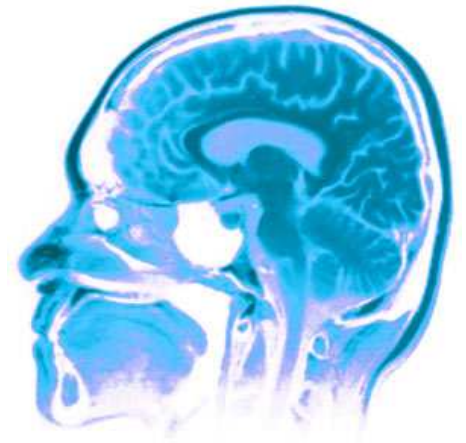




CPSC540



Gaussian Processes



Nando de Freitas

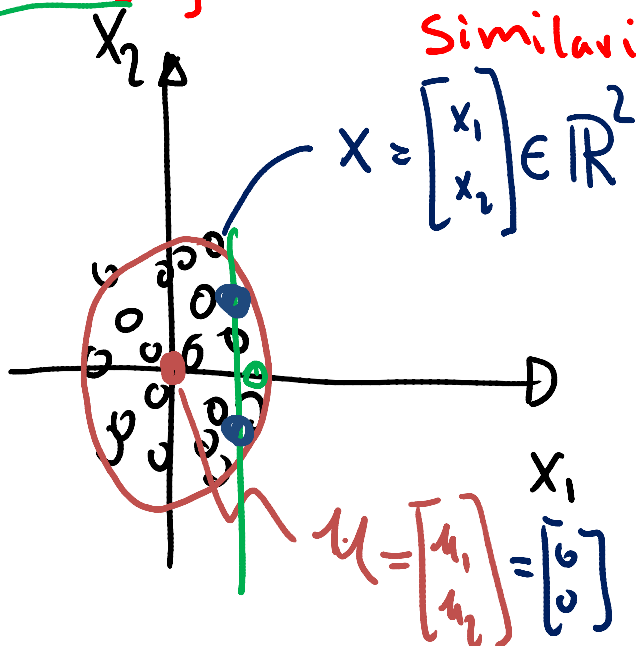
January 2013

KPM Book Sections 4.3 and 15.2.

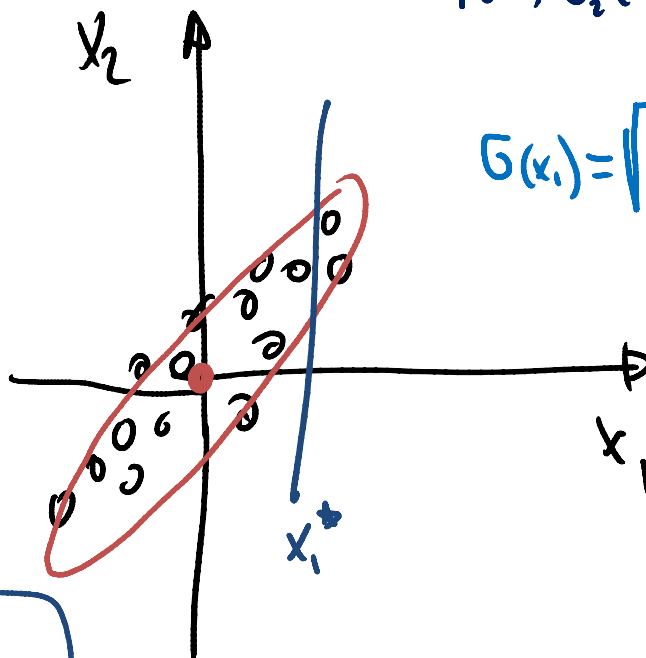
Gaussian basics

$$\left. \begin{aligned} [1 \ 0] \begin{bmatrix} 1 \\ 0 \end{bmatrix} &= 1 \\ [1 \ 0] \begin{bmatrix} 0 \\ 1 \end{bmatrix} &= 0 \end{aligned} \right\}$$

