

## Lec 15

Tuesday, November 5, 2019 11:03

Recap: Which linear separator to use?

Idea: maximize the margin

If it's possible to linearly separate the data then max margin classifier given by

$$\min_{\beta_0, \beta_1, \dots, \beta_p \in \mathbb{R}} \frac{1}{2} \|\beta\|_2^2$$

$$\text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1 \quad \forall i=1, \dots, n$$

(recall  $y_i \in \{-1, +1\}$ )

Nonseparable case      Soft margin

Allow some slack  $\xi_i \geq 0$

Only going to require

$$y_i (\beta^T x_i + \beta_0) \geq 1 - \xi_i \quad \forall i$$

But penalize positive slack

We get the

"Support Vector Classifier" (SVC)

$$\min_{\beta_0, \beta_1, \dots, \beta_p \in \mathbb{R}} \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$\xi_1, \dots, \xi_n \in \mathbb{R}$$

$$\text{s.t. } y_i (\beta^T x_i + \beta_0) \geq 1 - \xi_i \quad \forall i=1, \dots, n$$

$$\xi_i \geq 0 \quad \forall i=1, \dots, n$$

Trade off b/w margin ( $\|\beta\|_2^2$ )  
 & classification errors ( $\sum_i \xi_i$ )

Controlled by param  $c > 0$

E.g.  $c \rightarrow \infty \Rightarrow$  get hard margin, max margin classifier (if it's linearly separable)

Connect the differences to logistic regression

SVC can be equivalently written as:

$$\min \sum_i \max(0, 1 - y_i(\beta^T x_i + \beta_0)) + \frac{1}{2c} \|\beta\|_2^2$$

$$= \min \sum_i \text{loss}_{\text{Hinge}}(y_i, \beta^T x_i + \beta_0) + \lambda \|\beta\|_2^2$$

$$\text{loss}_{\text{Hinge}}(y, \hat{y}) = \max(0, 1 - y\hat{y}), \quad \lambda = \frac{1}{2c}$$

Logistic regression w/ ridge regularization

$$\min \sum_i \log(1 + e^{-y_i(\beta^T x_i + \beta_0)}) + \lambda \|\beta\|_2^2$$

- Both optimize a loss fn + (possible) regularization
- Different loss fn
  - Usually, for classification, prefer large loss

- For prob prediction, prefer logistic reg.  
 (SVC gives only a decision boundary,  
 not an estimate  $P(Y=1|X)$ )

(to get prob estimates from SVC  
 can use Platt scaling)

## Kernelizing the SVM:

What happens if we augment the  
 features  $x$  using a feature map:

$$\mathcal{Q} : \mathbb{R}^p \rightarrow \mathbb{R}^q \quad q \geq p$$

$$\text{e.g. } \mathcal{Q}(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix} \quad \begin{matrix} p=1 \\ q=2 \end{matrix}$$

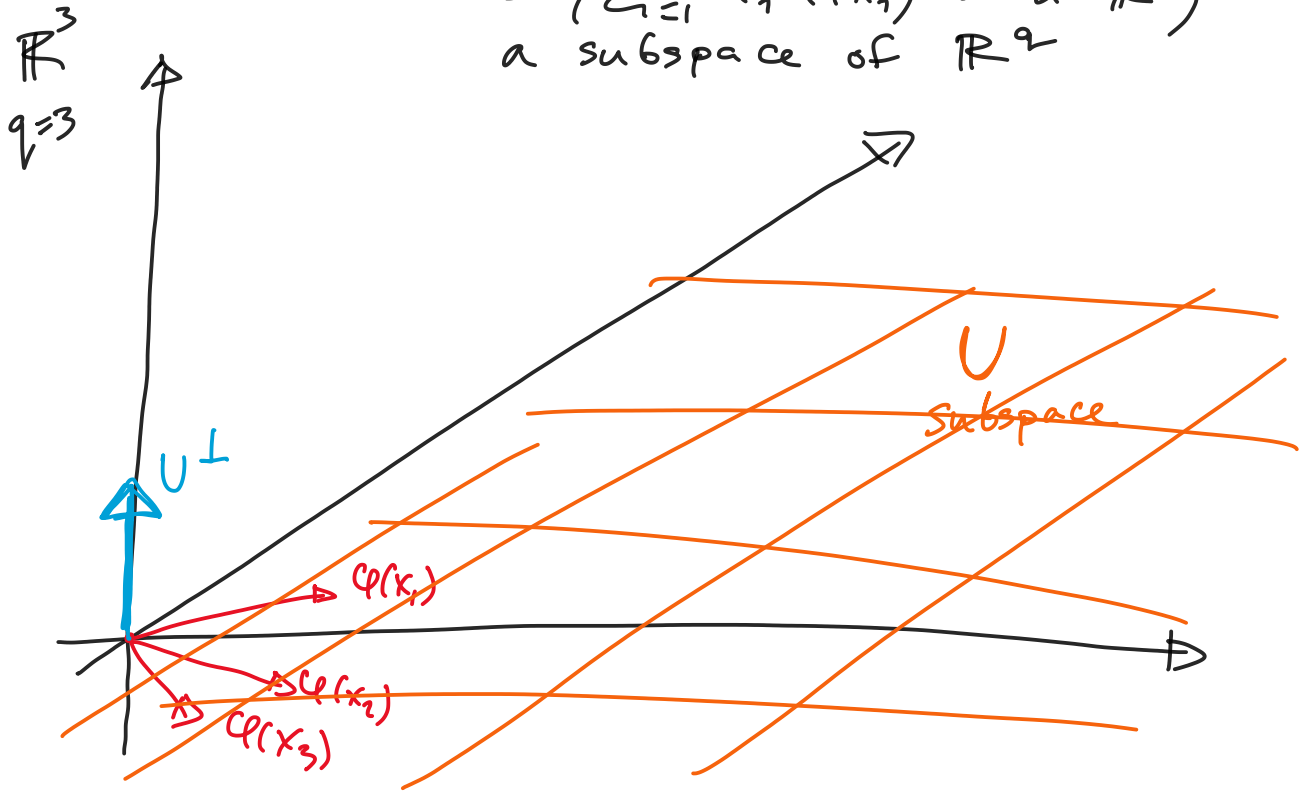
$$\mathcal{Q}\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{pmatrix} \quad \begin{matrix} p=2 \\ q=5 \end{matrix}$$

