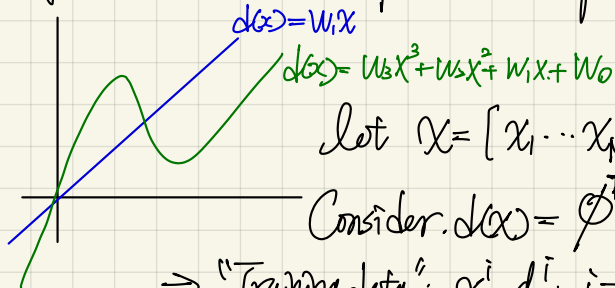


§ Kernel Regression.

(*) High-dim feature spaces

(*) Popular Kernels

- Higher dimensional feature spaces extend regressions.



$$\text{let } \mathbf{x} = [x_1 \dots x_M]' \in \mathbb{R}^M$$

$$\text{Consider } d(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \mathbf{w}, \text{ where } \boldsymbol{\phi}(\mathbf{x}) \in \mathbb{R}^p \ (p > M)$$

$$\Rightarrow \text{"Training data": } x^i, d^i, i=1, \dots, N$$

$$\Rightarrow \text{"Find } \mathbf{w}": \min_{\mathbf{w}} \sum_{i=1}^N (d^i - \boldsymbol{\phi}^T(x^i) \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2 \quad (\text{Ridge})$$

$$\text{Define } \mathbf{d} = [d^1 \dots d^N]'$$

$$\boldsymbol{\Phi} = [\boldsymbol{\phi}(x^1) \dots \boldsymbol{\phi}(x^N)]'$$

$$\Rightarrow \min_{\mathbf{w}} \|\mathbf{d} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\Rightarrow \mathbf{w} = (\boldsymbol{\Phi}' \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}' \mathbf{d} \quad (\text{Standard Ridge})$$

- Regression is a weighted sum of "kernels".

$$d(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \mathbf{w} = \boldsymbol{\phi}^T(\mathbf{x}) \left((\boldsymbol{\Phi}' \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}' \mathbf{d} \right)$$

$$\text{Since } (\boldsymbol{\Phi}' \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}' = \boldsymbol{\Phi}' (\boldsymbol{\Phi} \boldsymbol{\Phi}' + \lambda \mathbf{I})^{-1}$$

$$\text{we have } d(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\Phi} (\boldsymbol{\Phi} \boldsymbol{\Phi}' + \lambda \mathbf{I})^{-1} \mathbf{d}$$

$$\textcircled{1} [\boldsymbol{\Phi} \boldsymbol{\Phi}']_{ij} = \boldsymbol{\phi}^T(x^i) \boldsymbol{\phi}(x^j)$$

$$\textcircled{2} [\boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\Phi}]_j = \boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\phi}(x^j) \Rightarrow \text{define Kernel: } K(\mathbf{u}, \mathbf{v}) = \boldsymbol{\phi}^T(\mathbf{u}) \boldsymbol{\phi}(\mathbf{v})$$

$$\Rightarrow \text{let } \boldsymbol{\alpha} = [\alpha_1 \dots \alpha_N]' \text{ then } d(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\Phi} \boldsymbol{\alpha} = \sum_{i=1}^N \alpha_i \boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\phi}(x^i)$$

$$= (\boldsymbol{\Phi} \boldsymbol{\Phi}' + \lambda \mathbf{I})^{-1} \mathbf{d} \quad (\text{easier computation})$$

$$= \sum_{i=1}^N \alpha_i K(\mathbf{x}, x^i) \quad (*) \text{ weighted sum of } K(\mathbf{x}, x^i)$$

- Kernel methods find $d(x)$ without computing $\phi(x)$

We have $d(x) = \sum_{i=1}^N \alpha_i K(x, x_i)$ where $\alpha = (\Phi \Phi^T + \lambda I)^{-1} d$
 $\equiv (R + \lambda I)^{-1} d$

$$[R]_{ij} = \phi^T(x_i) \phi(x_j) = K(x_i, x_j)$$

$\Rightarrow R$ can be computed efficiently!

eg. R being "Monomials of degree q ": $\phi(x) \rightarrow x_1^q, x_1^{q-1}x_2, \dots, x_3^{q-5}x_6^2x_8^3$
 (power degree sum to q)

then $K(u, v) = \underbrace{\phi^T(u)}_{O(p)} \underbrace{\phi(v)}_{O(M)} = (\underbrace{u^T v}_{\text{inner prod.} \nearrow q \text{ power}})^q$

$$\Rightarrow p = \frac{(q+M-1)!}{q!(M-1)!} \text{ terms } (\uparrow \text{ when } M \text{ or } q \uparrow)$$

@ Suppose $\begin{cases} M=10, q=5 \Rightarrow p \sim 2000 \\ M=100, q=5 \Rightarrow p \sim 10^8 \end{cases}$

Computing: $O(p)$ vs $O(M)$
 Memory: $O(Np)$ vs $O(N^2)$

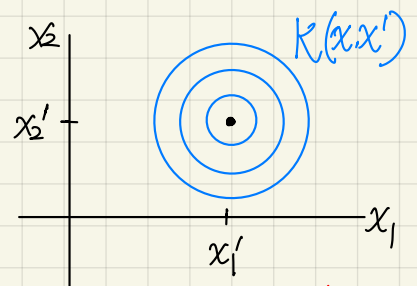
- Popular choices of kernels (depend on similarity of u, v)

Note that $u^T v = \|u\|_2 \|v\|_2 \cos \theta$

(i) Monomial of degree q : $K(u, v) = (u^T v)^q$

(ii) Polynomials up to degree q : $K(u, v) = (u^T v + 1)^q$

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



"how close u & v are"