dot products
measure
similarity



$$\rho_{x_1 x_2} = \frac{\text{cov}(x_1, x_2)}{\sigma_1(x_1) \sigma_2(x_2)}$$



$$\sigma(x_i) = \sqrt{\mathbb{E}(x_i^2)}$$

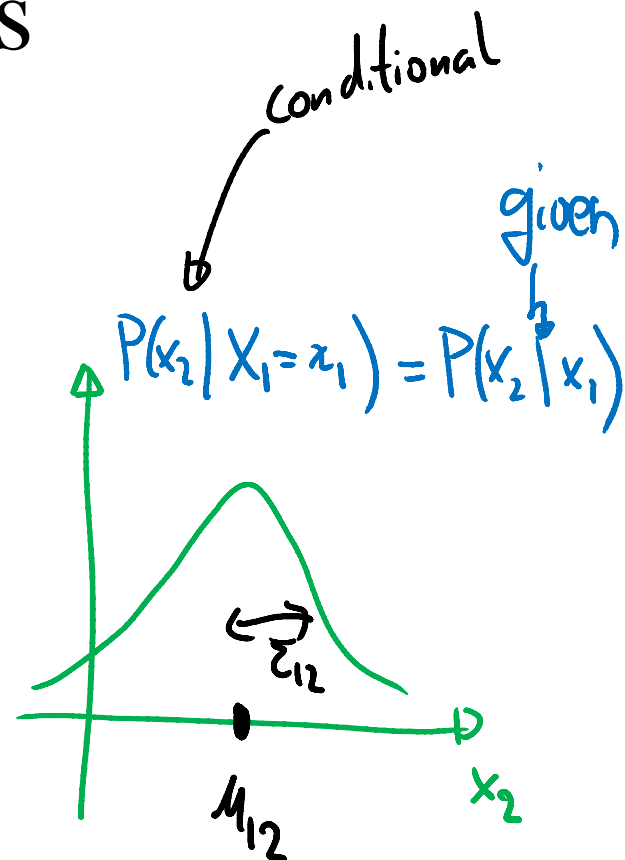
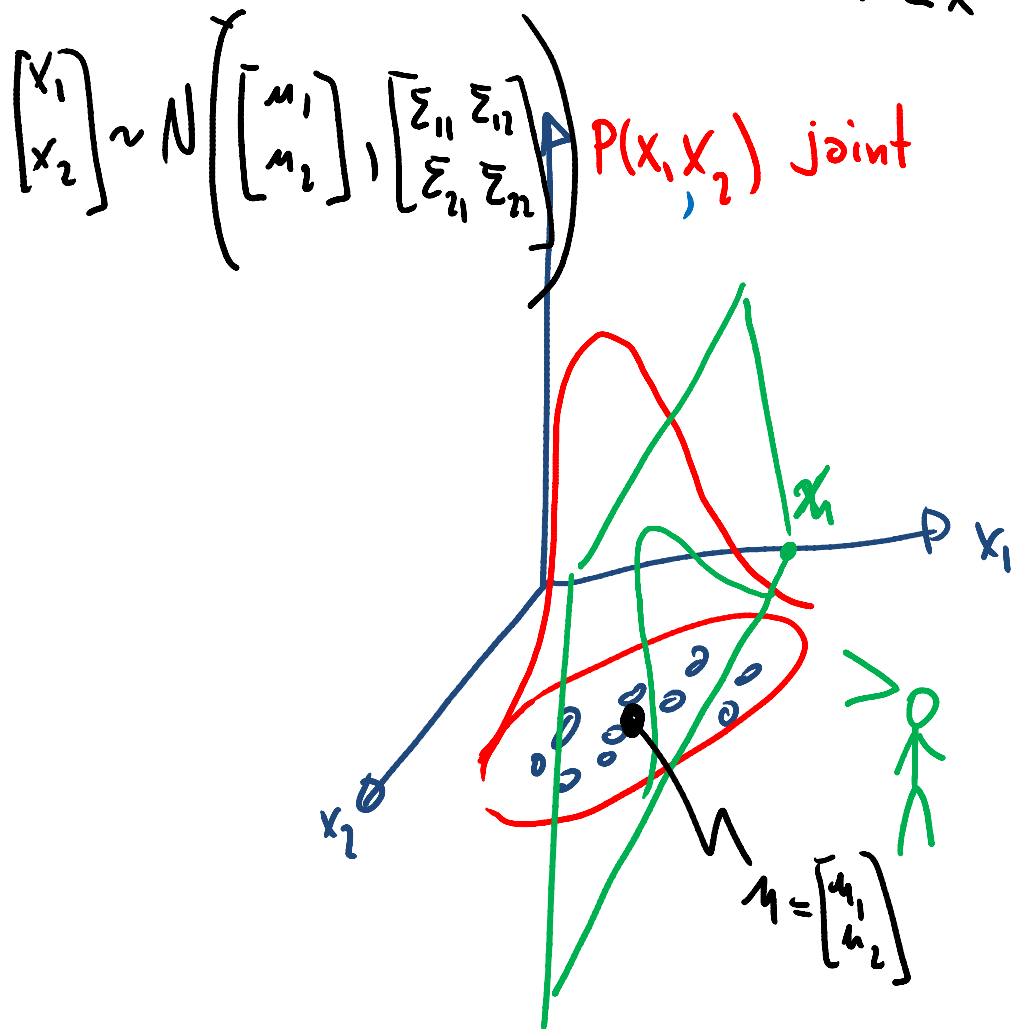
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

$\mathbb{E}(x_1 x_2)$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \right)$$

Gaussian basics

$$x^T \Sigma x$$

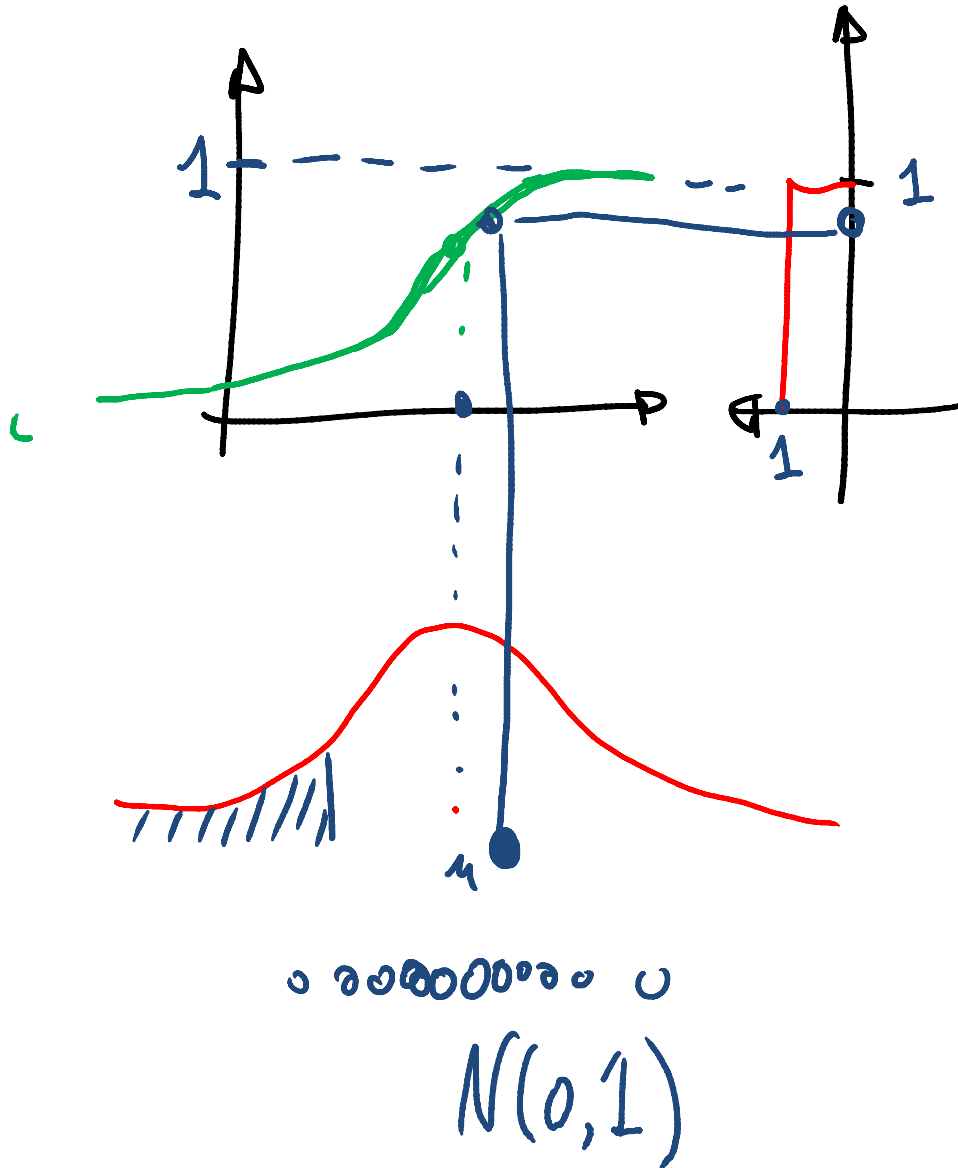


$$\mu_{12} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\Sigma_{12} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Cholesky
 $\Sigma = LL^T$

