

Optimizing Gaussian Processes

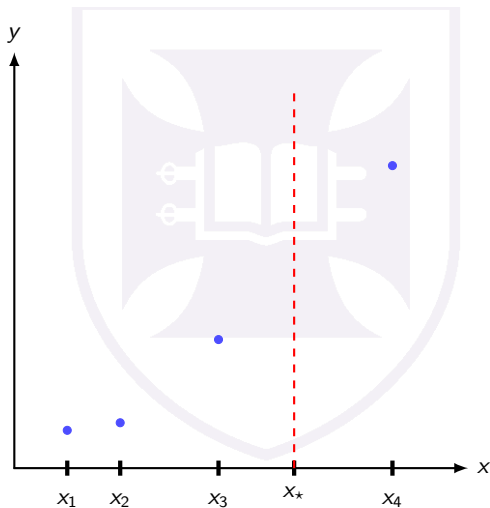
MSS Presentation

Michael Ciccotosto-Camp - 44302913

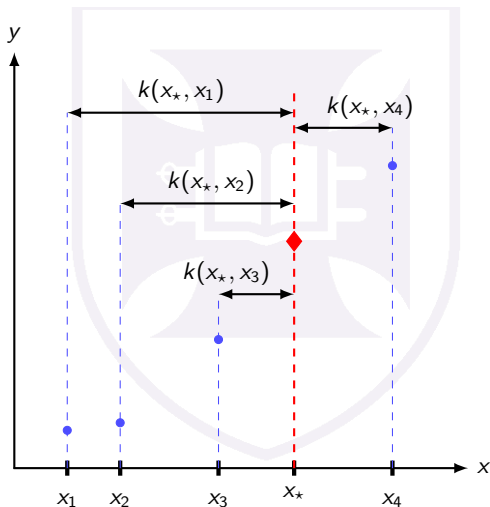


THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

Time Series Prediction

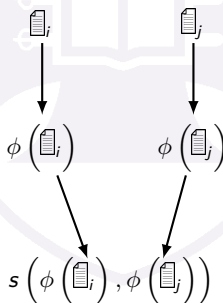


Time Series Prediction



The Kernel Trick

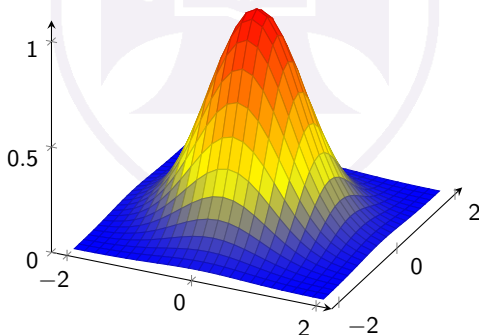
- Q: How do we get a suitable function k for computing similarity? A: Use the kernel trick!
- Suppose we have some inputs $\left[\mathbb{I}_1, \dots, \mathbb{I}_n\right]$ (with their corresponding experimental observations $[y_1, \dots, y_n]$), where \mathbb{I}_i can take a number of different of form (perhaps a tree data structure or vectors of values).



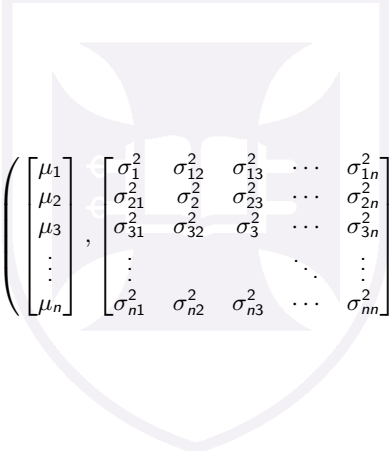
- The function s provides us with some notion of similarity between inputs after they've been "transformed" into a nicer form using a *feature map* ϕ .

The Kernel Trick

- The kernel function k does all this computation in one step so that $k(\mathbf{x}_i, \mathbf{x}_j) = s(\phi(\mathbf{x}_i), \phi(\mathbf{x}_j))$.
- Usually we have access to k , meaning we can avoid having to construct a feature map ϕ and similarity function s .
- A very common kernel function used is the RBF or Gaussian kernel.

