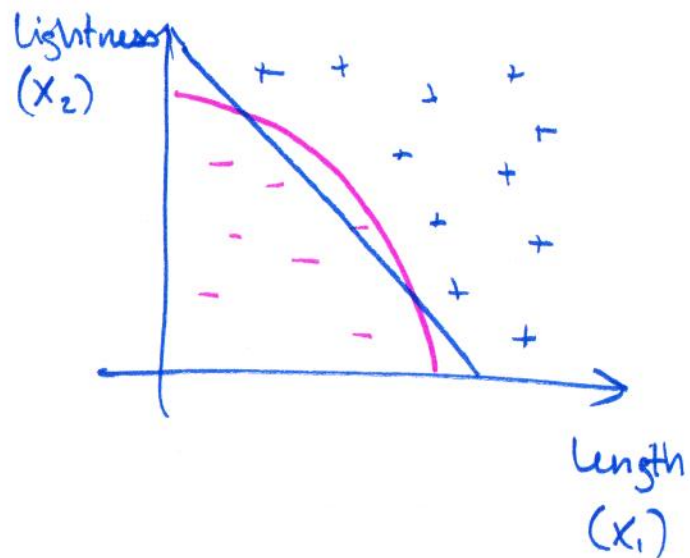


Lecture 25: Kernel Regression



$$\hat{y} = w_1 \cdot x_1 + w_2 \cdot x_2 + \cancel{w_0} = \underline{w}^T \underline{x}$$

$$\tilde{y} = \text{sign}(\hat{y})$$

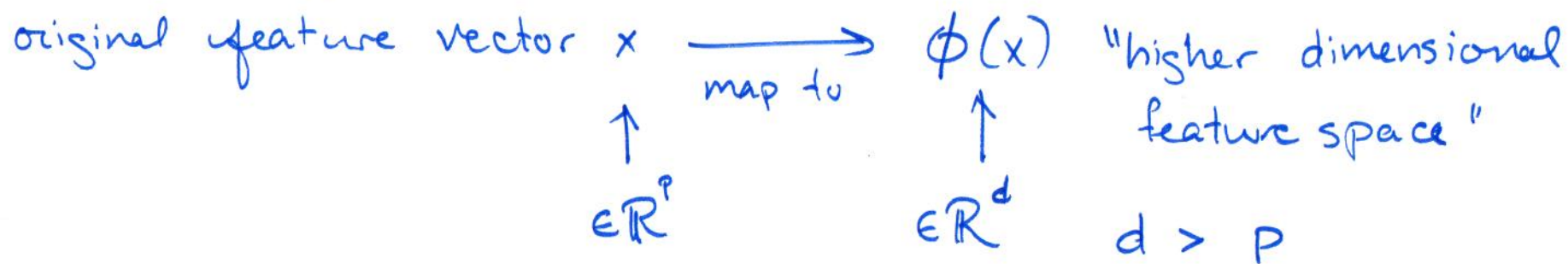
- or -

$$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2$$

$$\Rightarrow \begin{matrix} \tilde{x} \\ \phi(x) \end{matrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix} \Rightarrow \hat{y} = \underline{w}^T \tilde{x}$$

Kernel methods formalize this idea.

(2)



goal: find weights s.t. $\hat{y} = w^T \phi(x)$ gives good predictions. \otimes

when ϕ corresponds to quadratic monomials of p -dim vectors, then $d = O(p^2)$

kernel methods allow us to do \otimes without computing $\phi(x_i)$ explicitly for all training samples x_i .

Recall Ridge Regression

$$\min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

let $X = U\Sigma V^T$ (full SVD) $\Rightarrow U^T U = I, VV^T = V^T V = I$
 $\in \mathbb{R}^{n \times p}$

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\begin{aligned} \hat{w} &= (X^T X + \lambda I)^{-1} X^T y \\ &= (V \Sigma^T \Sigma V^T + \lambda V V^T)^{-1} V \Sigma^T U^T y \\ &= V (\Sigma^T \Sigma + \lambda I)^{-1} V^T V \Sigma^T U^T y \\ &= \underline{V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T U^T y} \end{aligned}$$

$\Rightarrow \hat{w}$ = weighted sum of columns of V
 = weighted sum of rows of $X \Rightarrow \hat{w} = X^T \alpha$ for some $\alpha \in \mathbb{R}^n$

First, note

$$(\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T = \Sigma^T (\Sigma \Sigma^T + \lambda I)^{-1}$$

