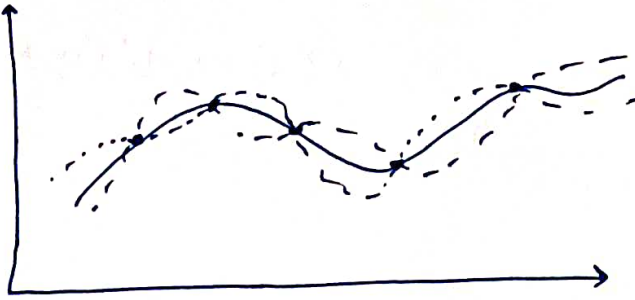# Gaussian Process Regression

## Objective



There can be many functions which pass through these points.

## Distribution over functions

$$\mu = K^* \cdot K^{-1} \cdot y \qquad\qquad V = K^{**} - K^* \cdot K^{-1} K^{*T}$$

## ① Kernel First

$$y = f(x) \quad \leftarrow \text{Noiseless.}$$

$$f(x) \;\, \sim\; GP\Big(m(x),\; K(x,x')\Big), \qquad m(x)=0$$

$$K(x,x') = e^{-\frac{1}{2}\|x-x'\|^2} \qquad \begin{array}{l}\text{If } x=x' \;\; K(x,x') = 1 \\ \quad x-x' \to \pm \infty \;\; K(x,x') \to 0\end{array}$$

$$K(X,X) = \begin{bmatrix} K(x_1,x_1) & \cdots & K(x_1,x_n) \\ \vdots & \ddots & \vdots \\ K(x_N,x_1) & \cdots & K(x_N,x_N) \end{bmatrix} = \begin{bmatrix} 1 & \cdots & K(x_1,x_N) \\ \vdots & \ddots & \vdots \\ K(x_N,x_1) & \cdots & 1 \end{bmatrix}$$

Prior to seeing any data = GP "prior"

$$\begin{pmatrix} f(x) \\ f(x^*) \end{pmatrix} \;\sim\; \mathcal{N}\left( 0, \begin{bmatrix} K[X,X] & K[X,X^*] \\ K[X^*,X] & K[X^*,X^*] \end{bmatrix} \right)$$

$$\text{If } \begin{pmatrix} x \\ y \end{pmatrix} \sim \mathcal{N}\left( \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right) \quad \text{then}$$

$$y|x \sim \mathcal{N}\Big( \mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(x-\mu_x),\; \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy} \Big)$$

$.1 \quad \{(x^*)| f(x) \sim \mathcal{N}\left( K(x^*, x) K(x, x)^{-1} f(x), \quad K(x^*, x^*) - \right.$

$$K(x^*, x) K(x, x)^{-1} K(x, x^*) \Big)$$

If we introduce noise,

$$y = f(x) + \epsilon \qquad \epsilon \sim \mathcal{N}(0, \sigma^2) \qquad f(x) \sim \mathcal{N}(0, K)$$

$$\therefore y \sim \mathcal{N}(0, K + \sigma^2 I)$$

$$\begin{pmatrix} y \\ f(x^*) \end{pmatrix} \sim \mathcal{N}\left( 0, \begin{bmatrix} K(x,x) + \sigma^2 I & K(x, x^*) \\ K(x^*, x) & K(x^*, x^*) \end{bmatrix} \right)$$

$$f(x^*) | y \sim \mathcal{N}(\bar{f}, \bar{V})$$

$$\bar{f} = K(x^*, x)\left( K(x, x) - \sigma^2 I \right)^{-1} y$$

$$\bar{V} = K(x^*, x^*) - K(x^*, x)\left( K(x, x) + \sigma^2 I \right)^{-1} K(x, x^*)$$

② <u>Prior</u>

$$y = f(x) + \epsilon \qquad f(x) = x^T w \qquad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$L(y | x, w) = \prod P(\gamma_i = y_i) = (2\pi \sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2} \|y - xw\|^2} = \mathcal{N}(x_w, \ldots)$$

Bayesian Prior on $w$ : $w \sim \mathcal{N}(0, \Sigma_p)$

Posterior : $w|y \propto L(y | x, w) \cdot \pi(w) = \ldots\ldots = \mathcal{N}(\cdot m, A^{-1})$

where $\quad$ A $= \frac{1}{\sigma^2} x^T x + \Sigma_p^{-1}, \quad m = \frac{1}{\sigma^2} A^{-1} x^T y$

$$f(x^*) = x^{*T} w \sim \mathcal{N}\left( x^{*T} m, \; x^{*T} A^{-1} x^* \right)$$

$$f(X^*) = X^{*T} w \sim \mathcal{N}\left( X^{*T} m, \; X^{*T} A^{-1} X^* \right)$$

