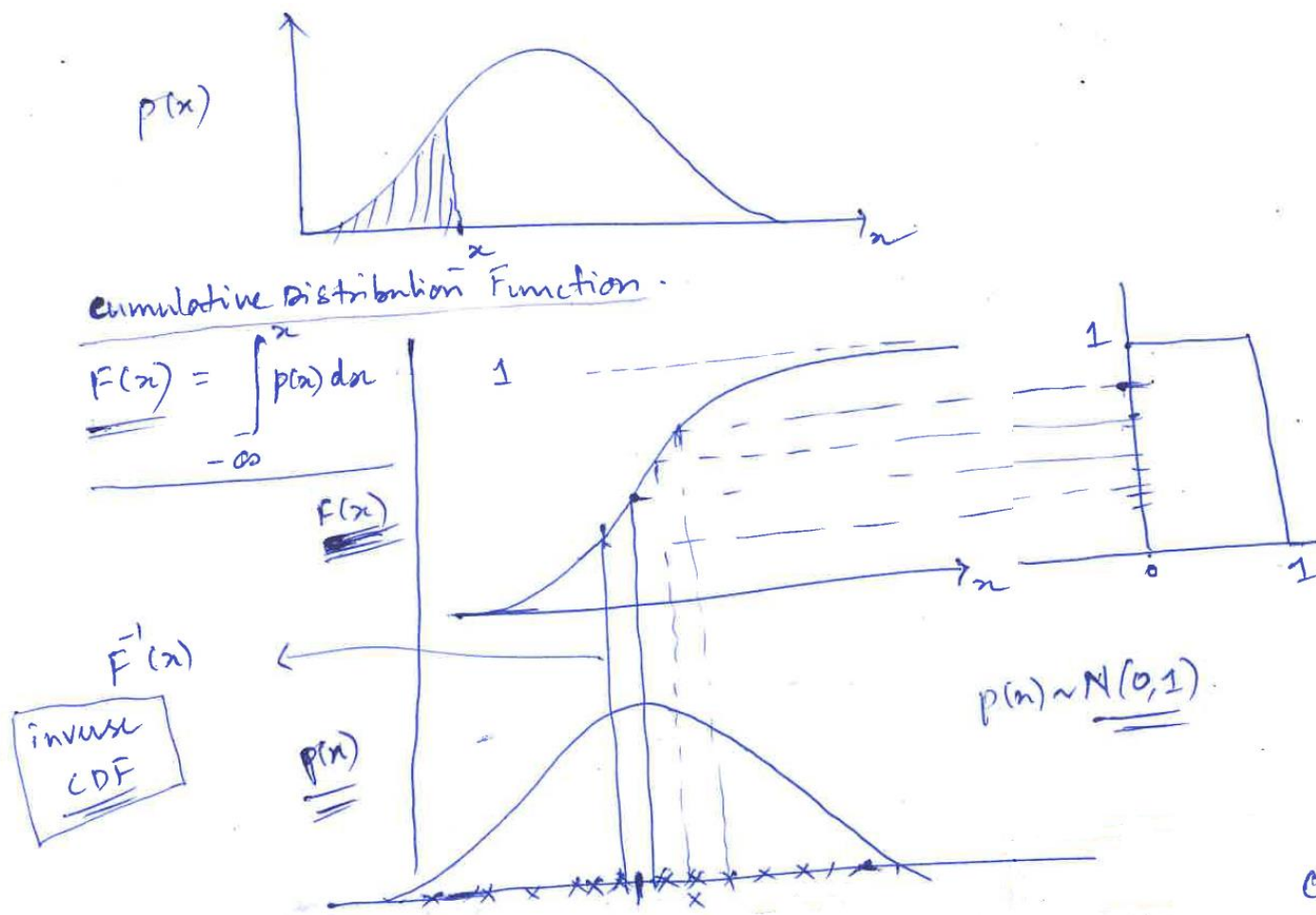
*Then the marginals are given by*

$$\begin{aligned} p(\mathbf{x}_1) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \\ p(\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_2 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) \end{aligned} \quad (4.68)$$

