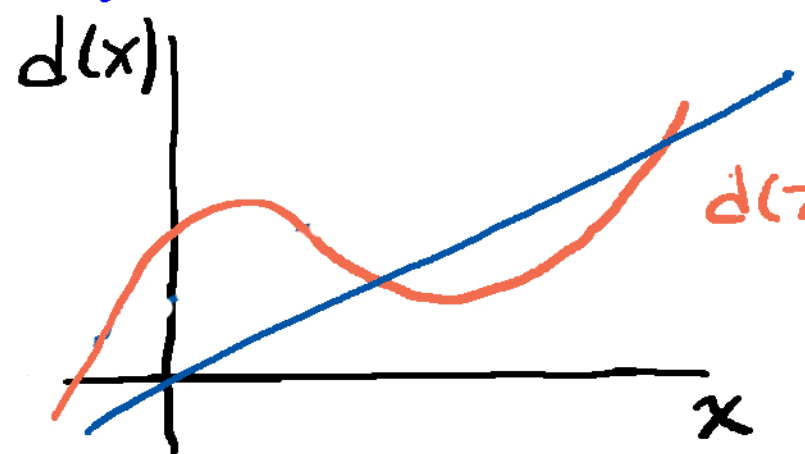


Kernel Regression

Objectives

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

Higher dimensional feature spaces extend ² regression



$$d(x) = w_1 x$$

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

Let $\underline{x} = [x_1, x_2, \dots, x_m]^T \in \mathbb{R}^m$

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

Example: $\underline{x} = [x_1, x_2]^T$, $\underline{\phi}^T(\underline{x}) = [x_1^2, x_2^2, \sqrt{2}x_1x_2, x_1, x_2, 1]$

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

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

$$\underline{d} = [d^1, d^2, \dots, d^N]^T$$

$$\underline{\Phi} = [\phi(\underline{x}^1), \phi(\underline{x}^2), \dots, \phi(\underline{x}^N)]^T \quad (N \times p) \Rightarrow \underline{w} = (\underline{\Phi}^T \underline{\Phi} + \lambda \underline{I})^{-1} \underline{\Phi}^T \underline{d}$$

Regression is a weighted sum of "kernels" 3

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

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 $d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underbrace{\underline{\Phi}^T (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d}}_{N \times N}$

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

Let $\underline{\alpha} = [\alpha_1 \dots \alpha_N]^T$
 $= (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d}$

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

Kernel methods find $d(\underline{x})$ without computing $\phi(\underline{x})$ 4

$$d(\underline{x}) = \sum_{i=1}^N \alpha_i K(\underline{x}, \underline{x}^i) \quad \underline{\alpha} = (\underline{\Phi} \underline{\Phi}^T + \lambda \underline{I})^{-1} \underline{d} \rightarrow \underline{K} = \underline{\Phi} \underline{\Phi}^T$$
$$= (\underline{K} + \lambda \underline{I})^{-1} \underline{d} \leftarrow$$

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

\underline{K} can be computed efficiently!

Ex: Monomials of degree q $\underline{\phi}(\underline{x}) \rightarrow \underbrace{x_1^q, x_1^{q-1} x_2, \dots, x_1^{q-5} x_2^2 x_3^3 \dots}_{\text{terms}}$

$$K(\underline{u}, \underline{v}) = \underbrace{\underline{\phi}^T(\underline{u}) \underline{\phi}(\underline{v})}_{O(P)} = \underbrace{(\underline{u}^T \underline{v})^q}_{O(M)} \quad (\text{activity}) \quad P = \frac{(q+M-1)!}{q! (M-1)!} \text{ terms}$$

Suppose $M=10, q=5 \rightarrow P \sim 2000$ computing $O(P)$ vs $O(M)$
 $M=100, q=5 \rightarrow P \sim 10^8$ memory $O(NP)$ vs $O(N^2)$

