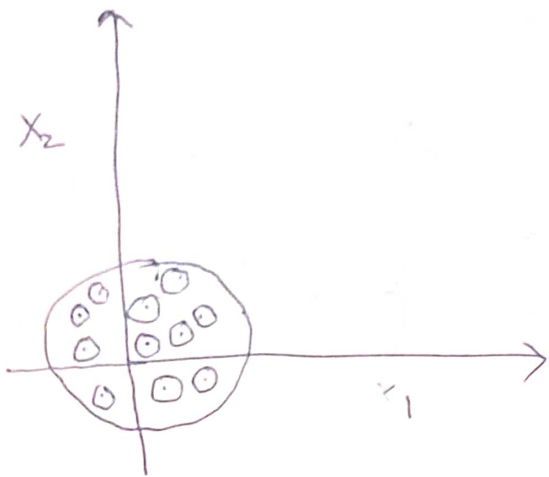


# Gaussian Basics

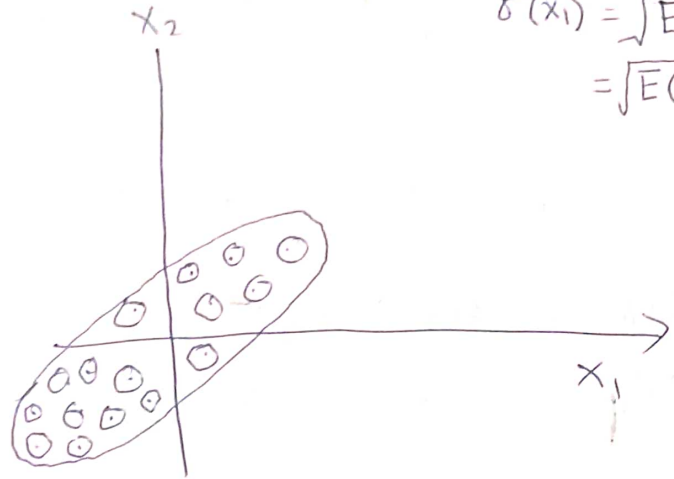
$$\rho_{X_1, X_2} = \frac{\text{Cov}(X_1, X_2)}{\sigma(X_1) \sigma(X_2)}$$

$$\sigma(X_1) = \sqrt{E(X_1 - \mu_1)^2}$$

$$= \sqrt{E(X_1^2)} \text{ if } \mu_1 = 0$$



$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

dot product  
is a measure  
of similarity

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

$E(X_1, X_2) \leftarrow$

$$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \right)$$

$P(X_2 | X_1 = x)$  can be obtained by cutting the  
joint density function.

$P(X_2 | X_1 = x) = P(X_2 | x)$  is also known as  
conditional distribution.

$$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Conditional Mean and Variance of  $p(X_1 | X_2)$

$$\left. \begin{aligned} \mu_{1|2} &= \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (X_2 - \mu_2) \\ \Sigma_{1|2} &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \end{aligned} \right\} p(X_1 | X_2)$$

### Theorem

Suppose  $(X_1, X_2)$  is a jointly Gaussian distribution with parameters

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12} & \Sigma_{22} \end{pmatrix}$$

$$\Lambda = \Sigma^{-1} = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12} & \Lambda_{22} \end{bmatrix}$$

The marginal distribution can be expressed as

$$p(X_1) = N(\mu_1, \Sigma_{11})$$

$$p(X_2) = N(\mu_2, \Sigma_{22})$$

The conditional distribution can be expressed as

$$p(X_1 | X_2) = N(\mu_{1|2}, \Sigma_{1|2}) \quad \text{where}$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (X_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

## Random Sample from multivariate Gaussian

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

then draw  $x_1 \sim N(0, 1)$   
                   $\&$   $x_2 \sim N(0, 1)$

the  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  jointly follow  $N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$

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

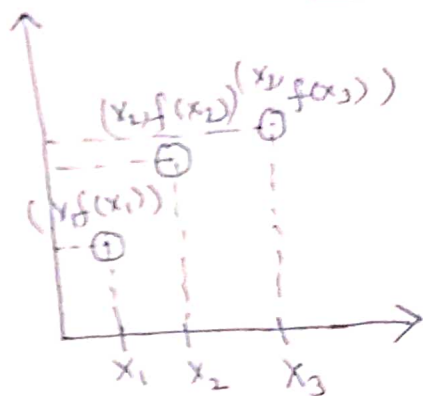
then  $X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + L X^*$

where  $X^* \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$

$\&$   $LL^T = \Sigma$

The matrix  $L$  can be obtained by  
cholesky decomposition.