*and the posterior conditional is given by*

$$\begin{aligned} p(\mathbf{x}_1 | \mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}) \\ \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1} \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\Sigma}_{1|2} (\boldsymbol{\Lambda}_{11} \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)) \\ \boldsymbol{\Sigma}_{1|2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1} \end{aligned} \quad (4.69)$$

# Sampling from a Gaussian Distribution



$$x \sim N(\mu, \sigma^2)$$

$$\boxed{p(x) = \frac{1}{\sigma} N\left(\frac{x - \mu}{\sigma}, 1\right) \sim N\left(\frac{x - \mu}{\sigma}, 1\right)}$$

$$P\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

$$p(x_1) \sim N(0, 1) \quad p(x_2) \sim N(0, 1)$$

$$P\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}\right)$$

$$p(x_1) = N(\mu_1, \sigma_1^2) \quad p(x_2) = N(\mu_2, \sigma_2^2)$$

$$P\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right) \rightarrow \Sigma$$

$\downarrow \mu$                        $\downarrow$  (sq root)

$$P\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) \sim \mu + L N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, I\right)$$

$\downarrow$  (sq root)  $\rightarrow \Sigma$

Cholesky Decomposition

$$\boxed{LL^T = \Sigma}$$

# Additional Sources – Gaussian Processes

Richard Turner, Univ. of Cambridge:

<https://www.youtube.com/watch?v=92-98SYOdIY>

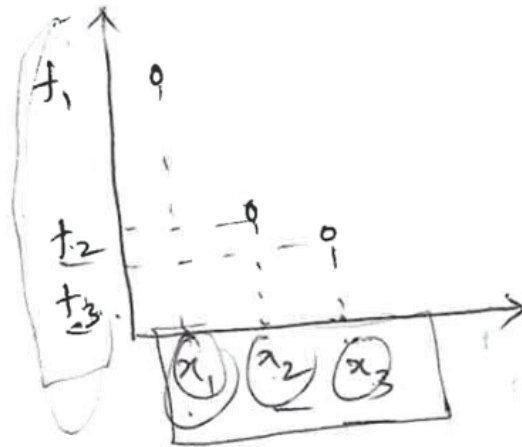
Slides: <http://cbl.eng.cam.ac.uk/pub/Public/Turner/News/imperial-gp-tutorial.pdf>

Nando de Freitas, Univ. of British Columbia

<https://www.youtube.com/watch?v=4vGiHC35j9s>

<https://www.youtube.com/watch?v=MfHKW5z-OOA>

# Gaussian Processes



$l_2 \rightarrow$  horizontal width  
 $\sigma_f \rightarrow$  vertical width  
Kernel function

Kernel

$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2l^2}(x - x')^2\right)$$

$$\begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{bmatrix} \right)$$
  

$$\begin{bmatrix} 1 & 0.7 & 0.5 \\ 0.7 & 1 & 0.4 \\ 0.5 & 0.4 & 1 \end{bmatrix} \rightarrow k(x_1, x_3)$$
  