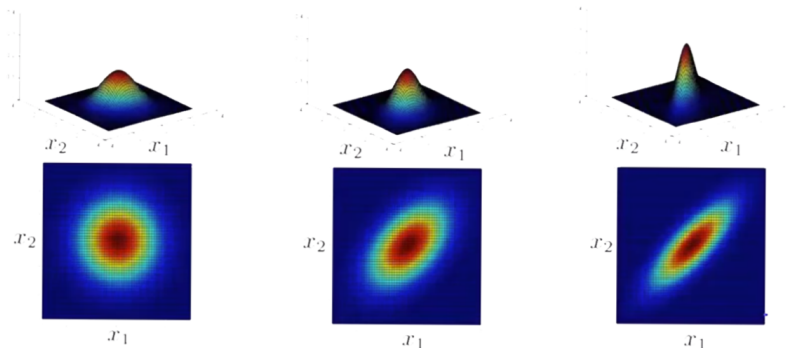


Multi-Variate Gaussian


$$\mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \vdots \\ \mu_n \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 & \sigma_{13}^2 & \cdots & \sigma_{1n}^2 \\ \sigma_{21}^2 & \sigma_2^2 & \sigma_{23}^2 & \cdots & \sigma_{2n}^2 \\ \sigma_{31}^2 & \sigma_{32}^2 & \sigma_3^2 & \cdots & \sigma_{3n}^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1}^2 & \sigma_{n2}^2 & \sigma_{n3}^2 & \cdots & \sigma_{nn}^2 \end{bmatrix} \right)$$

Predictions

- How do we use our data to make predictions with our kernel function?
- Within the Gaussian Process paradigm we assume that our data along with the novel point at which we would like to predict form a joint Gaussian distribution.



(Machine Learning, Stanford University, <https://www.coursera.org/learn/machine-learning>)

$$\mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \vdots \\ \mu_n \end{bmatrix}, \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \sigma_{13}^2 & \cdots & \sigma_{1n}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 & \sigma_{23}^2 & \cdots & \sigma_{2n}^2 \\ \sigma_{31}^2 & \sigma_{32}^2 & \sigma_{33}^2 & \cdots & \sigma_{3n}^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1}^2 & \sigma_{n2}^2 & \sigma_{n3}^2 & \cdots & \sigma_{nn}^2 \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \vdots \\ \mu_n \end{bmatrix}, \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \sigma_{13}^2 & \cdots & \sigma_{1n}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 & \sigma_{23}^2 & \cdots & \sigma_{2n}^2 \\ \sigma_{31}^2 & \sigma_{32}^2 & \sigma_{33}^2 & \cdots & \sigma_{3n}^2 \\ \vdots & & & \ddots & \vdots \\ \sigma_{n1}^2 & \sigma_{n2}^2 & \sigma_{n3}^2 & \cdots & \sigma_{nn}^2 \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ ? \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) & k(\mathbf{x}_1, \mathbf{x}_*) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) & k(\mathbf{x}_2, \mathbf{x}_*) \\ \vdots & & \ddots & & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) & k(\mathbf{x}_n, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{x}_1) & k(\mathbf{x}_*, \mathbf{x}_2) & \cdots & k(\mathbf{x}_*, \mathbf{x}_n) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

- How do we compute that missing mean within our MVN from the previous slide?



- How do we compute that missing mean within our MVN from the previous slide?
- We can use the following theorem: (Marginals and conditionals of an MVN), suppose $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2]$ is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$$

then the posterior conditional is given by

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

$$\mathcal{N} \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) & k(\mathbf{x}_1, \mathbf{x}_*) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) & k(\mathbf{x}_2, \mathbf{x}_*) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) & k(\mathbf{x}_n, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{x}_1) & k(\mathbf{x}_*, \mathbf{x}_2) & \cdots & k(\mathbf{x}_*, \mathbf{x}_n) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) & k(\mathbf{x}_1, \mathbf{x}_*) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) & k(\mathbf{x}_2, \mathbf{x}_*) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) & k(\mathbf{x}_n, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{x}_1) & k(\mathbf{x}_*, \mathbf{x}_2) & \cdots & k(\mathbf{x}_*, \mathbf{x}_n) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{xx}} & \mathbf{K}_{\mathbf{x}_* \mathbf{x}}^\top \\ \mathbf{K}_{\mathbf{x}_* \mathbf{x}} & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) & k(\mathbf{x}_1, \mathbf{x}_*) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) & k(\mathbf{x}_2, \mathbf{x}_*) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) & k(\mathbf{x}_n, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{x}_1) & k(\mathbf{x}_*, \mathbf{x}_2) & \cdots & k(\mathbf{x}_*, \mathbf{x}_n) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{X}\mathbf{X}} & \mathbf{K}_{\mathbf{X}_* \mathbf{X}}^\top \\ \mathbf{K}_{\mathbf{X}_* \mathbf{X}} & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