Get the augmented SVC:

$$\begin{aligned} \min_{\substack{\beta \in \mathbb{R}^q \\ \beta_0 \in \mathbb{R}}} & \quad \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n \xi_i \\ & \quad \xi_i \geq 0 \\ & \quad \xi_i \geq 1 - \gamma_i (\beta^\top \mathcal{Q}(x_i) + \beta_0) \end{aligned}$$

We can even take  $q \rightarrow \infty$

Consider  $U = \text{span}\{\varphi(x_1), \dots, \varphi(x_n)\}$   
 $= \left\{ \sum_{i=1}^n \alpha_i \varphi(x_i) : \alpha \in \mathbb{R}^n \right\}$   
 a subspace of  $\mathbb{R}^q$



Orthogonal complement

$$U^\perp = \{v \in \mathbb{R}^q : v^T u = 0 \quad \forall u \in U\}$$

$$= \{v \in \mathbb{R}^q : v^T \varphi(x_i) = 0 \quad \forall i\}$$

Any vector  $z \in \mathbb{R}^q$  can be written

$$\text{as } z = u + v$$

$$\text{where } u \in U$$

$$v \in U^\perp$$

$$\begin{aligned} \|z\|^2 &= z^T z = (u+v)^T (u+v) \\ &= u^T u + \cancel{2u^T v} + v^T v \\ &= \|u\|^2 + \|v\|^2 \end{aligned}$$

Write  $\beta = u + v$

$u \in \mathcal{U} = \text{span}\{\phi(x_1), \dots, \phi(x_n)\}$   
 $v \in \mathcal{U}^\perp$  ( $v^\top \phi(x_i) = 0$ )

$$\|\beta\|^2 = \|u\|^2 + \|v\|^2$$

$$\beta^\top \phi(x_i) = u^\top \phi(x_i)$$

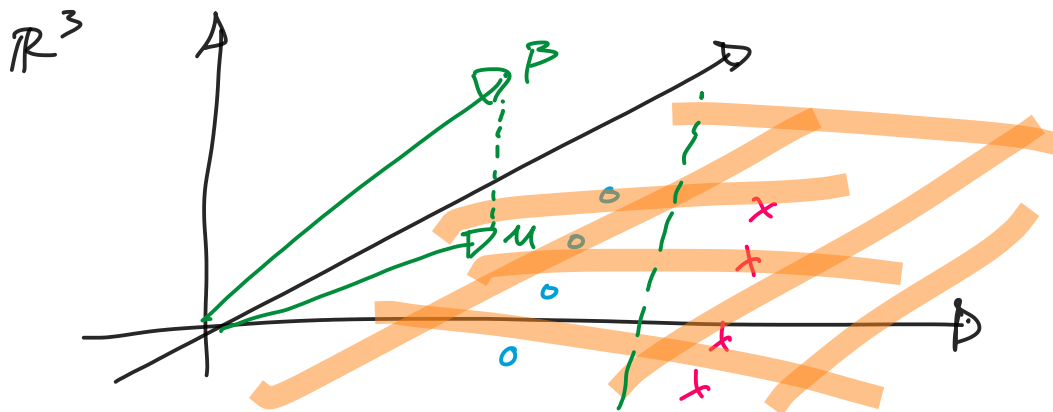
Get that the augmented SVC is:

$$\min_{\substack{u \in \mathcal{U}, \beta_0 \\ v \in \mathcal{U}^\perp}} \frac{1}{2} \|u\|^2 + \frac{1}{2} \|v\|^2 + C \sum_i \xi_i$$

$$\xi_i \geq 0$$

$$\xi_i \geq 1 - \gamma_i (u^\top \phi(x_i) + \beta_0)$$

At optimality get  $v=0$



Conclusion optimal  $\beta$  can be written as

$$\beta = \sum_i \alpha_i \phi(x_i)$$

(known "representer trick")

Let's do SVC & optimize over these  $\alpha$ 's next

$$\min \quad \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \underbrace{\varphi(x_i)^T \varphi(x_j)}_{K(x_i, x_j)} + C \sum_i \xi_i$$

$$\xi_i \geq 0$$

$$\xi_i \geq 1 - \left( \gamma_i \sum_j \alpha_j \varphi(x_j) \right)^T \varphi(x_i) + \beta_0$$

Define  $K(x, x') = \varphi(x)^T \varphi(x')$

$$\underline{K}_{ij} = K(x_i, x_j) \quad \underline{K} \in \mathbb{R}^{n \times n}$$

Can write augmented SVC as

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

$$\xi_i \geq 0$$

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

"kernel trick"