$$k_{ij} = \sigma_f^2 \exp \left\| \frac{1}{2l^2} x - x' \right\|_2^2$$
  

$$k_{ij} = 0$$
  

$$k_{ij} = 1$$
  

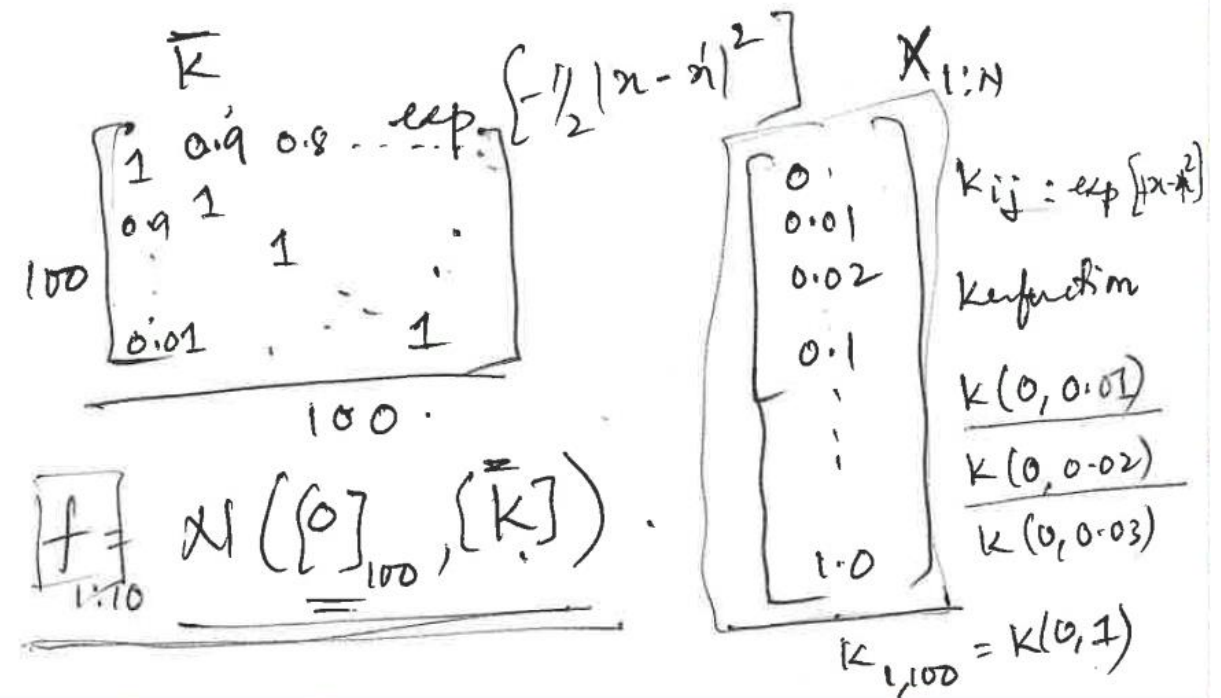
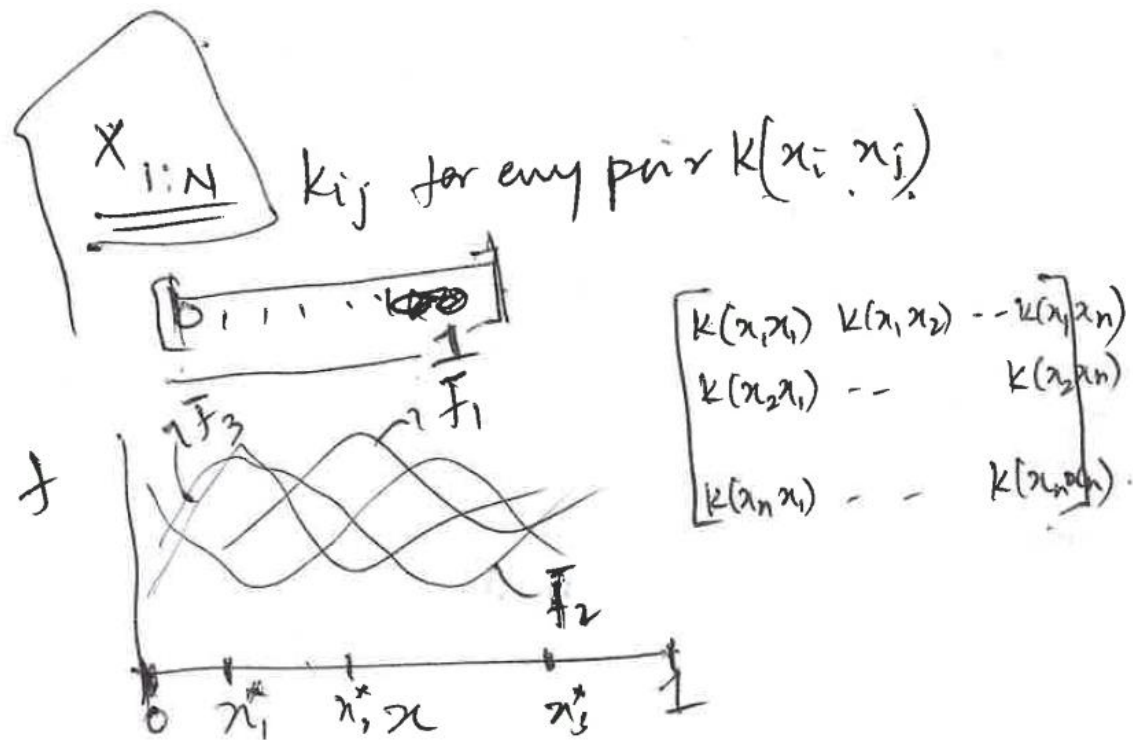
if  $\|x - x'\|^2 \rightarrow \infty$ .