- The mean and covariance can then be computed using the theorem from before as

$$y_* = \mathbf{K}_{\mathbf{X}_* \mathbf{X}} \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{y}$$

$$\mathbb{V}[y_*] = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{K}_{\mathbf{X}_* \mathbf{X}} \mathbf{K}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{K}_{\mathbf{X}_* \mathbf{X}}^\top.$$

$$\mathcal{N} \left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \\ y_* \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) & k(\mathbf{x}_1, \mathbf{x}_*) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) & k(\mathbf{x}_2, \mathbf{x}_*) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) & k(\mathbf{x}_n, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{x}_1) & k(\mathbf{x}_*, \mathbf{x}_2) & \cdots & k(\mathbf{x}_*, \mathbf{x}_n) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$\mathcal{N} \left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{XX}} & \mathbf{K}_{\mathbf{X}_* \mathbf{X}}^\top \\ \mathbf{K}_{\mathbf{X}_* \mathbf{X}} & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

- The mean and covariance can then be computed using the theorem from before as

$$y_* = \mathbf{K}_{\mathbf{X}_* \mathbf{X}} \mathbf{K}_{\mathbf{XX}}^{-1} \mathbf{y}$$

$$\mathbb{V}[y_*] = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{K}_{\mathbf{X}_* \mathbf{X}} \mathbf{K}_{\mathbf{XX}}^{-1} \mathbf{K}_{\mathbf{X}_* \mathbf{X}}^\top$$

- Another way of looking at the prediction is seeing it as a linear combination of n kernel evaluations centered at the input \mathbf{x}_*

$$y_* = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_*)$$

Unoptimized GPR

Algorithm 1: Unoptimized GPR

input : Observations \mathbf{X}, \mathbf{y} and a test input \mathbf{x}_* .

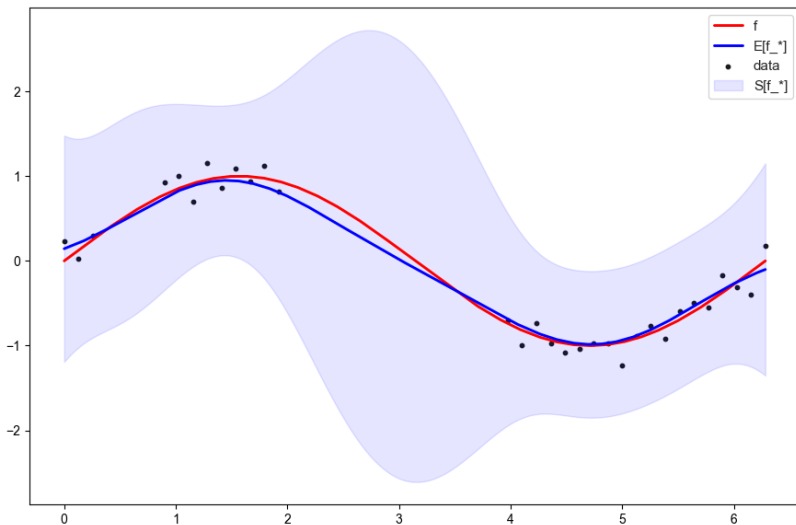
output: A prediction y_* with its corresponding variance $\mathbb{V}[y_*]$.

- 1 $\mathbf{L} = \text{cholesky}(\mathbf{K}_{\mathbf{X}\mathbf{X}})$
 - 2 $\boldsymbol{\alpha} = \text{lin-solve}(\mathbf{L}^\top, \text{lin-solve}(\mathbf{L}, \mathbf{y}))$
 - 3 $y_* = \mathbf{K}_{\mathbf{x}_* \mathbf{X}} \boldsymbol{\alpha}$
 - 4 $\mathbf{v} = \text{lin-solve}(\mathbf{L}, \mathbf{K}_{\mathbf{x}_* \mathbf{X}})$
 - 5 $\mathbb{V}[y_*] = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$
 - 6 **return** $y_*, \mathbb{V}[y_*]$
-

Implementation

```
def gp_reg_pred(X_train, Y_train, x_pred, sigma):  
    n, d = X_train.shape  
    # Create the Gram matrix corresponding to the training data set.  
    K = exact_kernel(X_train, sigma=sigma)  
    # Noise variance of labels.  
    s = np.var(Y_train.squeeze())  
    L = np.linalg.cholesky(K + s*np.eye(n))  
    # Compute the mean at our test points.  
    Lk = np.linalg.solve(L, exact_kernel(X_train, x_pred, sigma=sigma))  
    Ef = np.dot(Lk.T, np.linalg.solve(L, Y_train))  
    # Compute the variance at our test points.  
    K_ = exact_kernel(x_pred, sigma=sigma)  
    Vf = np.diag(K_) - np.sum(Lk**2, axis=0)  
    return Ef, Vf
```

sin Function Prediction with Added Noise



Stock Market Prediction

