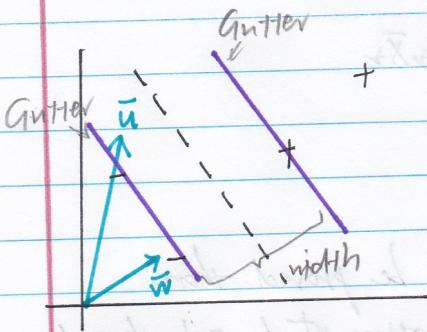


## Support Vector Machines



$$\bar{w} \cdot \bar{u} \geq c \quad c = -b$$

$\boxed{\text{① } \bar{w} \cdot \bar{u} + b \geq 0 \text{ then the unknown is a "+"}}$   
this is our decision rule

$\bar{w}$  is a vector perpendicular to the middle line  
 $\bar{u}$  represent the unknown to be classified.

Although we defined the decision rule, we yet have to find out what  $\bar{w}$  and  $b$  is, the goal is to maximize the width

$$\bar{w} \cdot \bar{x}_+ + b \geq 1 \quad (\bar{x}_+ \text{ represent a known "+" sample})$$

$$\bar{w} \cdot \bar{x}_- + b \leq -1 \quad ("-" \text{ sample})$$

for our convenience, we introduce a new variable  $y_i$  such that  $y_i = \begin{cases} 1 & \text{for "+" samples} \\ -1 & \text{for "-" samples} \end{cases}$

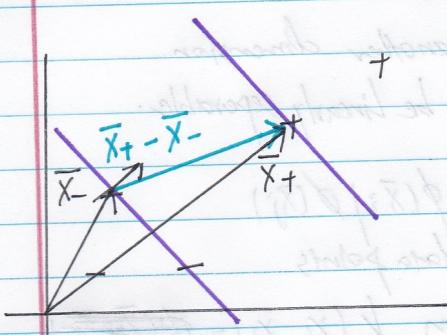
we multiply  $y_i$  to both sides of the two inequalities

$$\Rightarrow y_i(\bar{w} \cdot \bar{x}_+ + b) \geq 1 \cdot 1$$

$$y_i(\bar{w} \cdot \bar{x}_- + b) \geq (-1)(-1) \quad \begin{matrix} \text{(since multiple negative)} \\ \text{values to both sides} \\ \text{flips the sign} \end{matrix}$$

the two inequalities above are actually the same.

$$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 \geq 0, \quad y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0 \quad \begin{matrix} \text{for } x_i \\ \text{in gutter} \end{matrix}$$



$$\text{width} = (\bar{x}_+ - \bar{x}_-) \cdot \frac{\bar{w}}{\|\bar{w}\|} \quad \begin{matrix} \text{this is a unit} \\ \text{vector of the same} \\ \text{direction of } \bar{w} \end{matrix}$$

$$\text{③ } 1-b \quad 1+b$$

$$\Rightarrow \text{width} = (1-b+1+b) \cdot \frac{\bar{w}}{\|\bar{w}\|} = \frac{2}{\|\bar{w}\|}$$

so to maximize  $\frac{2}{\|\bar{w}\|}$ , we want to minimize  $\|\bar{w}\|$

for our convenience, we will instead find the min of  $\frac{1}{2} \|\bar{w}\|^2$

we use Lagrange multipliers here to find the minimum of a term with constraints (the constraint is  $y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$ )

$$L = \frac{1}{2} \|\bar{w}\|^2 - \sum_{i=1}^n [y_i(\bar{w} \cdot \bar{x}_i + b) - 1]$$

$$\frac{\partial L}{\partial \bar{w}} = \bar{w} - \sum_{i=1}^n y_i \bar{x}_i = 0 \Rightarrow \bar{w} = \sum_{i=1}^n y_i \bar{x}_i$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^n y_i = 0 \Rightarrow \sum_{i=1}^n y_i = 0$$

plugin  $\bar{w} = \sum_{i=1}^n y_i \bar{x}_i$  in  $L$

$$L = \frac{1}{2} \left( \sum_{i=1}^n y_i \bar{x}_i \right) \cdot \left( \sum_{i=1}^n y_i \bar{x}_i \right)$$

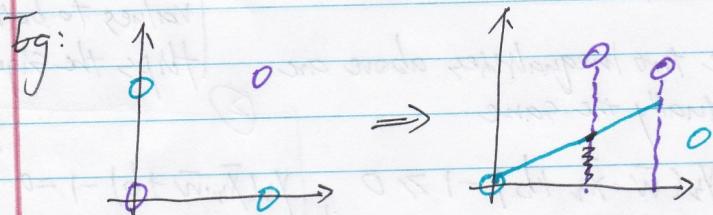
$$\text{we will eventually get } L = \sum_{i=1}^n -\frac{1}{2} \sum_j \sum_{i=1}^n y_i y_j \bar{x}_i \cdot \bar{x}_j$$

plug  $\bar{w} = \sum_{i=1}^n y_i \bar{x}_i$  in the decision rule

if  $\sum_{i=1}^n y_i \bar{x}_i \cdot \bar{w} + b \geq 0$  Then classify  $\bar{w}$  as "+"

The decision rule depends on the dot product too  
SVM only works for linear classification.

but if we transform the non-linear data  
to another perspective



he can transform a data point  $x_i$  to another dimension using function  $\phi(x)$  for the data to be linearly separable.

then the decision rule will depend on  $\phi(\bar{x}_i) \cdot \phi(\bar{x}_j)$

the dot product of two transformed data points

to simplify, we can have a kernel function  $k(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$

Thereof's two popular choice of the kernel function

$$① (\bar{w} \cdot \bar{v} + 1)^n$$

$$② e^{-\frac{\|\bar{x}_i - \bar{x}_j\|}{\sigma}}$$

$$\phi(\bar{x}_i) \cdot \phi(\bar{x}_j)$$