Gaussian basics



$$X_i \sim N(0, 1)$$

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

$$\sim \mu + \sigma N(0, 1)$$

$$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix}_i \sim N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

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

$\uparrow \mu$ $\uparrow \Sigma$

$$X \sim N(\mu, \Sigma)$$

$$X \sim \mu + L N(0, I)$$

Multivariate Gaussian Theorem (see KPM)

Theorem 4.2.1 (Marginals and conditionals of an MVN). Suppose $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{pmatrix} \quad (4.12)$$

Then the marginals are given by

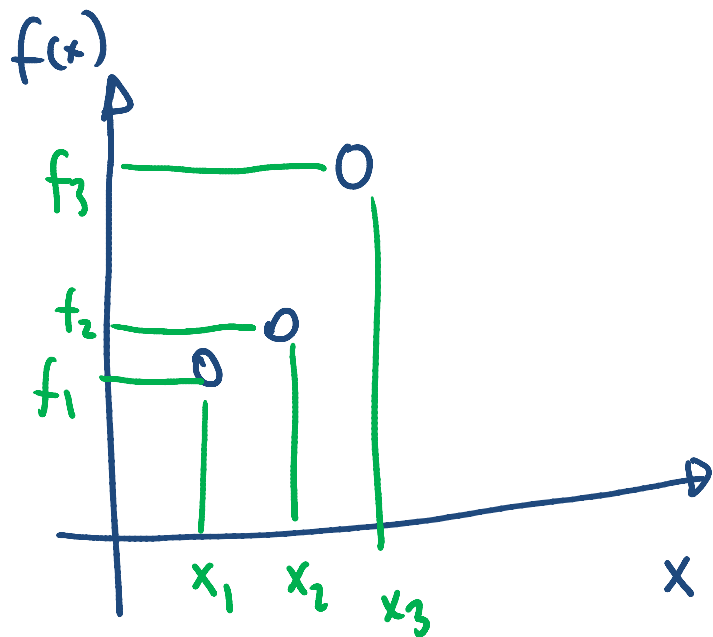
$$\begin{aligned} p(\mathbf{x}_1) &= \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \\ p(\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_2 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) \end{aligned}$$

and the posterior conditional is given by

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

Gaussian basics

x 's given
want to model f 's



$$\begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \sim \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{bmatrix} \right)$$

$$\sim \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.7 & 0.2 \\ 0.7 & 1 & 0.6 \\ 0.2 & 0.6 & 1 \end{bmatrix} \right)$$

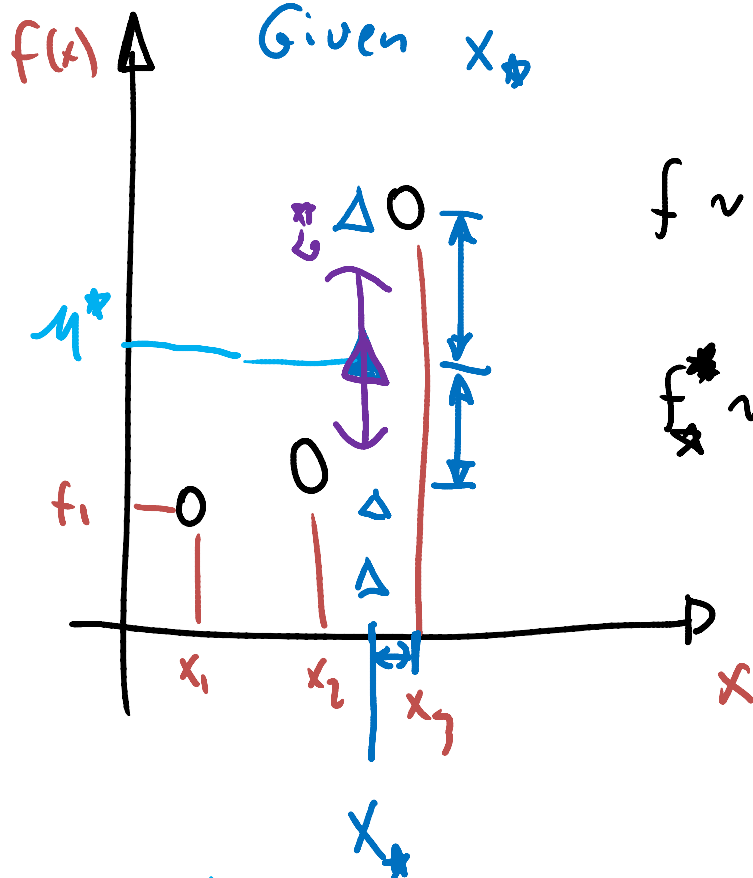
$\nearrow K$

$$k_{ij} = e^{-\lambda \|x_i - x_j\|^2} = \begin{cases} 0 & \|x_i - x_j\| \rightarrow \infty \\ 1 & x_i = x_j \end{cases}$$

$$f \sim N(0, K)$$

Gaussian basics

Given Data $D = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$ $\Rightarrow f_* = ?$