# Gaussian process



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

Let  $K_{ij}$  can be given by a measure of similarity

$$K_{ij} = e^{-\|x_i - x_j\|^2}$$

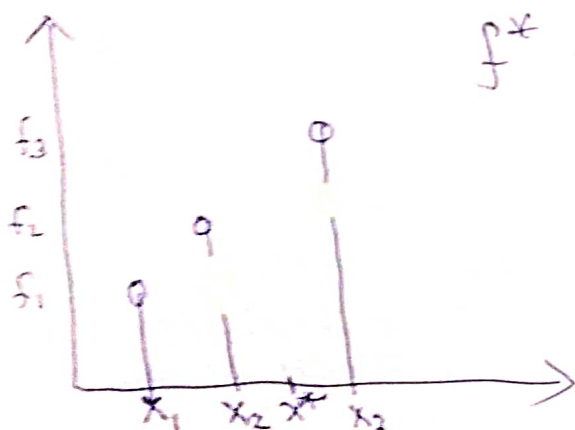
$$= \begin{cases} 0 & \text{if } \|x_i - x_j\| \rightarrow \infty \\ 1 & \text{if } x_i = x_j \end{cases}$$

Given Data

$$\mathcal{D} = \{ (x_1, f_1), (x_2, f_2), (x_3, f_3) \}$$

and  $x^*$

$$f^* = ?$$



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

$$f^* \sim N(0, K(x_* x_*))$$

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

If you increase  $x^*$ ,  $f^*$  should increase but, we assume that  $x^*$  is uncorrelated. So, to implement correlation, we can write

$$\begin{bmatrix} f \\ f_x \end{bmatrix} \sim N \left( 0, \begin{bmatrix} K_{11} & K_{12} & K_{13} & K_{1x} \\ & K_{22} & K_{23} & K_{2x} \\ & & K_{33} & K_{3x} \\ K_{x1} & K_{x2} & K_{x3} & K_{xx} \end{bmatrix} \right)$$

$$K_x = \begin{bmatrix} K_{1x} \\ K_{2x} \\ K_{3x} \end{bmatrix}$$

$$E(f^*) = K_x^T K^{-1} f = \mu^*$$

$$\sigma^* = K_x^T K^{-1} K_x + K_{xx}$$

So, we can estimate the mean and confidence interval.

## Gaussian Process: A distribution over functions

A Gaussian process is a Gaussian distribution over function

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = E[f(x)]$$

$$k(x, x') = E[(f(x) - m(x))(f(x') - m(x'))^T]$$

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

## Simulation from Gaussian Process

(I) Create  $x_{1:N}$

(II) Create  $\mu = 0_N$  &  $K_{N \times N} = \left(\exp\left(-\frac{1}{2}\|x_i - x_j\|^2\right)\right)_{i,j}$

(III) Decompose  $K = LL^T$  by Cholesky Method  $N \times N$

(III) Generate  $f_N^* \sim N(0_N, I_{N \times N})$

(IV)  $f = L \cdot f_N^*$

## Gaussian Process Prior and Posterior

Let  $D = \{(x_i, f_i), i=1, 2, \dots, N\}$

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

$$P(f) = N(\mu(x), K(x, x'))$$

## Noiseless Gaussian Process Regression

We observe a training set  $D = \{(x_i, f_i), i=1, \dots, n\}$  where  $f_i = f(x_i)$ . Given a test set  $X_*$  of size  $N_* \times D$ , we want to predict the outputs  $f_*$

$$\begin{bmatrix} f \\ f_* \end{bmatrix} \sim N \left( \begin{bmatrix} \mu \\ \mu_* \end{bmatrix}, \begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix} \right)$$

