

# Kernel Regression

# Objectives

1

- Why use higher-dimensional feature spaces
- Reformulate regression in terms of kernels
- Popular kernels
- Cautions and considerations

# Higher dimensional feature spaces extend regression



数据点很复杂，传统的线性回归没法完美拟合数据  
→ 我们只能引入 x 的平方等高次项

$$d(x) = w_1 x$$

$d(x)$  是决策函数而不是 cost function (拿一个 x 算出来我们要的预测值)

$$d(x) = w_3 x^3 + w_2 x^2 + w_1 x + w_0$$

Let  $\underline{x} = [x, x_2 \dots x_m]^T \in \mathbb{R}^m$

$\phi^T(\underline{x})$  目的就是为了升维 (这样数据加到高维空间就会好些)

Consider  $d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w}$ ,  $\phi(\underline{x}) \in \mathbb{R}^P$   
 $P > m$

Example:  $\underline{x} = [x, x_2]^T$ ,  $\underline{\phi}^T(\underline{x}) = [x^2, x^2, \sqrt{2}x_1 x_2, x_1, x_2, 1]$

Finding  $\underline{w}$ : "training" data  $\underline{x}^i, d^i, i=1, 2, \dots N$

$$\min_{\underline{w}} \sum_{i=1}^N (d^i - \underline{\phi}^T(\underline{x}^i) \underline{w})^2 + \lambda \|\underline{w}\|_2^2 \quad (\text{Ridge})$$

$$\underline{d} = [d^1 \ d^2 \ \dots \ d^N]^T \quad \xrightarrow{N \text{ samples}}$$

$$\underline{\Phi} = [\phi(\underline{x}^1) \ \phi(\underline{x}^2) \ \dots \ \phi(\underline{x}^N)]^T \quad \xrightarrow{(N \times P)}$$

$$\min_{\underline{w}} \|\underline{d} - \underline{\Phi} \underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

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

Regression is a weighted sum of "Kernels" 3

$$d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w} = \underline{\phi}^T(\underline{x}) (\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I})^{-1} \underline{\Phi}^T \underline{d}$$

$\underline{\Phi}^T$  is  $P \times N$   $\rightarrow$   $\underline{\Phi}^T \underline{\Phi}$   $P \times P$   
 $\underline{I}$  is  $N \times P$

Matrix identity:  $(\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I})^{-1} \underline{\Phi}^T = \underline{\Phi}^T (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1}$  (activity)

Thus  $d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{\Phi}^T (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d}$

$\rightarrow$  样本数  $\times$  样本数 (如果大的话计算量就非常小了)

Note:  $[\underline{\Phi} \underline{\Phi}^T]_{i,j} = \underline{\phi}^T(x^i) \underline{\phi}(x^j)$  Define "Kernel"

$\rightarrow$  核函数 (向量点积代表相似度)

$$[\underline{\phi}^T(\underline{x}) \underline{\Phi}^T]_j = \underline{\phi}^T(\underline{x}) \underline{\phi}(x^j)$$

Let  $\underline{\alpha} = [\alpha_1, \dots, \alpha_N]^T$   $\rightarrow N \times 1$

$$= (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d}$$

$\downarrow$  新数据和矩阵的乘积 = 新数据和 j-th training sample  $x^j$ 's dot product

$d(\underline{x}) = \sum_{i=1}^N \alpha_i \underline{\phi}^T(\underline{x}) \underline{\phi}(x^i) = \sum_{i=1}^N \alpha_i K(x, x^i)$

$\downarrow$   $d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{\Phi}^T \underline{\alpha}$  (新数据和训练数据的相似度)

Kernel methods find  $d(\underline{x})$  without computing  $\phi(\underline{x})$

$$d(\underline{x}) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i)$$

$$[\underline{K}]_{ij} = \underline{\phi}^T(\underline{x}^i) \underline{\phi}(\underline{x}^j) = K(\underline{x}^i, \underline{x}^j)$$

$$\begin{aligned} \underline{\alpha} &= (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d} \rightarrow \underline{K} = \underline{\Phi} \underline{\Phi}^T \\ &= (\underline{K} + \lambda \underline{I})^{-1} \underline{d} \end{aligned}$$

↓ Gram matrix

$\underline{K}$  can be computed efficiently!

Ex: Monomials of degree  $g$   $\underline{\phi}(\underline{x}) \rightarrow x_1^g, x_1^{g-1}x_2, \dots, x_3^{g-5}x_6^2x_8^3 \dots$

$$K(\underline{u}, \underline{v}) = \underbrace{\underline{\phi}^T(\underline{u})}_{O(P)} \underline{\phi}(\underline{v}) = \underbrace{(\underline{u}^T \underline{v})}_O(M)^g \quad (\text{activity})$$

$$P = \frac{(g+M-1)!}{g! (M-1)!} \text{ terms}$$

↑ 扩展后的特征总数  
↓ 扩展的次数  
↓ 原始特征数

Suppose  $M=10, g=5 \rightarrow P \sim 2000$

$M=100, g=5 \rightarrow P \sim 10^8$  computing  $O(P)$  vs  $O(M)$   
memory  $O(NP)$  vs  $O(N^2)$