$$f \sim N(0, K)$$

$$K(x_*, x_*) = e^{-\|x_* - x_*\|^2} = 1$$

$$f_* \sim N(0, K(x_*, x_*))$$

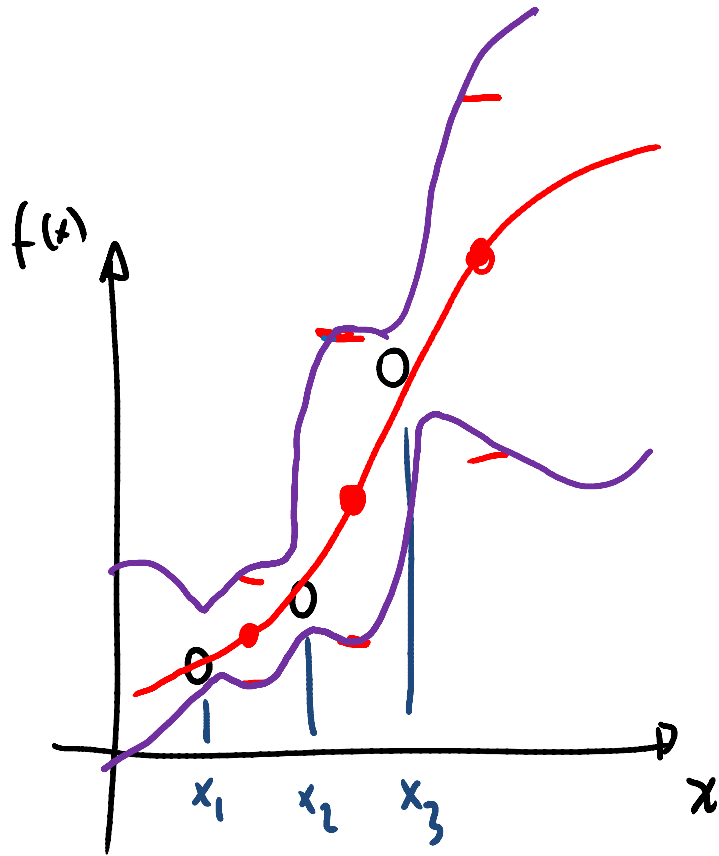
$$k_{i*} = K(x_i, x_*)$$

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim N\left(0, \begin{bmatrix} \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{bmatrix} & \begin{bmatrix} k_{1*} \\ k_{2*} \\ k_{3*} \end{bmatrix} \\ \begin{bmatrix} k_{*1} & k_{*2} & k_{*3} \end{bmatrix} & \begin{bmatrix} k_{**} \end{bmatrix} \end{bmatrix}\right)$$

$$\mu_* = \mathbb{E}(f_*) = K_*^T K^{-1} f$$

$$\sigma_*^2 = K_*^T K^{-1} K_* + K_{**}$$

Gaussian basics



GP: a distribution over functions

A GP is a Gaussian distribution over functions:

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

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

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

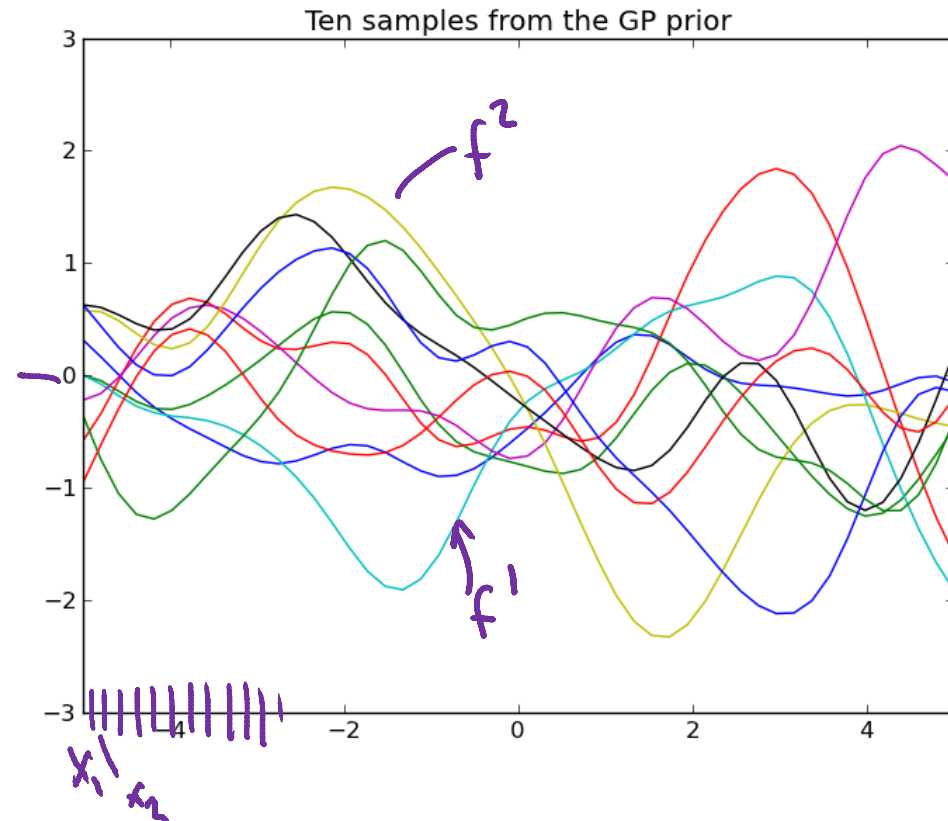
$$k(x, x') = \exp\left(-\frac{1}{2}(x - x')^2\right)$$

Create $x_{1:N}$

Create $\mu = 0_N$, $K_{N \times N}$

$$K = LL^T$$

$$f^i \sim \mathcal{N}(0_N, K) \\ \sim \mathcal{N}(0, \mathbf{I})L$$



Sampling from $P(f)$

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

```
def kernel(a, b):
```

```
    """ GP squared exponential kernel """
```

```
    sqdist = np.sum(a**2,1).reshape(-1,1) + np.sum(b**2,1) - 2*np.dot(a, b.T)
```

```
    return np.exp(-.5 * sqdist)
```

```
n = 50
```

```
# number of test points.
```

```
Xtest = np.linspace(-5, 5, n).reshape(-1,1)
```

```
# Test points.
```

```
K_ = kernel(Xtest, Xtest)
```

```
# Kernel at test points.
```

```
# 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)))
```