if  $\|x - x'\|^2 \rightarrow 0$ .

RBF



# Gaussian Processes Prior





# Gaussian Processes Prior

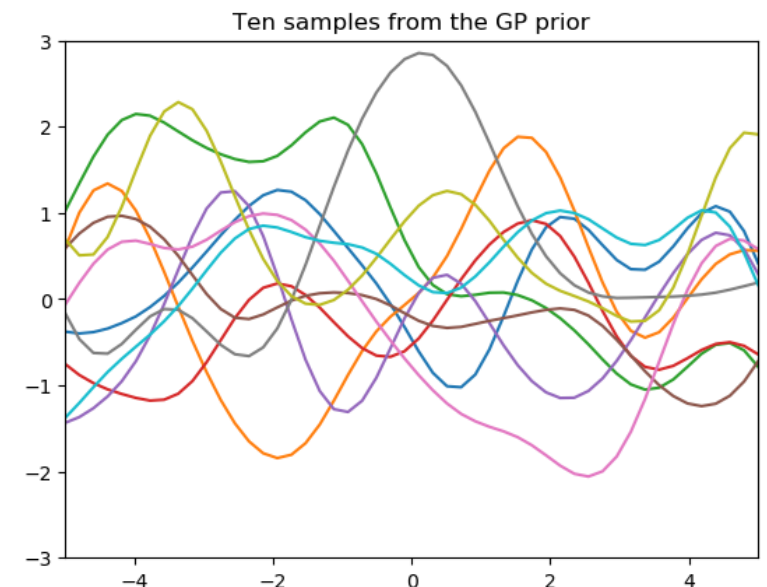
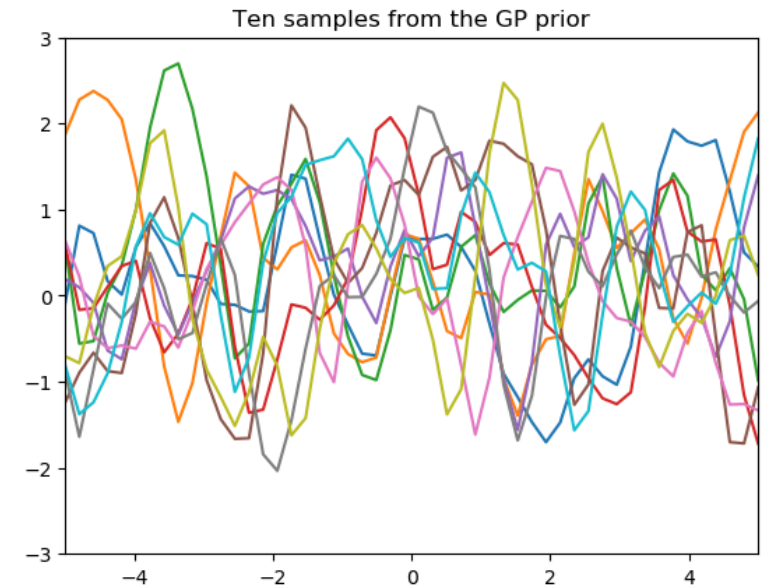
```
from __future__ import division
import numpy as np
import matplotlib.pyplot as plt

def kernel(a, b):
    kernelParameter = 0.1
    sqdist = np.sum(a**2,1).reshape(-1,1) + np.sum(b**2,1) - 2*np.dot(a, b.T)
    return np.exp(-.5 * (1/kernelParameter) * sqdist)

n = 50
Xtest = np.linspace(-5, 5, n).reshape(-1,1)
K_ = kernel(Xtest, Xtest)

# draw samples from the prior at our test points.
L = np.linalg.cholesky(K_ + 1e-6*np.eye(n))
f_prior = np.dot(L, np.random.normal(size=(n,10)))

plt.plot(Xtest, f_prior)
```