Then $\hat{w} = V \Sigma^T (\Sigma \Sigma^T + \lambda I)^{-1} U^T y$

λ
 $U^T U = I$

$$= V \underbrace{\Sigma^T U^T U}_{= X^T} (\Sigma \Sigma^T + \lambda I)^{-1} \underbrace{U^T y}_{= \alpha}$$
$$= X^T \underbrace{U}_{\rightarrow} (\Sigma \Sigma^T + \lambda I)^{-1} \underbrace{U^T y}_{\leftarrow} = y$$
$$= X^T (U \Sigma \Sigma^T U^T + \lambda I)^{-1} y$$
$$= X^T (X X^T + \lambda I)^{-1} y$$

$$\Rightarrow \hat{w} = X^T \alpha, \text{ where } \alpha = (X X^T + \lambda I)^{-1} y$$

④

$$\begin{aligned} \hat{y} &= X^T \hat{\Sigma} \\ &= X^T X^T \alpha \\ &= \sum_{i=1}^n \alpha_i \underline{X^T X_i} \end{aligned}$$

Let $K = XX^T$ (linear)
(kernel matrix)

$$K_{ij} = x_i^T x_j$$

$$\Rightarrow \hat{w} = (K + \lambda I)^{-1} y$$

Now work in "feature space"

Let ϕ be a nonlinear ~~linear~~ function
mapping from \mathbb{R}^p to \mathbb{R}^d

define matrix $\Phi = [\phi(x_1) \quad \phi(x_2) \quad \dots \quad \phi(x_n)]$

e.g. if $\phi(x) = x$, then $\Phi = X$

e.g. if $\phi\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ x_2^2 \\ x_2 x_3 \\ x_3^2 \end{bmatrix}$

Regression in kernel space:

$$\hat{w} = \underset{w}{\operatorname{argmin}} \|y - \Phi w\|_2^2 + \lambda \|w\|_2^2$$

$$\boxed{\hat{w}} = \Phi^T \alpha, \quad \alpha = (\Phi \Phi^T + \lambda I)^{-1} y$$

define $K = \Phi \Phi^T$, $K_{ij} = \underline{\phi(x_i)^T \phi(x_j)}$

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

$$K = (X X^T)^{\odot \beta}$$

$\in \mathbb{R}^{n \times n}$

new sample prediction

$$\hat{y} = \sum_{i=1}^n \alpha_i \underline{\phi(x)^T \phi(x_i)}$$

$K(x, x_i)$

e.g. $\underline{x} \in \mathbb{R}^2 \Rightarrow x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{bmatrix}$$

$$\phi(x_i)^T \phi(x_j) = x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{i2} x_{j1} x_{j2} + x_{i2}^2 x_{j2}^2$$

$$= (x_{i1} x_{j1} + x_{i2} x_{j2})^2$$

$$= (x_i^T x_j)^2$$

if $x_i \in \mathbb{R}^p$ then computing this takes $O(p)$ time

if $\phi(x) \in \mathbb{R}^d$,
then computing this takes $O(d)$ time

Recall $p \ll d$

for ϕ above, $d \approx O(p^2)$

Popular kernels:

Kernel function : $\phi(x_i)^T \phi(x_j) = K(x_i, x_j)$

- Polynomials of degree q :

$$K(x_i, x_j) = (x_i^T x_j)^q$$

- Polynomials of degree up to q :

$$K(x_i, x_j) = (x_i^T x_j + 1)^q$$

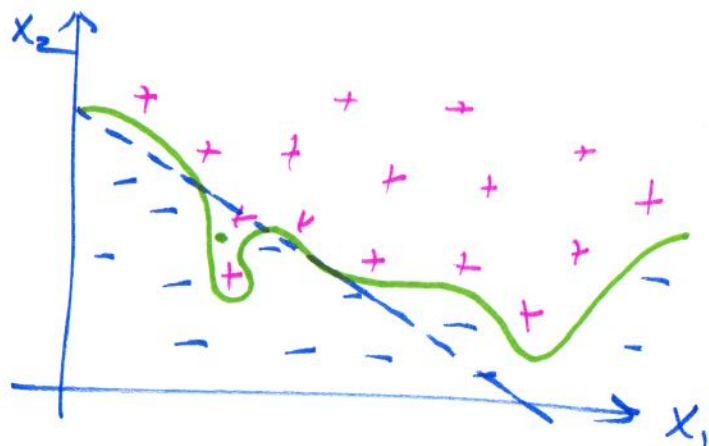
- Gaussian/Radial kernels

$$K(x_i, x_j) = \exp\left(\frac{-\|x_i - x_j\|_2^2}{2\sigma^2}\right), \quad d = \infty$$

Choosing a kernel:

- Ⓐ. Kernel regression gives linear boundary in high-dimensional feature space
↔ nonlinear boundary in original space.

- Ⓑ. Looking for perfect classifier kernel can be bad



risk "overfitting"