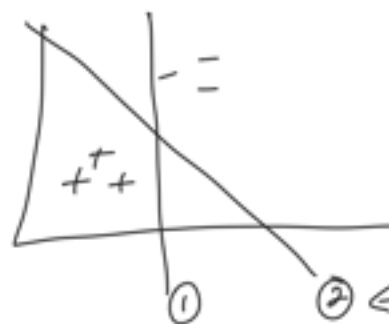


Last time: max margin



← because it's farther from the data.  
generalizes in the face of perturbations

$$\text{margin}(x_n, y_n) = \frac{y_n (w^T x_n + w_0)}{\|w\|_2} \quad \left. \begin{array}{l} \text{relative distance} \\ \text{absolute distance} \end{array} \right\}$$

two formulations:

hard-margin (separable)

$$\begin{aligned} \min_{w, w_0} \quad & \frac{1}{2} \|w\|_2^2 \\ \text{s.t.} \quad & y_n (w^T x_n + w_0) \geq 1 \end{aligned}$$

soft margin

$$\begin{aligned} \min_{w, w_0} \quad & \frac{1}{2} \|w\|_2^2 + c \sum \xi_n \\ \text{s.t.} \quad & y_n (w^T x_n + w_0) \geq 1 - \xi_n, \quad \xi_n \geq 0 \end{aligned}$$

↳ both convex! in high dimensions, solving these can still be slow  
today: convert these into "dual" formulation [trick of optimization]  
↳ incorporating kernels/bases

apply Lagrangian:

$$\min_{w, w_0} \max_{\alpha_n} \left[ \frac{1}{2} \|w\|_2^2 - \sum_n \alpha_n [y_n (w^T x_n + w_0) - 1] \right]$$

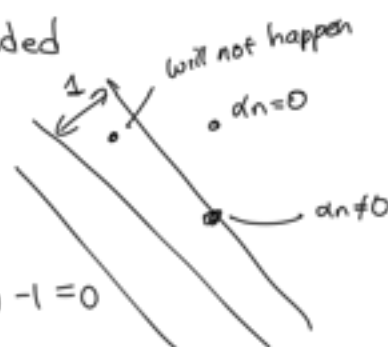
thought experiment:  
consider new obj  
old obj + dII (constraint met)  
↕  
old obj + d amount of violation  
∇ old obj + ∇ violation = 0

1) if constraint is violated,  $y_n (w^T x_n + w_0) - 1 < 0$   
then  $\alpha$  can become very large & make obj unbounded

2) if the constraint is met & has some slack,  
 $y_n (w^T x_n + w_0) - 1 > 0$  (NOT equal to 0)

$\alpha_n = 0$

only case where  $\alpha_n \neq 0$  will be when  $y_n (w^T x_n + w_0) - 1 = 0$



Now notice

3) given (2), at optimality, we recover the original objective

Next steps: by property known as strong duality, we can swap the max & min

$$\max_{\alpha} \min_{w, w_0} \left[ \frac{1}{2} \|w\|_2^2 - \sum_n \alpha_n [y_n (w^T x_n + w_0) - 1] \right]$$

goal: 1st let's solve the min. problem, next deal w/ the max

$$\begin{aligned} \cdot \frac{\partial \mathcal{L}}{\partial w} &= w - \sum_n a_n y_n x_n = 0 \Rightarrow \boxed{w = \sum_n a_n y_n x_n} \\ \cdot \frac{\partial \mathcal{L}}{\partial w_0} &= -\sum_n a_n y_n = 0 \quad \left. \vphantom{\frac{\partial \mathcal{L}}{\partial w_0}} \right\} \text{new constraint} \end{aligned}$$

$$\max_{\alpha} \underbrace{\left( \frac{1}{2} (\sum_n a_n y_n x_n)^T (\sum_{n'} a_{n'} y_{n'} x_{n'}) - \sum_n a_n y_n (\sum_{n'} a_{n'} y_{n'} x_{n'})^T x_n \right)}_{w^T w} \quad \left. \vphantom{\max_{\alpha}} \right\} \text{sub. for } w$$

$$-\frac{\sum_n a_n y_n w_0 + \sum_n a_n}{w_0 \sum_n a_n y_n} \downarrow 0$$

$$\boxed{\max_{\alpha} -\frac{1}{2} \sum_{n,n'} a_n a_{n'} y_n y_{n'} \underbrace{x_n^T x_{n'}}_{\text{kernel}} + \sum_n a_n \quad \text{s.t.} \quad \sum_n a_n y_n = 0, a_n \geq 0}$$

new optimization problem in terms of  $\alpha$  (instead of  $w$ !)  
easy to solve! (convex)