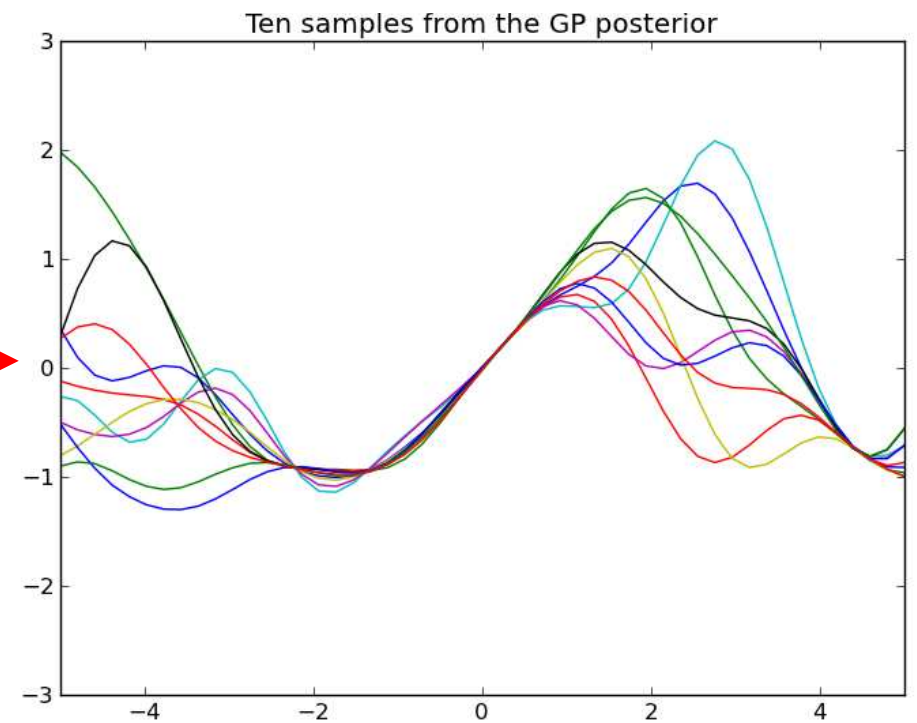
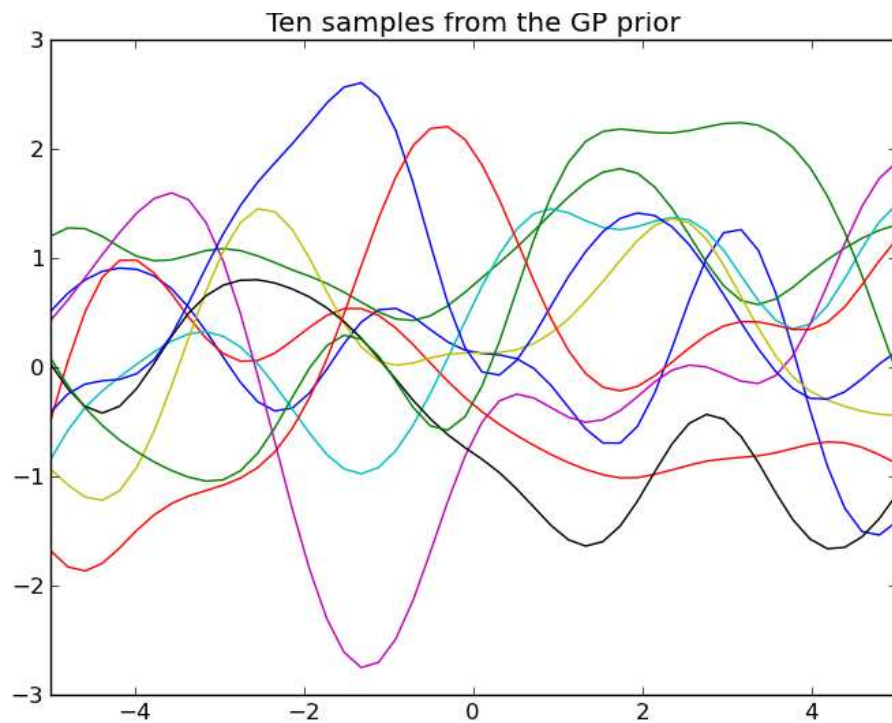
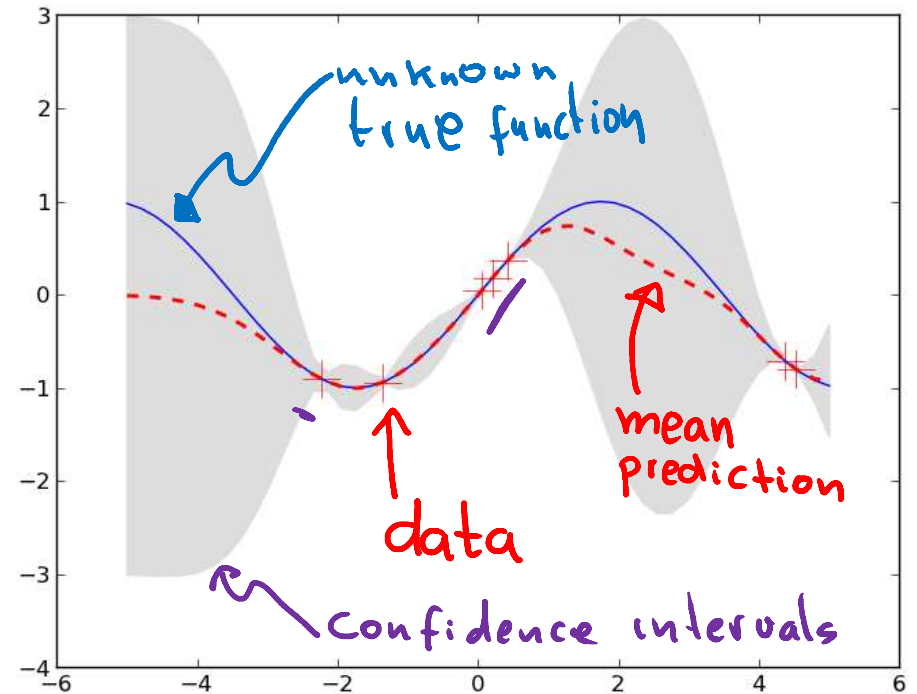
```
plt.plot(Xtest, f_prior)
```