Code Source: Nando de Freitas: <https://www.youtube.com/watch?v=4vGiHC35j9s>

# Gaussian Processes: A Distribution Over Functions

In this section, we discuss GPs for regression. Let the prior on the regression function be a GP, denoted by

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x}')) \quad (15.2)$$

where  $m(\mathbf{x})$  is the mean function and  $\kappa(\mathbf{x}, \mathbf{x}')$  is the kernel or covariance function, i.e.,

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \quad (15.3)$$

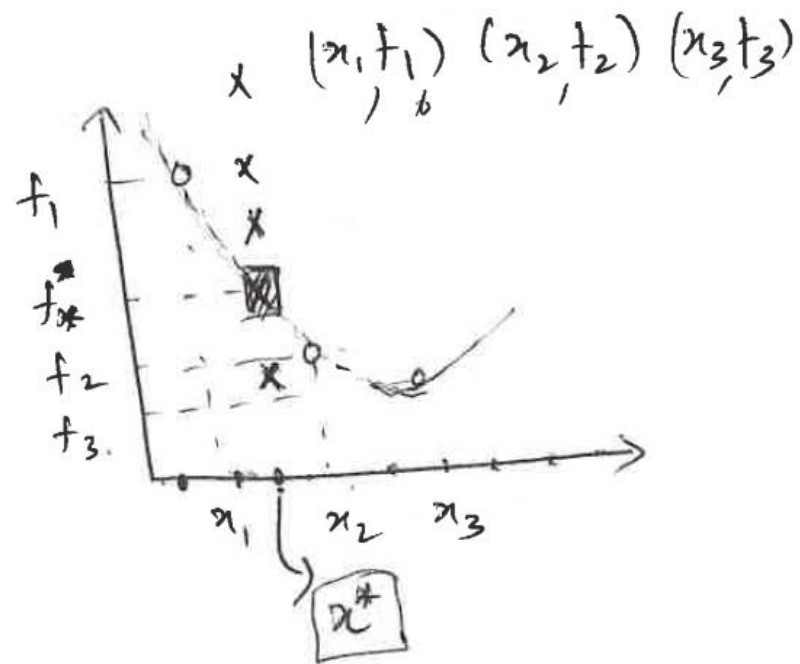
$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))^T] \quad (15.4)$$

We obviously require that  $\kappa()$  be a positive definite kernel. For any finite set of points, this process defines a joint Gaussian:

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \mathbf{K}) \quad (15.5)$$

where  $K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$  and  $\boldsymbol{\mu} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_N))$ .

# Prediction of $f_*$ for sample $x_*$



$$\begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} = \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \right)$$

$$\begin{bmatrix} k_{1*} & k_{2*} & k_{3*} \end{bmatrix}$$

$$p(f_* | \{f_1, f_2, f_3\})$$

$$\begin{bmatrix} k_{11} & k_{12} & k_{13} & k_{1*} \\ k_{21} & k_{22} & k_{23} & k_{2*} \\ k_{31} & k_{32} & k_{33} & k_{3*} \\ k_{*1} & k_{*2} & k_{*3} & k_{**} \end{bmatrix}$$

Labels:  $K$  (top left),  $K^T$  (bottom left),  $K_{**}$  (bottom right),  $K_{*}$  (middle right).