$$K_* = K(X, X^*)$$

$$K_{**} = K(X_* X_*)$$

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

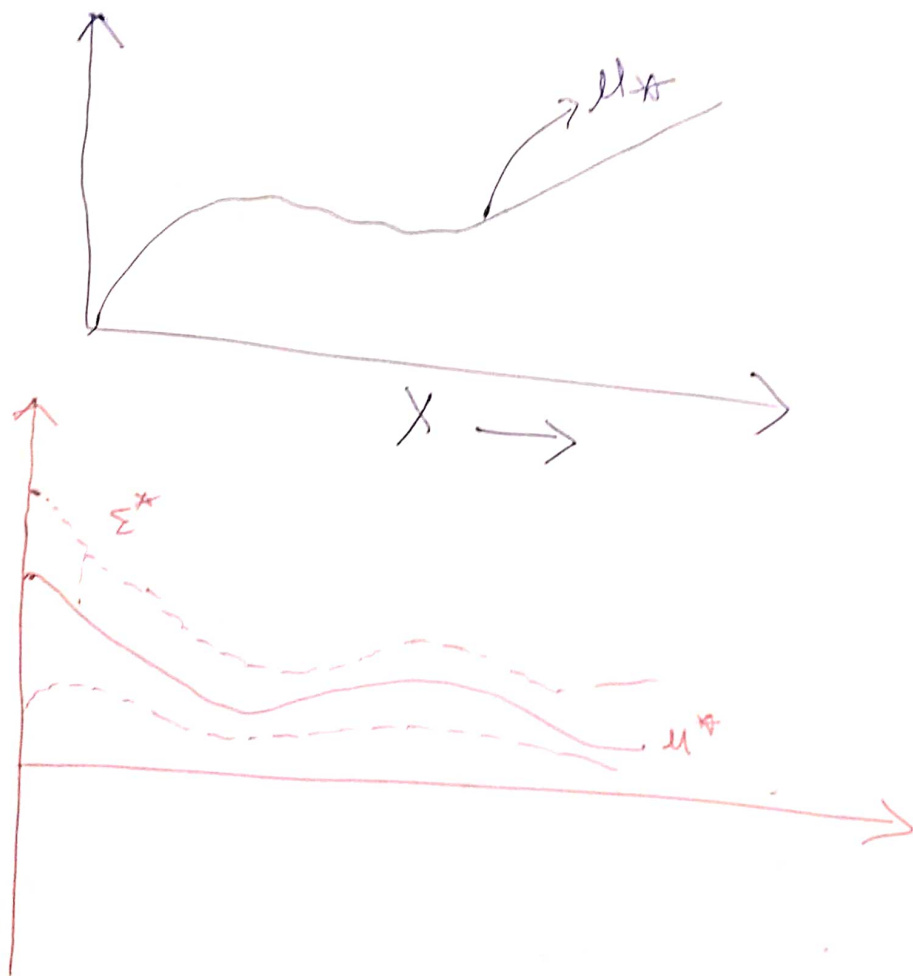


Noiseless Gaussian Process Regression can be estimated as

$$p(f_* | X_*, X, f) = N(\mu_*, \Sigma_*)$$

$$\mu_* = \mu(X_*) + K_*^T K^{-1} (f - \mu(X))$$

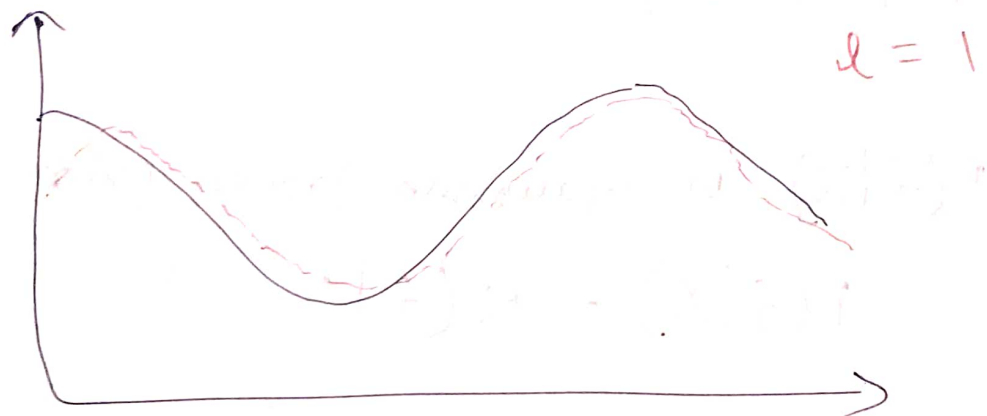
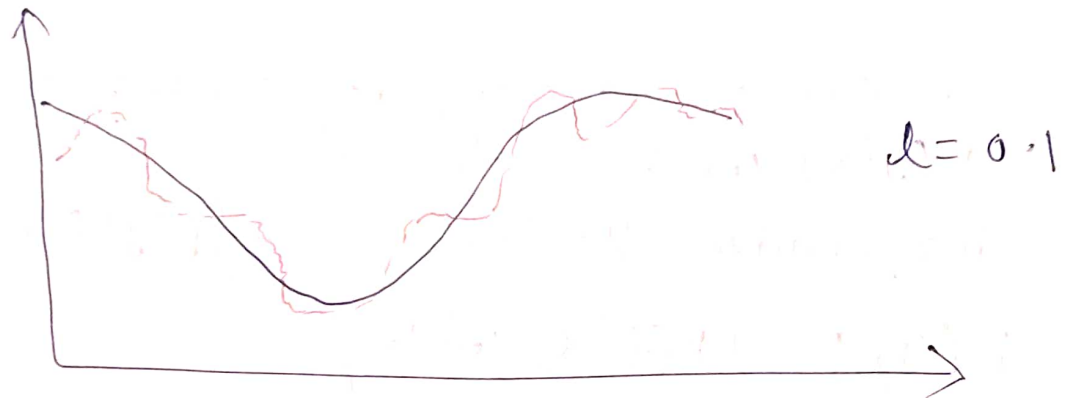
$$\Sigma_* = K_{**} - K_*^T K^{-1} K_*$$