* Not required explicit form of $\phi(x)$

* An infinite polynomial orders

* Smoothness: Controlled by σ

\Rightarrow desired choice of kernel

Since it generalizes regression

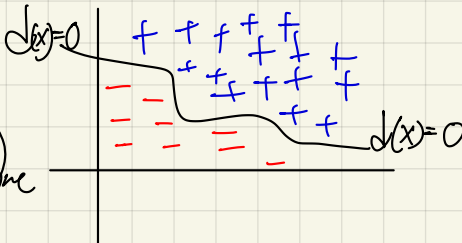
with high-dim spaces. (semiparametric)

- Kernel Regression Considerations.

$$d(x) = \phi^T(x) \mathbf{w} \longleftrightarrow d(x) = \sum_{i=1}^N \alpha_i \mathcal{K}(x, x_i) \quad (\text{Kernel})$$

(*) Store & compute: $\mathbf{W} (p \times 1)$ vs. $\alpha (N \times 1)$

(*) Binary classifier: $\text{sign}(d(x))$ (more complex. with high-dim space)



(*) **WANT**: Avoid overfitting with high-dim feature spaces
 [challenge of working with high-dim]

⇒ Choose σ^2 parameter.

⇒ with CV.

§ Kernel-based SVM.

(*) Linear max margin classifier.

(*) Kernel version of Hinge Loss w/ Ridge Regression.

• SVM defines "max-margin" classifier.

$$\Rightarrow \begin{cases} \tilde{x}'\tilde{w} + b = 0 \\ \text{or } X'W = 0 \end{cases} \quad \text{where } W = [\tilde{w}' \ b'] \text{ depends only on Support Vectors}$$

Const.

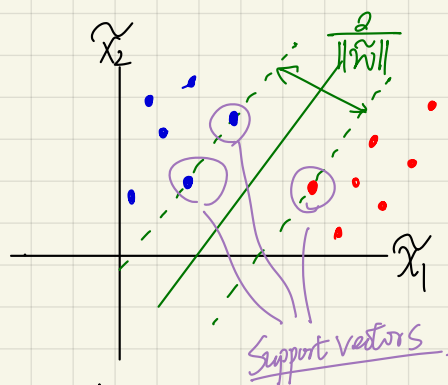
Recall: kernel regression.

$$d(x) = \underbrace{\phi^T(x)}_{\text{"high-dim feature space"}} W = \sum_{j=1}^N \alpha_j \underbrace{K(x, x^j)}_{\text{"kernel"}}$$

$$\text{Let } \phi^T(x) = x \Rightarrow K(x, x^j) = x'x^j$$

$$\Rightarrow d(x) = \sum_{j=1}^N \alpha_j K(x, x^j) = \sum_{j=1}^N \alpha_j x'x^j = x' \underbrace{\sum_{j=1}^N \alpha_j x^j}_W$$

$$\text{So, } W = \sum_{j=1}^N \alpha_j x^j \Rightarrow \text{all } \alpha_j = 0 \text{ except for the } \underline{\text{Support Vectors!}}$$



• Use kernels for "nonlinear" decision boundaries.

\Rightarrow high dim feature space: $x \rightarrow \phi(x)$

$$\boxed{\text{eg}} \quad \phi(x) = [x_1^2 \ x_2^2 \ \dots \ x_2 x_4 \ \dots \ x_{m-1} \ x_m \ 1]$$

$$\Rightarrow \hat{d}(x) = \text{sign}(\phi^T(x) W)$$

Then, "Hinge loss with Ridge": $\min_w \left(\sum_{i=1}^N (1 - d^i \phi^T(x^i) W)_+ \right) + \lambda \|W\|_2^2$

Claim. $W = \sum_{j=1}^N \phi(x^j) \alpha_j$ (difference: $x^j \rightarrow \phi(x^j)$)

- "Kernel trick" replaces $\phi(x^j) \phi(x^i)$ with $\mathcal{K}(x^j, x^i)$

From $(*)$: $w = \sum_{j=1}^N \phi(x^j) \alpha_j$.

\Rightarrow hinge loss w/ Ridge: $\min_w \left(\sum_{i=1}^N (1 - d^i \phi^T(x^i) w)_+ \right) + \lambda \|w\|_2^2$

$\Rightarrow \min_{\alpha} \left[\sum_{i=1}^N (1 - d^i \phi^T(x^i) \sum_{j=1}^N \alpha_j \phi(x^j))_+ \right] + \lambda \sum_{i=1}^N \alpha_i \phi^T(x^i) \sum_{j=1}^N \alpha_j \phi(x^j)$

$= \min_{\alpha} \sum_{i=1}^N (1 - d^i \sum_{j=1}^N \alpha_j \underbrace{\phi^T(x^i) \phi(x^j)}_{\mathcal{K}(x^i, x^j)})_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \underbrace{\phi^T(x^i) \phi(x^j)}_{\mathcal{K}(x^i, x^j)}$

\Rightarrow SVM: $\min_{\alpha} \sum_{i=1}^N (1 - d^i \sum_{j=1}^N \alpha_j \mathcal{K}(x^i, x^j))_+ + \lambda \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \mathcal{K}(x^i, x^j)$

\otimes It turns out SVM has sparse α (in linear case + high-dim space)

• Decision boundary: $d(x) = 0 = \phi^T(x) w = \sum_{j=1}^N \alpha_j \mathcal{K}(x, x^j)$

• Boundary (hinge loss) depends only on Support Vectors

• Kernels ($\mathcal{K}(u, v)$) measures similarity of vectors u, v

\square eq. Gaussian kernels: $\mathcal{K}(u, v) = \exp\left(\frac{-\|u - v\|_2^2}{2\sigma^2}\right)$

\rightarrow Solve α via Gradient Descent