$K = U^T$
 $LN(0, I)$

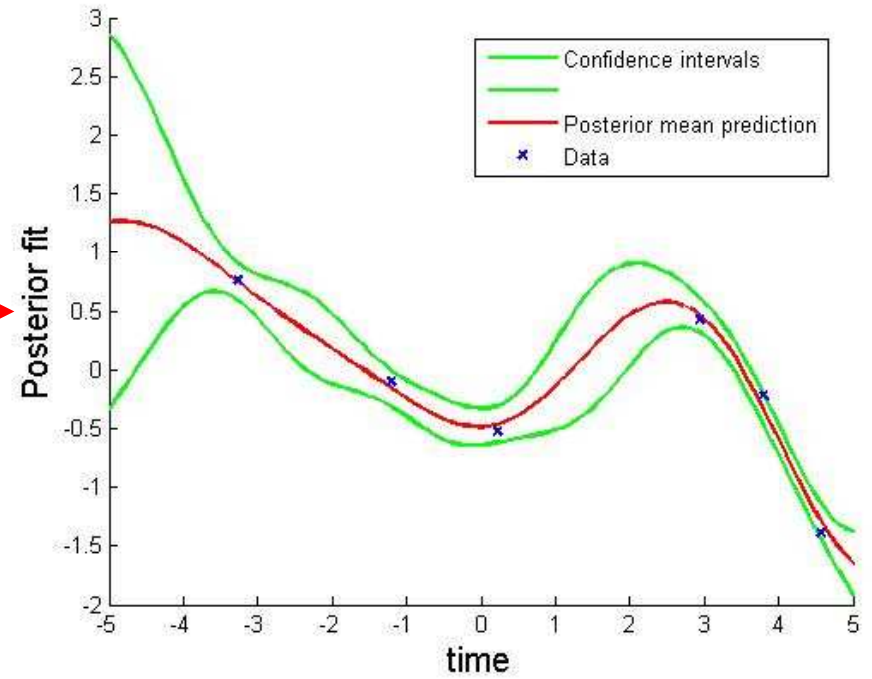
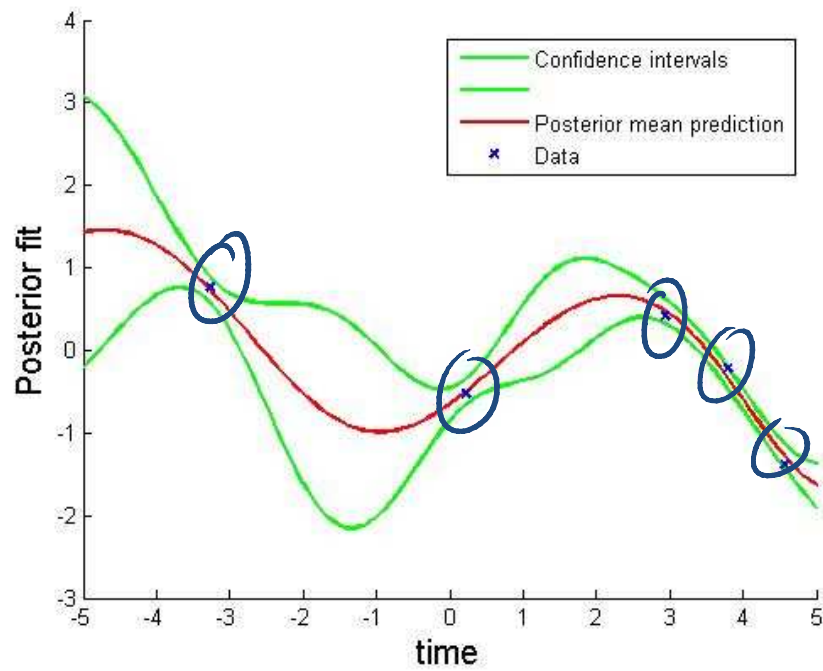
GP posterior

$$\mathcal{D} = \{(\mathbf{x}_i, f_i), i = 1 : N\}$$

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



Active learning with GPs



Noiseless GP regression

we observe a training set $\mathcal{D} = \{(\underline{\mathbf{x}}_i, \underline{f}_i), i = 1 : N\}$, where $\underline{f}_i = \underline{f}(\underline{\mathbf{x}}_i)$

Given a test set $\underline{\mathbf{X}}_*$ of size $N_* \times D$, we want to predict the function outputs $\underline{\mathbf{f}}_*$.

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

where $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$ is $N \times N$, $\mathbf{K}_* = \kappa(\mathbf{X}, \mathbf{X}_*)$ is $N \times N_*$, and $\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$ is $N_* \times N_*$.

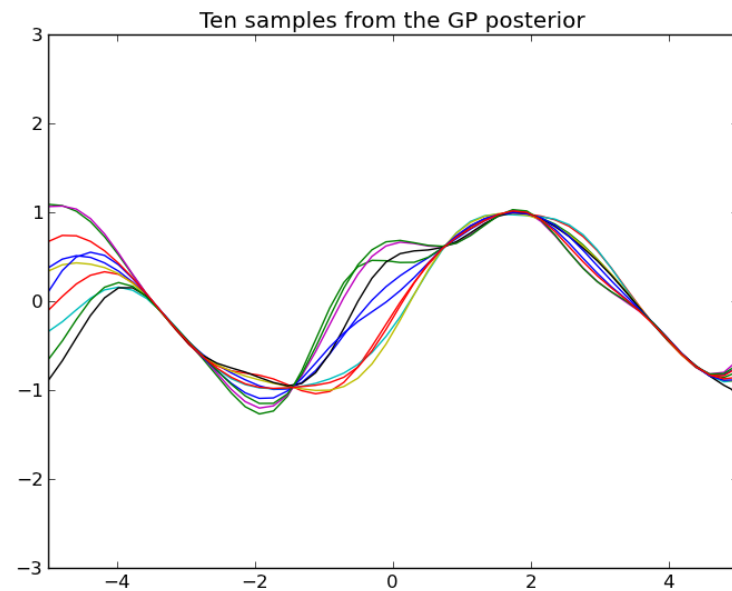
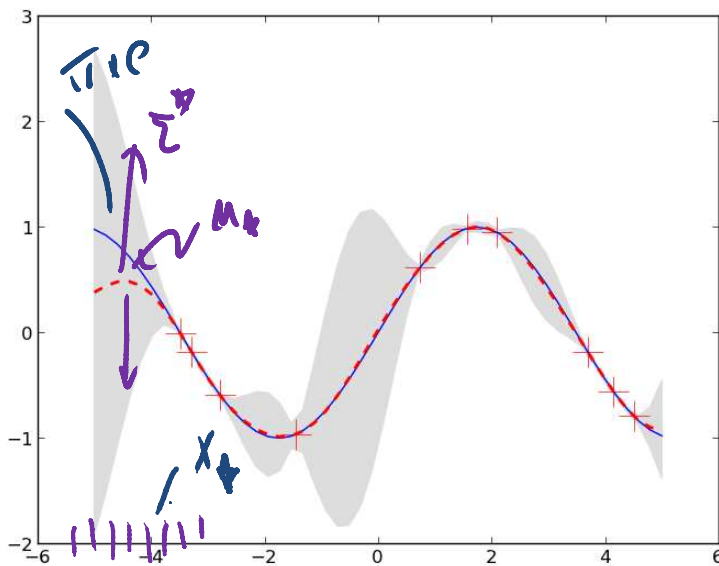
$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$

Noiseless GP regression

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

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{f}) = \mathcal{N}(\mathbf{f}_* | \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \leftarrow$$

$$\begin{aligned} \boldsymbol{\mu}_* &= \boldsymbol{\mu}(\mathbf{X}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (\mathbf{f} - \boldsymbol{\mu}(\mathbf{X})) \\ \boldsymbol{\Sigma}_* &= \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_* \end{aligned}$$