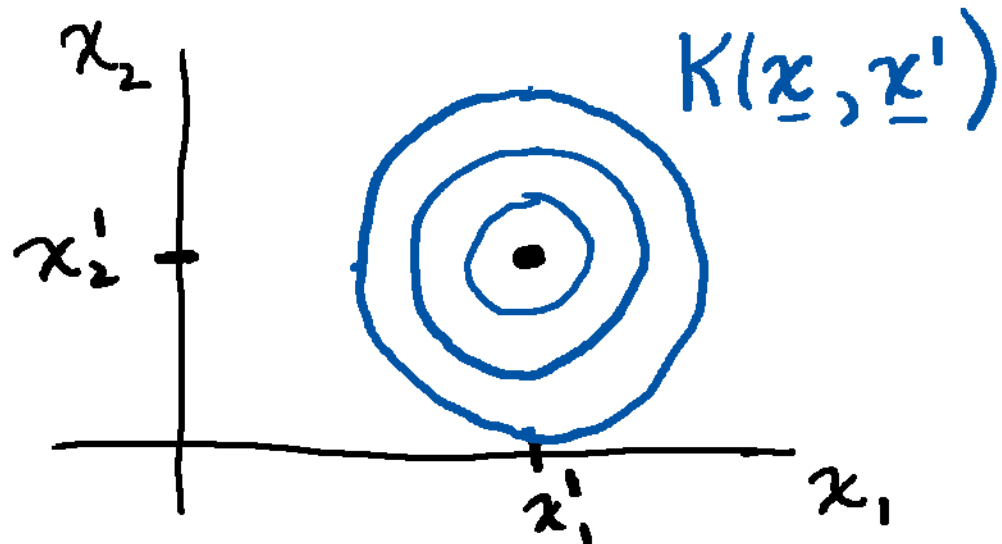
Popular kernels depend on similarity of $\underline{u}, \underline{v}$ 5

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

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 σ

Kernel regression considerations 6

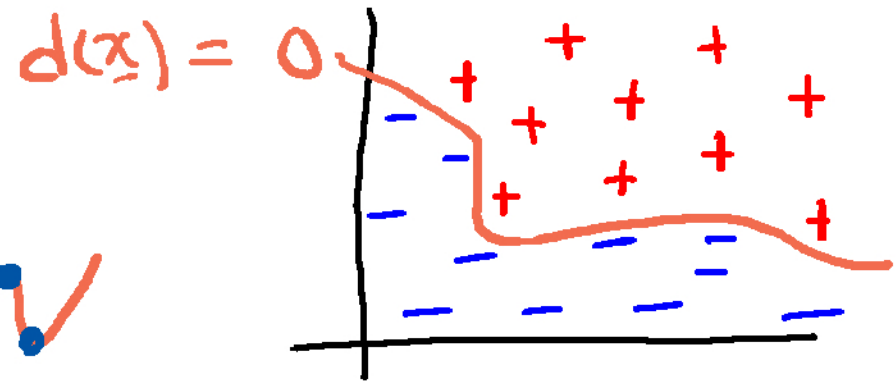
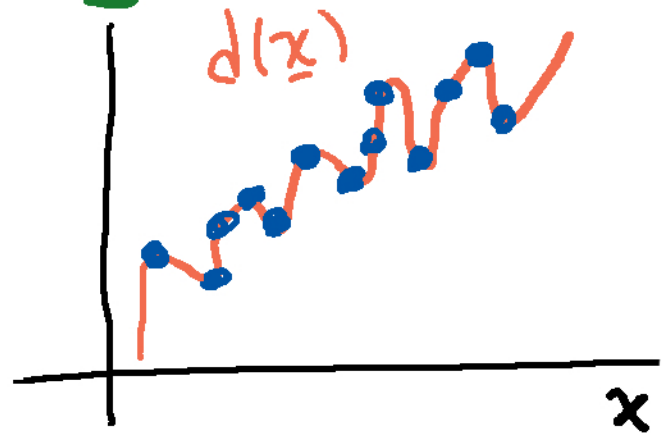
$$d(\underline{x}) = \underline{\phi}^T(\underline{x}) \underline{w} \quad \text{vs} \quad d(\underline{x}) = \sum_{i=1}^N \alpha_i k(\underline{x}, \underline{x}^i)$$

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

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

- Avoid "overfitting" with
high-D feature
spaces

(cross-validation)



Copyright 2019
Barry Van Veen