

Kernel Regression

Recall

Binary classification: $\hat{y} = \text{sign}(x^T w)$

Linear regression: $\hat{y} = x^T w$

Kernels (in Machine Learning)

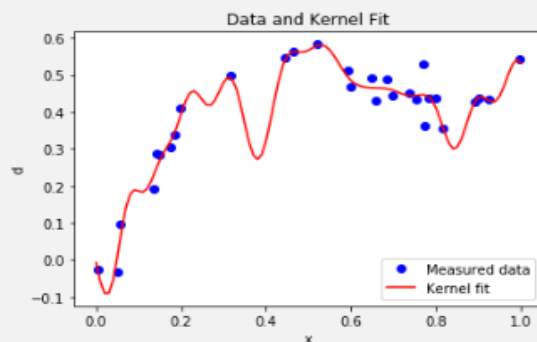
w depends on $x_1, y_1, x_2, y_2, \dots$

Kernel

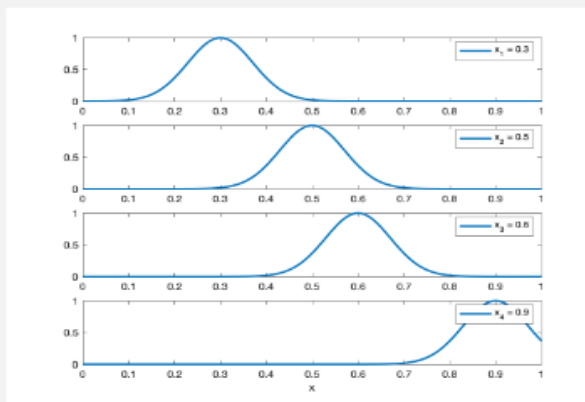
$$\hat{y} = \phi(x)^T w \rightarrow \hat{y} = \sum_i \alpha_i K(x, x_i)$$

High dimensional

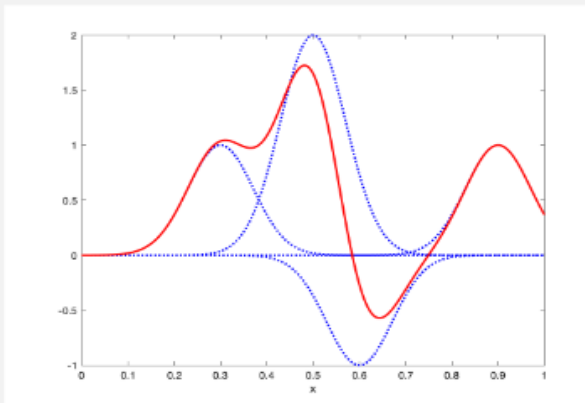
weighted sum of similarities between feature vector and each training point



$$K(x, x_i) = e^{-\frac{(x-x_i)^2}{0.01}}$$



$$\hat{y}(x) = K(x, x_1) + 2K(x, x_2) - K(x, x_3) + K(x, x_4)$$



Prediction: $\hat{y} = \sum_i \alpha_i K(x, x_i)$

How do we find good α_i ?

start by finding w using ridge regression

$$w^* = \arg \min_w \|\Phi w - y\| + \lambda \|w\|^2$$

$$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

↓ Activity 23, problem 2

$$w^* = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y$$

$$\hat{y} = \phi(x)^T w^* \rightarrow \hat{y} = \sum_i \alpha_i K(x, x_i)$$

$$\alpha = (\Phi \Phi^T + \lambda I)^{-1} y$$

where $\Phi \Phi^T$ has i, j entry $K(x_i, x_j)$

No need to compute $\phi(\cdot)$ to compute $K(\cdot, \cdot)$ or \hat{y} !

Kernel can have a few hyper-parameters.

Danger! → Overfitting

Use cross-validation to choose params.