Effect of kernel width parameter

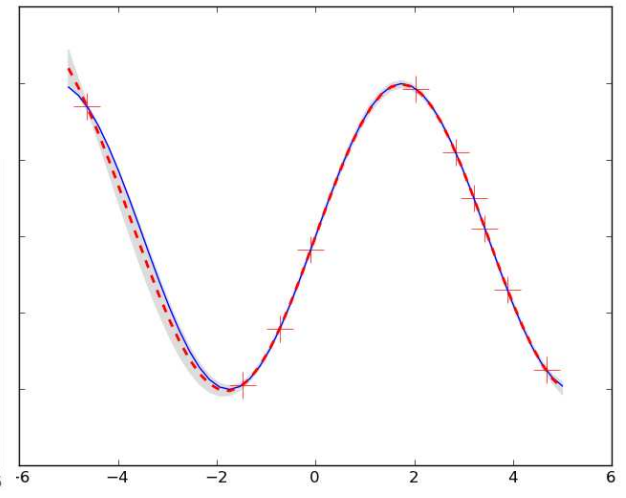
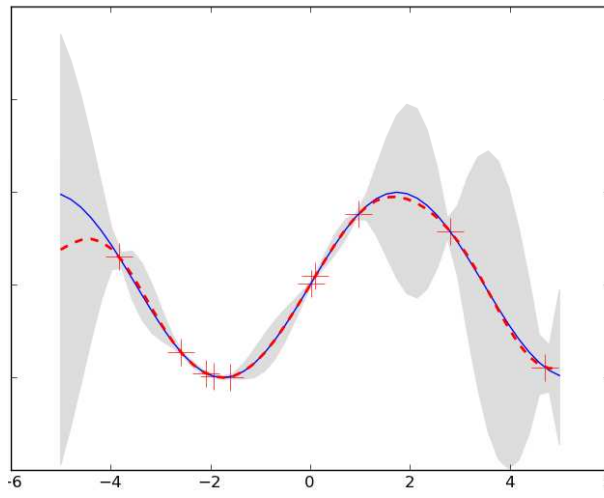
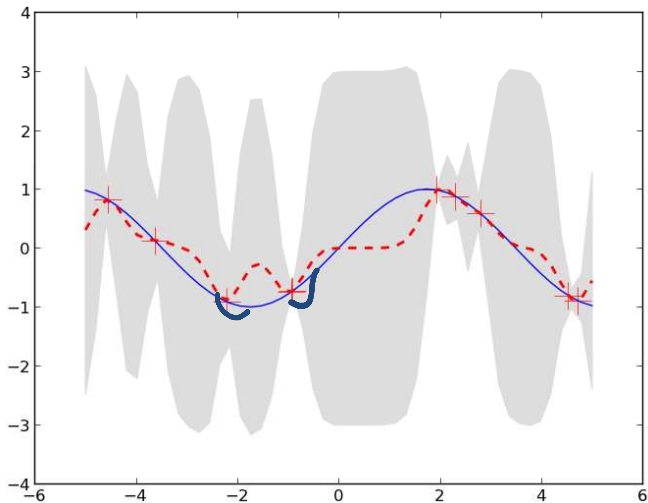
$$\kappa(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x - x')^2\right)$$

Let $\sigma_f = 1$

$\ell^2 = 0.1$ /

$\ell^2 = 1$

$\ell^2 = 10$



Noisy GP regression

noisy

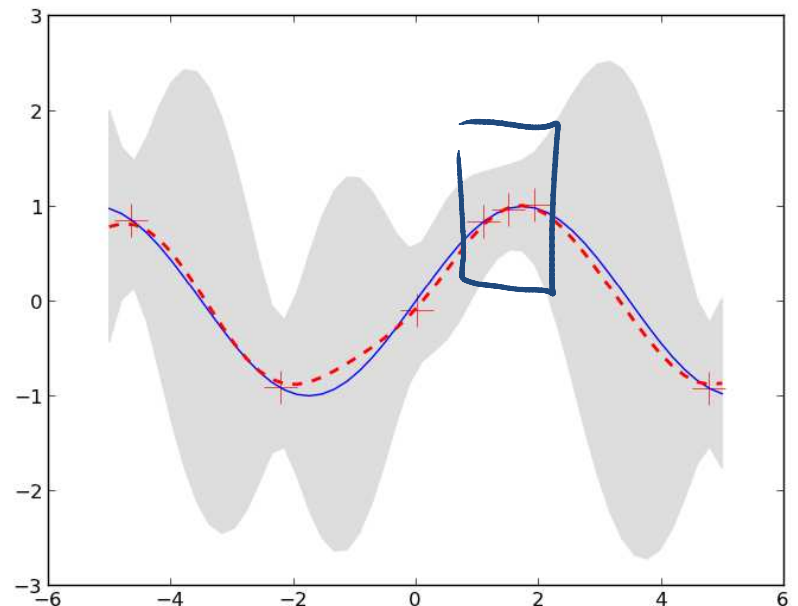
$$y = \underline{f(\mathbf{x})} + \underline{\epsilon}, \text{ where } \epsilon \sim \mathcal{N}(0, \sigma_y^2)$$

$$p(\mathbf{y}|\underline{\mathbf{X}}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f}$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$$

$$p(\mathbf{y}|\mathbf{f}) = \prod_i \mathcal{N}(y_i|f_i, \sigma_y^2)$$

$$\text{cov}[\underline{\mathbf{y}}|\underline{\mathbf{X}}] = \underline{\mathbf{K}} + \underline{\sigma_y^2 \mathbf{I}_N} \triangleq \underline{\mathbf{K}_y}$$



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

thm

$$p(\mathbf{f}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}_*|\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\boldsymbol{\mu}_* = \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{K}_*$$

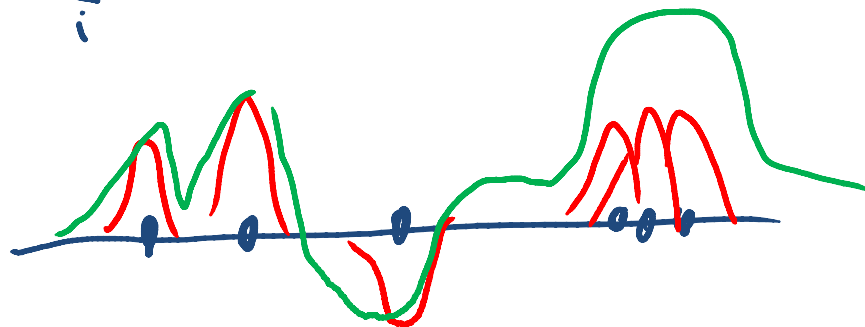
Noisy GP regression

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

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

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

$$\begin{aligned} \bar{f}_* &= \underbrace{\mathbf{k}_*^T}_{1 \times N} \underbrace{\mathbf{K}_y^{-1}}_{N \times N} \underbrace{\mathbf{y}}_{N \times 1} = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}_*) & \alpha &= \underbrace{\mathbf{K}_y^{-1} \mathbf{y}}_{\text{training}} \\ & \text{mean} & & \\ & 1 \times 1 & & \\ & = \mathbf{K}_*^T \alpha & & \\ & = \sum_i \alpha_i e^{-\frac{1}{2} \|\mathbf{x}_i - \mathbf{x}_*\|^2} & & \end{aligned}$$



