

# Kernels for classification

Kernel regression (last class):

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

High-dimensional feature transformation

Kernel: measures similarity between  $x$  and  $x_i$ .

Kernels for binary classification (today):

Classification, after feature map:

$$\hat{y} = \text{sign}(\phi(x)^T w) \quad (1)$$

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

Kernel methods – re-write above as:

$$\hat{y} = \text{sign} \left( \sum_i \alpha_i K(x, x_i) \right) \quad (2)$$

weighted sum of similarities between feature vector and each training point

Representer Theorem: (1) and (2) are the same, when

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

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

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

where  $K$  has  $\ell, m$  entry  $K(x_\ell, x_m)$



[Kimeldorf 1970]

Example of kernel classification

How do we predict class of  $x$ ?

$$x = \begin{bmatrix} 0.3 \\ 0.1 \end{bmatrix}$$

$$K(x, x_i) = \exp(-\|x - x_i\|^2)$$

$$\hat{y} = \text{sign} \left( \sum_i \alpha_i \exp \left( -\left\| \begin{bmatrix} 0.3 \\ 0.1 \end{bmatrix} - x_i \right\|^2 \right) \right)$$