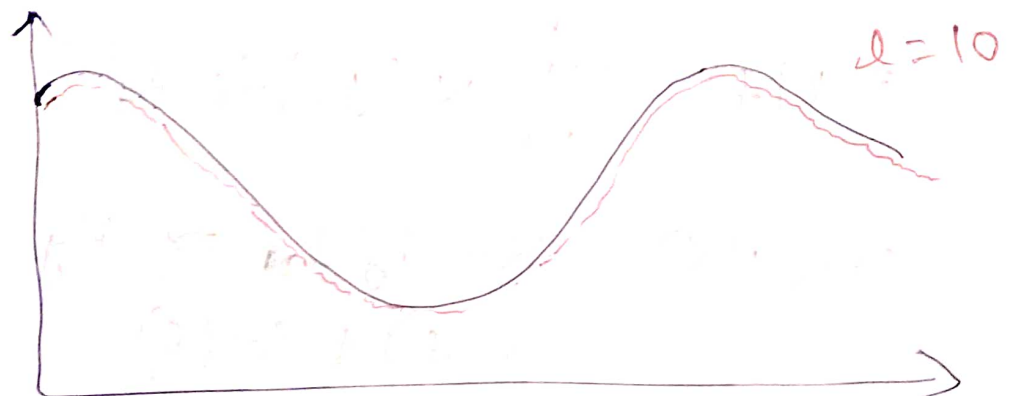


## Effect of Choosing Kernel Width parameter

If  $\lambda^2$  is very small, means Kernel is very thin only points that are nearby similar to each other and we can expect more wiggly function.



As we increase the kernel width the function will become smooth.



One way of choose kernel width is cross-validation.  
The blue curve is true function.

# Noisy Gaussian Process Regression

$$y = f(x) + \epsilon$$

$$\text{where } \epsilon \sim N(0, \sigma_y^2)$$

$y$  is called noisy function measurement

In order to compute  $p(y|x)$  we need to marginalize the joint  $p(y|f, x) p(f|x)$  with respect to  $f$ .

$$p(y|x) = \int p(y|f, x) p(f|x) df$$

$p(f|x)$  is Gaussian process prior

$$p(f|x) = N(f|0, k)$$

Similarity Kernel

$$p(y|f) = \prod_{i=1}^n N(y_i | f_i, \sigma_y^2)$$

$$\begin{aligned} \text{cov}(y|x) &= K + \sigma_y^2 I_N \approx K_y \\ &= \text{cov}(f) + \text{cov}(\epsilon) \end{aligned}$$

$$\begin{bmatrix} y \\ f_* \end{bmatrix} \sim N \left( 0, \begin{bmatrix} K_y & K_* \\ K_*^T & K_{**} \end{bmatrix} \right)$$

$$p(f_* | X_*, X, y) = N(f_* | \mu_*, \Sigma_*)$$

$$\mu_* = K_*^T K_y^{-1} y$$

$$\Sigma_* = K_{**} - K_*^T K_y^{-1} K_*$$

If you want to fit Gaussian process to data you just construct a matrix  $K$  using the similarity kernel (If you have noise, you just need to add  $\sigma$ ) and obtain  $\mu_*$  and  $\Sigma_*$ . This method is sometimes called ML-II method.