Noisy GP regression and Ridge

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \delta^2 \|\boldsymbol{\theta}\|_2^2$$

$$\mathbf{X} \in \mathbb{R}^{n \times d}$$

$$\mathbf{y} \in \mathbb{R}^n$$

$$(\mathbf{X}^T \mathbf{X} + \delta^2 \mathbf{I}_d) \boldsymbol{\theta} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$x_i \in \mathbb{R}^{1 \times d}$$

solution can be written as $\boldsymbol{\theta} = \mathbf{X}^T \boldsymbol{\alpha}$, where $\boldsymbol{\alpha} = \delta^{-2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} + \delta^2 \boldsymbol{\theta} = \mathbf{X}^T \mathbf{y}$$

$$\delta^2 \boldsymbol{\theta} = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

$$\boldsymbol{\theta} = \mathbf{X}^T \delta^{-2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \mathbf{X}^T \boldsymbol{\alpha}$$

Noisy GP regression and Ridge

$$(\underbrace{\mathbf{X}^T \mathbf{X}}_{d \times d} + \delta^2 \mathbf{I}_d) \underbrace{\boldsymbol{\theta}}_{d \times 1} = \mathbf{X}^T \mathbf{y}$$

solution can be written as $\boldsymbol{\theta} = \mathbf{X}^T \boldsymbol{\alpha}$, where $\boldsymbol{\alpha} = \delta^{-2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$

$\boldsymbol{\alpha}$ can also be written as follows: $\boldsymbol{\alpha} = (\underbrace{\mathbf{X}\mathbf{X}^T}_{n \times n} + \delta^2 \mathbf{I}_n)^{-1} \mathbf{y}$

$$\delta^2 \boldsymbol{\alpha} = \mathbf{y} - \mathbf{X}\boldsymbol{\theta}$$

$$\delta^2 \boldsymbol{\alpha} = \mathbf{y} - \mathbf{X}\mathbf{X}^T \boldsymbol{\alpha}$$

$$\mathbf{X}\mathbf{X}^T \boldsymbol{\alpha} + \delta^2 \mathbf{I}_n \boldsymbol{\alpha} = \mathbf{y}$$

$$\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^T + \delta^2 \mathbf{I}_n)^{-1} \mathbf{y}$$

$$\mathbf{y}^* = \mathbf{X}^* \boldsymbol{\theta} = \mathbf{X}^* \mathbf{X}^T \boldsymbol{\alpha}$$

Noisy GP regression and Ridge

$$\begin{aligned}
 \underset{|x|}{y^*} &= \underset{|x| \times d}{x^*} \underset{d \times 1}{\theta} \\
 &= x^* X^T \alpha
 \end{aligned}$$

$$\bar{f}_* = \mathbf{k}_*^T \mathbf{K}_y^{-1} \mathbf{y}$$

$$\begin{aligned}
 &= \underbrace{x^* X^T}_{\mathbf{k}_*^T} \underbrace{\left[X X^T + \underset{G_y}{\sigma^2 I_n} \right]^{-1}}_{\mathbf{K}_y^{-1}} \mathbf{y} \\
 &= \mathbf{k}_*^T \mathbf{K}_y^{-1} \mathbf{y}
 \end{aligned}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \text{1} \times d$$

$$\mathbf{k}_*^T = \underbrace{\begin{bmatrix} x^* x_1^T & x^* x_2^T & \dots & x^* x_n^T \end{bmatrix}}_{1 \times n}$$

$$\mathbf{K}_y = \underbrace{X X^T}_{n \times n} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} x_1^T & \dots & x_n^T \end{bmatrix} = \begin{bmatrix} x_1 x_1^T & \dots \\ \vdots & \ddots \\ x_n x_n^T & \dots \end{bmatrix}$$

Learning the kernel parameters

marginal likelihood

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f}$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$$

$$\theta = \ell$$

$$p(\mathbf{y}|\mathbf{f}) = \prod_i \mathcal{N}(y_i|f_i, \sigma_y^2)$$

$$\log p(\mathbf{y}|\mathbf{X}) = \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_y) = -\frac{1}{2}\mathbf{y}\mathbf{K}_y^{-1}\mathbf{y} - \frac{1}{2}\log|\mathbf{K}_y| - \frac{N}{2}\log(2\pi)$$

$$\frac{\partial}{\partial\theta_j}\log p(\mathbf{y}|\mathbf{X}) = \frac{1}{2}\mathbf{y}^T\mathbf{K}_y^{-1}\frac{\partial\mathbf{K}_y}{\partial\theta_j}\mathbf{K}_y^{-1}\mathbf{y} - \frac{1}{2}\text{tr}(\mathbf{K}_y^{-1}\frac{\partial\mathbf{K}_y}{\partial\theta_j}) \leftarrow$$

Numerical computation considerations

$$\mu_{f_*} = \overline{f_*} = \mathbf{k}_*^T \underbrace{\mathbf{K}_y^{-1} \mathbf{y}}_{\alpha}$$

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

$$\alpha = \mathbf{K}_y^{-1} \mathbf{y} = \underbrace{\mathbf{L}^{-T} \mathbf{L}^{-1} \mathbf{y}}_m$$

$$m = \mathbf{L}^{-1} \mathbf{y}$$

$$\mathbf{L}m = \mathbf{y}$$

$$\mathbf{L}^T m = \alpha$$

Algorithm 15.1: GP regression

1 $\mathbf{L} = \text{cholesky}(\mathbf{K} + \sigma_y^2 \mathbf{I}); \leftarrow$

2 $\alpha = \mathbf{L}^T \setminus (\mathbf{L} \setminus \mathbf{y});$

3 $\mathbb{E}[f_*] = \mathbf{k}_*^T \alpha;$

4 $\mathbf{v} = \mathbf{L} \setminus \mathbf{k}_*;$

5 $\text{var}[f_*] = \kappa(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v};$

6 $\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^T \alpha - \sum_i \log L_{ii} - \frac{N}{2} \log(2\pi)$

$$\mathbf{L}^T \alpha = m$$

Next lecture

In the next lecture, we capitalize on GPs to introduce **active learning**, **Bayesian optimization** and **GP bandits**.

For a recent article on bandits and **Thompson sampling** at work at Google, see:

<http://analytics.blogspot.ca/2013/01/multi-armed-bandit-experiments.html>

For an article on Bayesian optimization, see:

<http://arxiv.org/abs/1012.2599>