



**Hochschule  
Bonn-Rhein-Sieg**  
University of Applied Sciences



# Gaussian Process

Manoj Kolpe Lingappa

# Introduction

- In parametric model  $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$  to explain the data and find optimal value of  $\boldsymbol{\theta}$  using MLE.
- Posterior can be obtained using the formula  $p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y})$
- Gaussian processes is a non parametric method.
- Gaussian process can be used to infer a distribution over function directly.
- Gaussian process define a prior over function and once some new data points are observed, a posterior can be inferred.
- Inference of continuous function values in this context is known as GP regression.

## Gaussian process

- A Gaussian process is a random process where any points  $\mathbf{x} \in \mathbb{R}^d$  is assigned a random variable  $f(\mathbf{x})$  where the joint distribution of a finite number of these variables  $p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))$  is itself Gaussian:

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

- $\mathbf{f} = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))$ ,  $\boldsymbol{\mu} = (\mu(\mathbf{x}_1), \dots, \mu(\mathbf{x}_n))$  and  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  is the kernel, it is common to use  $\mu(\mathbf{x}) = 0$
- Gaussian process is a distribution over function and shape of distribution is defined by  $\mathbf{K}$ .

# Introduction

- Let's consider a noise free function  $\mathbf{f}$  for input  $\mathbf{X}$ .
- A GP prior can be converted into a GP posterior  $p(\mathbf{f}_*|\mathbf{X}_*,\mathbf{X},\mathbf{f})$  which can be used to make prediction  $\mathbf{f}_*$  at new inputs  $\mathbf{X}_*$
- By definition of marginalization, the joint distribution of observed values  $\mathbf{f}$  and predictions  $\mathbf{f}_*$  is again a Gaussian which can be partitioned into

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

- where  $\mathbf{K}_* = \kappa(\mathbf{X}, \mathbf{X}_*)$  and  $\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$
- With  $N$  training data and  $N_*$  new input data  $\mathbf{K}$  is a  $N \times N$  matrix,  $\mathbf{K}_*$  a  $N \times N_*$  matrix and  $\mathbf{K}_{**}$  a  $N_* \times N_*$  matrix.
- The predicted distribution is given by,

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

# Introduction

- where  $K_y = K + \sigma_y^2 \mathbf{I}$
- To additionally include noise  $\epsilon$  into predictions  $y_*$  we have to add  $\sigma_y^2$  to the diagonal of  $\Sigma_*$

$$p(\mathbf{y}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{y}_* | \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_* + \sigma_y^2 \mathbf{I})$$

- Let's take a RBF kernel

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

⊗ RBF kernels



The **squared exponential kernel** (SE kernel) or **Gaussian kernel** is defined by

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \mathbf{x}')\right)$$

If  $\boldsymbol{\Sigma}$  is diagonal, this can be written as

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2} \sum_{j=1}^D \frac{1}{\sigma_j^2} (x_j - x'_j)^2\right)$$

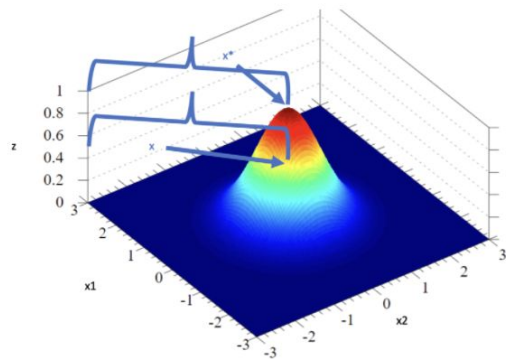
$$k(x, x^*) = \exp\left(-\frac{\|x - x^*\|^2}{2\sigma^2}\right)$$

If you take the function apart, you'll see that part of the function is a measure of the squared distance between  $x$  and  $x^*$

$$\|x - x^*\|^2 \qquad k(x, x^*) = \exp\left(-\frac{\sqrt{2}}{2 \times 2^2}\right) = 0.84$$

For example, our vectors  $x$  and  $x^*$  would approximately look like this on our Gaussian kernel.

$$x^* = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad x = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$



First, we have to calculate the Euclidian distance between the two vectors.

$$\|x - x^*\|^2 = \sqrt{(-1 - 0)^2 + (-1 - 0)^2} = \sqrt{2}$$

Then we plug this result into the Gaussian kernel, assuming that the standard deviation  $\sigma$  equals 1.

$$k(x, x^*) = \exp\left(-\frac{\sqrt{2}}{2 \times 1^2}\right) = 0.49$$

The similarity is pretty close to 0.5