## ③ MSE

We want the best "linear" estimator — linear in terms of $Y$

$$\hat{y}(x) = \sum \lambda_i y_i = y^T \lambda \qquad y^* - \wedge y$$

$$\min_\lambda \; \mathbb{E}\left[\hat{y}(x) - y(x)\right]^2 = \mathbb{E}\left[\lambda^T y \, y^T \lambda - 2\lambda^T y \cdot y^* - y^{*2}\right]$$

$$\text{"} \; \mathbb{E}\left[\hat{y}(x) - y(x^T)\right]^2 \text{"}$$

$$= \lambda^T K \lambda - 2\lambda^T K^* - \tau^2$$

$$\frac{\partial}{\partial \lambda} = 2K\lambda - 2K^* = 0 \quad \Rightarrow \quad \lambda = K^{-1} K^*$$

$$\hat{y}(x) = y^T K^{-1} K^* = K^{*T} K^{-1} y$$

$$\wedge y \Rightarrow (y^* - \wedge y)^T (y^* - \wedge y) = y^{*T} y^* - 2 y^{*T} \wedge y + y^T \wedge^T \wedge y$$

$$\frac{\partial}{\partial \wedge} = 2(y^* - \wedge y) y^T = 0 \Rightarrow \wedge (y y^T) = y^* y^T \Rightarrow \wedge K = K^*$$

$$\Rightarrow \wedge = \cancel{K^0} K^* K^{-1}$$

$$\hat{y} = K^* K^{-1} y$$

## SVD:

In linear Regression, we know that $\hat{\beta} = (X^T X)^{-1} X^T y$

Let's decompose $X$ via SVD: $X = U S V^T$

$$\hat{\beta} = (V S^T U^T U S V^T)^{-1} V S^T U^T y = (V S^2 V^T)^{-1} V S^T U^T y$$

$$= (V^T)^{-1} (S^2)^{-1} V^{-1} V \; S^T U^T y$$

$$= V (S^2)^{-1} V^T V \; S^T U^T y$$

$$= V \underbrace{S^T U^T . U S (S^2)^{-1}}_{} (S^2)^{-1} S^T U^T y$$

$$= X^T B$$

$$\hat{y}(X) = X \hat{\beta} = X X^T B$$

What if linear function is too limited?
Project to a higher dimensional feature space $\phi(x)$

$$f(X^*) = \phi(X^*) \cdot W \sim \mathcal{N}\left(\phi(X^*) \, m, \; \phi(X^*) A^{-1} \phi(x^*)^T\right)$$

Matrix $A$ is $p \times p$ matrix, it will be difficult to invert if $p \to \infty$.

Let's define:
$$K(X, X) = X \Sigma_p X^T = K$$

$$\frac{1}{\sigma^2} X^T (K + \sigma^2 I) = \frac{1}{\sigma^2} X^T (X \Sigma_p X^T + \sigma^2 I) = A \Sigma_p X^T$$

$$A^{-1} \frac{1}{\sigma^2} X^T (K + \sigma^2 I)(K + \sigma^2 I)^{-1} = A^{-1} A \Sigma_p X^T (K + \sigma^2 I)^{-1}$$

$$\boxed{\frac{1}{\sigma^2} A^{-1} X^T = \Sigma_p X^T (K + \sigma^2 I)^{-1}} \quad -$$

$$\mathbb{E}(f(x^*)) = X^* m = X^* \frac{1}{\sigma^2} A^{-1} X^T y = X^* \Sigma_p X^T (K + \sigma^2 I)^{-1} y$$

$$X^* \Sigma_p X^T = K[X^*, X)$$

$$\therefore \mathbb{E}(f(x^*)) = K[X^*, X](K(X, X) + \sigma^2 I)^{-1} y$$

$$V(f(x^*)) = K(X^*, X^*) - K[X^*, X)(K(X, X) + \sigma^2 I)^{-1} K(X, X^*)$$

$$\left(Z + U W V^T\right)^{-1} = Z^{-1} - Z^{-1} U \left(W^{-1} + V^T Z^{-1} U\right)^{-1} V^T Z^{-1}$$

↗ Matrix Inversion Lemma

Now, if we replace $Z = \Sigma_p$ $\qquad W^{-1} = \sigma^2 I \qquad U = V = X$

**Minimal MSE**

$$\left(y - XX^T B\right)^T \left(y - XX^T B\right) = y^T y - 2y^T X X^T B + B^T X X^T X X^T B.$$

$$\frac{\partial}{\partial B} = -2y^T X X^T + 2 X X^T X X^T B = 0 \qquad K := X X^T$$

$$\Rightarrow K \cdot K \cdot B = K y \qquad \Rightarrow \qquad B = K^{-1} y$$

$$\hat{y}(x^*) = X^* \hat{B} = X^* X^T K^{-1} y \qquad K^* = X^* X^T \qquad \hat{y}(x^*) = K^* K^{-1} y$$