# Gaussian Prior with zero mean

The joint density of the observed data and the latent, noise-free function on the test points is given by

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{pmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix} \right)$$

where we are assuming the mean is zero, for notational simplicity. Hence the posterior predictive density is

$$\begin{aligned} p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \\ \boldsymbol{\mu}_* &= \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{y} \\ \boldsymbol{\Sigma}_* &= \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_* \end{aligned}$$

In the case of a single test input, this simplifies as follows

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(f_* | \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{y}, k_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*)$$

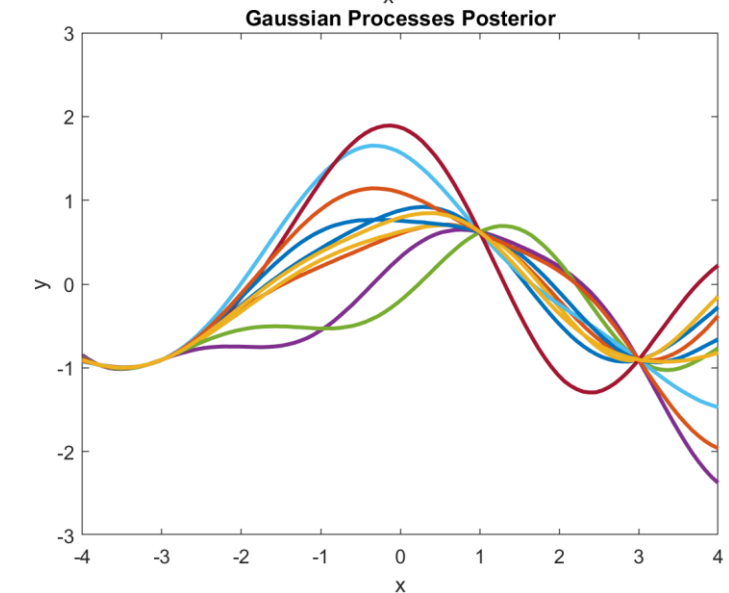
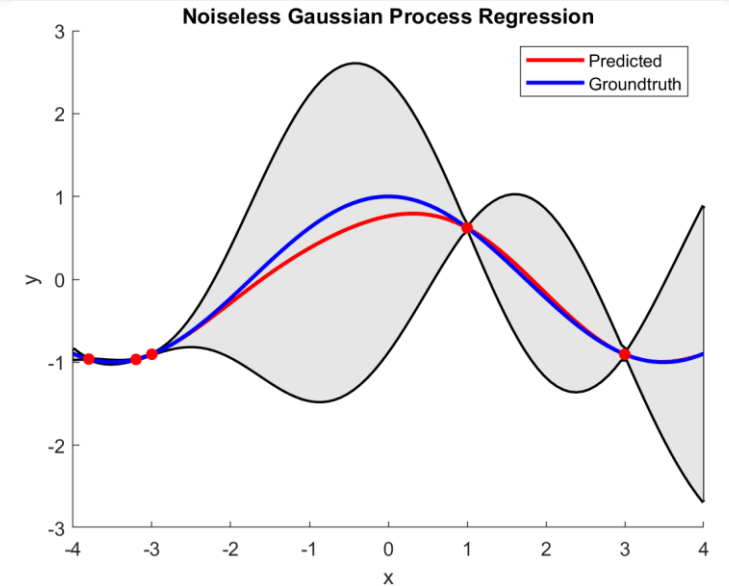
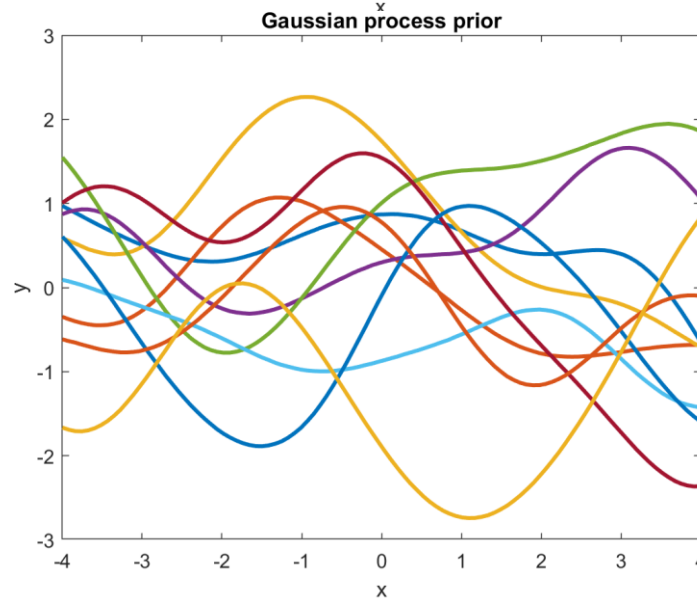
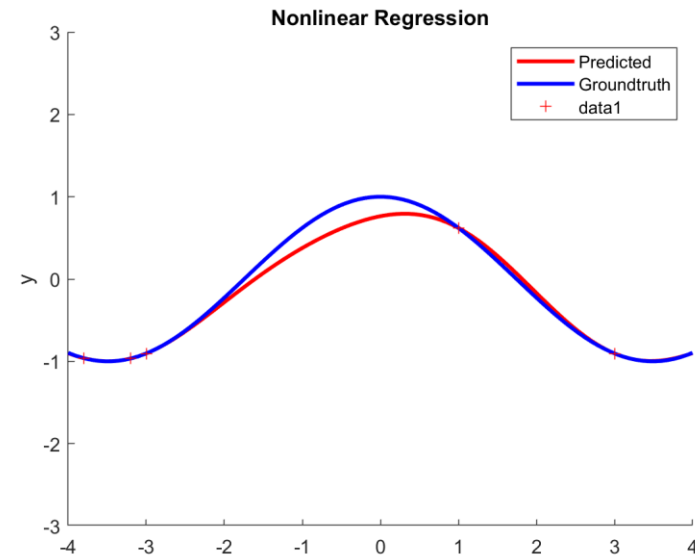
where  $\mathbf{k}_* = [\kappa(\mathbf{x}_*, \mathbf{x}_1), \dots, \kappa(\mathbf{x}_*, \mathbf{x}_N)]$  and  $k_{**} = \kappa(\mathbf{x}_*, \mathbf{x}_*)$ .

