# Lec 16

Thursday, November 7, 2019   11:10

$\underline{Recap}$: Augmenting the data

$$X \longmapsto \varphi(X)$$

& then plug in to SVC

Called this "augmented SVC"

$$\min_{\beta, \beta_0} \quad \frac{1}{2} \|\beta\|_2^2 + c \sum_{i=1}^{n} \xi_i$$

$$\xi_i \geq 0 \qquad \xi_i \geq 1 - Y_i (\beta^T \varphi(X_i) + \beta_0)$$

Last time rephrased as:

$$\min_{\substack{\alpha \in \mathbb{R}^n \\ \beta_0}} \quad \frac{1}{2} \alpha^T \underline{K} \alpha + c \sum_{i=1}^{n} \xi_i$$

$$\xi_i \geq 0 \qquad \xi_i \geq 1 - Y_i (\underline{K}_i^T \alpha + \beta_0)$$

Where $\underline{K}_{ij} = \underbrace{\varphi(X_i)^T \varphi(X_j)}_{\kappa(X_i, X_j)}$ kernel fn

Conclusion: the augmented SVC only

depends on the $\underset{\wedge}{\text{example}}$ feature data

via the kernel Gram matrix $\underline{K}$

$\underline{Idea}$: Instead of defining $\varphi$,

define kernel $\kappa$

$\underline{Ex}$: Linear kernel: $\kappa(x, x') = x^T x' \longrightarrow$ gives the orignl

$k(x) = x$      SVC

# Ex: Polynomial kernel

First consider $p=1$

$$k(x, x') = (1 + x x')^d$$

E.g. for $d=3$

$$k(x, x') = 1 + 3 x \cdot x' + 3 (x^2)(x'^2) + (x)^3 (x')^3$$

Q: What $\varphi$ gives $k(x, x') = \varphi(x)^T \varphi(x')$?

$$\varphi(x) = \begin{pmatrix} 1 \\ \sqrt{3} \, x \\ \sqrt{3} \, x^2 \\ x^3 \end{pmatrix}$$

For general $d$:

$$k(x, x') = 1 + d \, x x' + \cdots + d (x x')^{d-1} + (x x')^d$$

$$\varphi(x) = \left( 1, \; \sqrt{d} \, x, \; \cdots, \; \sqrt{d} \, x^{d-1}, \; x^d \right)$$

For general $d, p$:

$$k(x, x') = (1 + x^T x')^d$$

$\varphi(x) = $ vector of all monomials of $(x_1, \cdots, x_p)$ of degree $\leq d$

e.g. $\underbrace{x_1 x_2 x_3^2}_{\text{monomial}}$ of degree 4

Let's get crazy:

$$k(x, x') = \exp\left( - \|x - x'\|_2^2 / \sigma^2 \right)$$

Radial Basis Function (RBF) kernel

Corresponding $\varphi$ sent $x$ to $\infty$ space

of functions (infinite dim)

We can even do this with unstructured data.

e.g. $X =$ string of "$A, C, G, T$"

$\varphi(x) =$ 〜〜〜〜

z.f. $k(x, x') = e^{-\text{edit dist}(x,x')^2 / \sigma^2}$

Note: Not everything is a kernel

Need: ① $k(x, x') = k(x', x)$ symmetric
　　　② $k(x, x') \geq 0$
　　　③ PSD
　　　　　$\forall x_1, \ldots, x_n$
　　　　　$K_{ij} = k(x_i, x_j) \Rightarrow K$ is PSD

If we have ①-③ then $\exists \varphi$
　　　s.f. $k(x, x') = \varphi(x)^\top \varphi(x')$
Key insight: don't even need to know $\varphi$


Kernel trick generalizes to other problems
___

## Kernel Ridge Regression

Recall ridge regression

$$\min_{\beta, \beta_0} \sum_i (y_i - \beta^\top x_i - \beta_0)^2 + \lambda \|\beta\|_2^2$$

$\xrightarrow{\text{augment}}$ $\sum_i (y_i - \beta^\top \varphi(x_i) - \beta_0)^2 + \lambda \|\beta\|^2$

$\longrightarrow$ similar arguments as before
(representer theorem + kernel trick)
yield that this also just
depends on the data
via the kernel matrix

# Kernel PCA

# Recall PCA

Amounts to finding the eigen decomp of

$$\left[\overline{X}\,\overline{X}^{\top}\right]_{ij} = x_i^{\top} x_j$$

$X \xrightarrow[\text{linear mapping}]{}$ lower dim space that explains most of $x$'s variance

Kernelize this

$X \xrightarrow{\varphi}$ hi-dim (possibly $\infty$-dim) Space $\xrightarrow[\text{projection}]{\text{linear}}$ low-dim space that explains most of variance in $\varphi$

$\underbrace{\phantom{XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX}}$

$\underline{\underline{\text{non}}}$ linear mapping into low-dim space
that explains most of $x$'s variance

Apply eigen decomp to the kernel gram
matrix $\underline{K}_{ij} = k(x_i, x_j)$
Centered version: apply the eigen decomp

on the centered kernel gram matrix

$$\underline{\underline{K}}^c = \underline{\underline{K}} - \frac{\mathbb{1}_{n\times n}}{n}\underline{\underline{K}} - \bar{K}\frac{\mathbb{1}_{n\times n}}{n} + \mathbb{1}_{n\times n}\underline{\underline{K}}\,\mathbb{1}_{n\times n}$$

(← $n\times n$ matrix of all 1s)

Linear PCA: find the $q$ linear fns s.t.

$$Z = (u_1^T x, u_2^T x, \cdots, u_q^T x) \quad \Leftarrow PC's$$

explains the most variance

Kernel PCA: find the $q$ fns $f_j \in \mathcal{H} = \text{Span}(\{K(x, \cdot) : x \in$

s.t.

$$Z = (f_1(x), f_2(x), \cdots, f_q(x)) \quad \Leftarrow \text{kernel } PC's$$

explains the most ~~the~~ variance

# Neural Networks

The "Vanilla" neural network

Add an "embedding step" to the linear model

$$X \xrightarrow[A]{} Z = \varphi(x) \longrightarrow \hat{f}(x) = \beta^T Z$$

big diff for NN:
also learn $\varphi$

What if $Z \in \mathbb{R}^q$   $Z = Ax$   $A \in \mathbb{R}^{q \times p}$

$$\hat{f}(x) = \beta^T z(x) = \beta^T A x = (A^T \beta)' x$$

$\underbrace{\qquad}$ Still a linear model

So the NN approach:

apply an activation/nonlinearity.

$$z_j = \sigma(A_j^T x) \quad \text{where } \sigma \text{ is } \underline{\text{nonlinear}}$$

$$\text{e.g.} \quad \sigma(u) = \frac{1}{1+e^{-u}}$$

Now $\beta^T z(x)$ is $\underline{\underline{\text{not}}}$ linear in $x$