# Popular kernels depend on similarity of $\underline{u}, \underline{v}$ 5

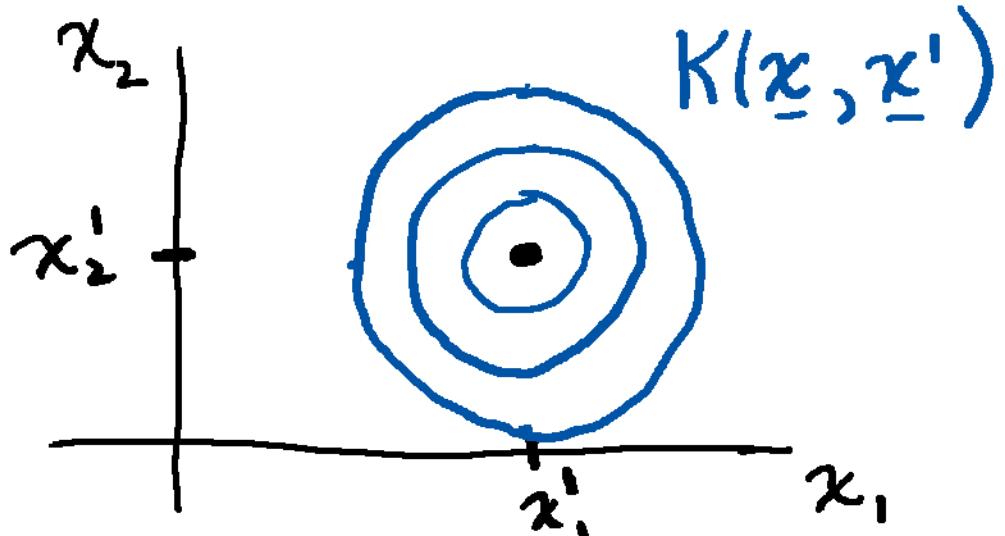
$$\underline{u}^T \underline{v} = \|\underline{u}\|_2 \|\underline{v}\|_2 \cos \theta$$

↓ 定义好的，直接用

Monomials of degree  $q$ :  $K(\underline{u}, \underline{v}) = (\underline{u}^T \underline{v})^q$

Polynomials up to degree  $q$ :  $K(\underline{u}, \underline{v}) = (\underline{u}^T \underline{v} + 1)^q$

Gaussian/radial Kernel:  $K(\underline{u}, \underline{v}) = \exp\left\{-\frac{\|\underline{u} - \underline{v}\|_2^2}{2\sigma^2}\right\}$



- No explicit  $\phi(\underline{x})$
  - All polynomial orders
  - smoothness controlled by  $\sigma$
- ↓ sigma

# Kernel regression considerations

6

$$d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w} \quad \text{vs} \quad d(\underline{x}) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i)$$

原始视角  
核方法视角

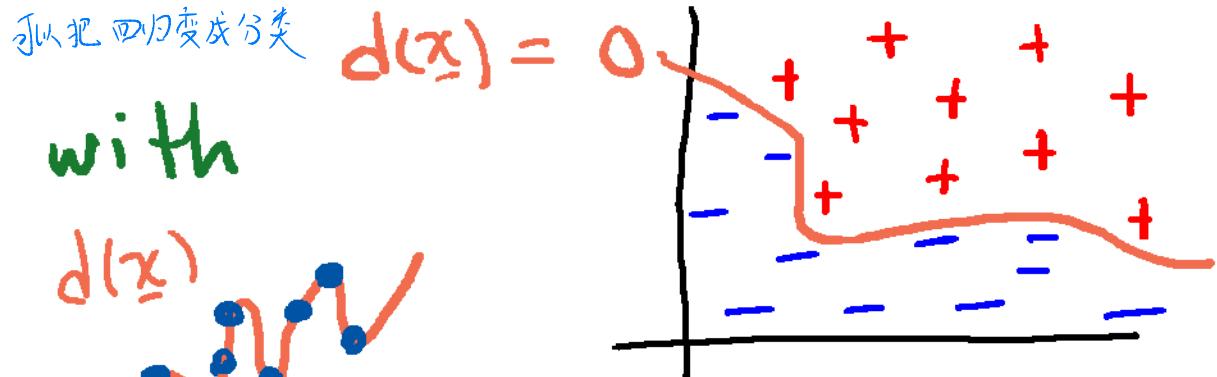
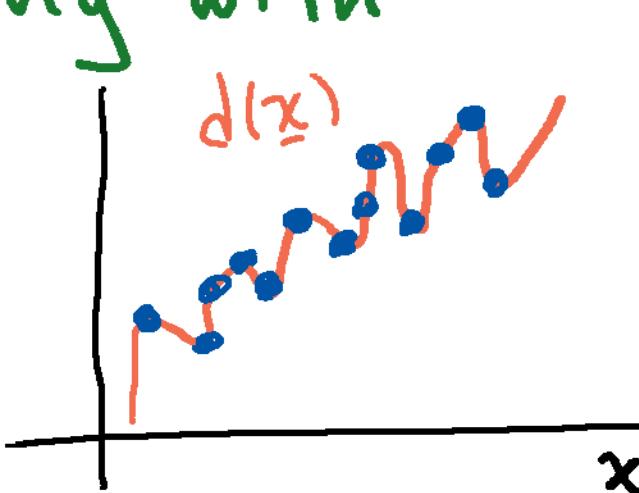
- Store and compute  $\underline{\alpha}$  ( $N \times 1$ ) vs  $\underline{w}$  ( $P \times 1$ )

- Binary classification  $\text{sign}\{d(\underline{x})\}$

从把回归变成分类

- Avoid "overfitting" with high-D feature spaces

(cross-validation)



Copyright 2019  
Barry Van Veen