# Gaussian Processes Nonlinear Regression

Gaussian Processes Regression:  
Probabilistic Interpretation

$$D = \{(x_i, f_i), i = 1:N\}$$

$$p(f|D) = \frac{p(D|f)p(f)}{p(D)}$$



# GP Algorithm – Zero Mean Prior

$$\bar{f}_* = \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{y} \qquad \boldsymbol{\alpha} = \mathbf{K}^{-1} \mathbf{y} = \mathbf{L}^{-T} \mathbf{L}^{-1} \mathbf{y}$$

$$\mathbf{K} = \mathbf{L} \mathbf{L}^T$$

Solve for  $m$  in System of Linear Equations  $Lm = y$

This is more stable than estimating  $L^{-1}y$

Similarly, solve for  $\alpha$  in system of linear equations  $L^T \alpha = m$

This is more stable than estimating  $L^{-T}m$

---

## Algorithm 15.1: GP regression

---

- 1  $\mathbf{L} = \text{cholesky}(\mathbf{K} + \sigma_y^2 \mathbf{I});$
  - 2  $\boldsymbol{\alpha} = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y});$
  - 3  $\mathbb{E}[f_*] = \mathbf{k}_*^T \boldsymbol{\alpha};$
  - 4  $\mathbf{v} = \mathbf{L} \setminus \mathbf{k}_*;$
  - 5  $\text{var}[f_*] = \kappa(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v};$
  - 6  $\log p(\mathbf{y}|\mathbf{X}) = -\frac{1}{2} \mathbf{y}^T \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{N}{2} \log(2\pi)$
-