Note! soft margin: constraint  $a_n \geq 0$  becomes  $\bullet \quad \geq \quad \geq \quad 0$

once we solve for  $\alpha^*$ , we can recover  $w, w_0$ :

$$\boxed{w^* = \sum_n \alpha_n^* y_n x_n}$$

$\hookrightarrow$  implication:

$$\Rightarrow \underbrace{\sum_n \alpha_n^*}_{\substack{\text{weight} \\ \text{per} \\ \text{datum} \\ \text{in} \\ \text{training} \\ \text{set}}} \underbrace{y_n}_{\substack{\text{vote} \\ \text{of that} \\ \text{pt.}}} \underbrace{x_n^T x}_{\substack{\text{how similar} \\ \text{are } x_n \text{ to} \\ \text{test pt. } x}} + w_0$$

typically # of  $\alpha_n > 0$  is small... so sparse classifier / fast at test time.

finally  $w_0$ : take any  $n$  s.t.  $\alpha_n > 0 \Rightarrow y_n (\underbrace{w^*}_{\text{sub in}}^T x_n + w_0) = 1 \quad \uparrow \text{solve for } w_0$

$$w_0 = \frac{1}{y_n} - w^{*T} x_n$$

$\mathcal{D}$  doesn't seem to be gone...

if we kernelize:  $x \rightarrow \phi(x)$ , we still need to deal w/  
the dimensionality of  $\phi$

Notice! All computations ONLY require  $x^T x'$  or  $\phi(x)^T \phi(x')$

idea: we can define  $\underbrace{k(x, x')}_{\text{kernel function}} = \phi(x)^T \phi(x')$

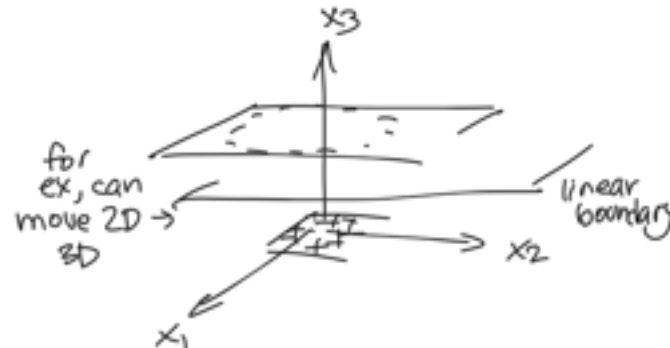
reminder: why bases help:



hard  $\rightarrow$  soft



NOT linearly sep.



how to choose  $\phi(x)$ ? expert, NN,  
today: highly expressive spaces that correspond to  
our notions of similarity

Let's talk about  $K(x, x')$ :

ex:  $K(x, x') = x^T x'$  (trivial)

ex:  $K(x, x') \Rightarrow \phi: [x_1^2, x_1 x_2, \dots, x_1 x_D, x_2 x_1, x_2^2, x_2 x_D, \dots, x_D^2]$

ex:  $K(x, x') = \sum_{d, d'} x_d x_{d'} x'_d x'_{d'}$

ex:  $\boxed{K(x, x')} = (1 + x^T x')^q$  polynomial kernel (corresponds to a  $\phi$  that is exponential in  $q$ )

ex:  $K(x, x') = \exp\left(\frac{-\|x - x'\|_2^2}{2\sigma^2}\right)$  RBF (radial basis function) corresponds to an  $\infty$ -D  $\phi$

Can we just "create" some  $K(x, x')$ ? "similarity function"  
how do we know that it is valid?

Mercer's theorem:

- Matrix form
- let matrix  $K$  be an  $N \times N$  matrix where  $K_{nn'} = K(x_n, x_{n'}) \Rightarrow \boxed{K(x, x')}$
  - if  $K$  corresponds to a valid kernel, then  $K$  is positive semi-definite ( $\sum_i K_{ii} z_i^2 \geq 0 \forall z$ ) for any choice of input set  $\{x_1, \dots, x_N\}$

$\Rightarrow \int_{x, x'} f(x) K(x, x') f(x') dx dx' \geq 0 \quad \forall f$  (nice)

then  $K(x, x')$  is a valid kernel  $\Leftrightarrow$  is an inner product in some space.

[C] Thought exercise: consider RBF kernel w/ diff  $\sigma$ 's

$\exp\left(\frac{\|x - x'\|_2^2}{-2\sigma^2}\right)$

as well as linear kernel  $K(x, x') = x^T x'$

... profit and how?

b can these overfit?



$$\sum \alpha y_n \underline{K(x_n, x)} + w_0$$

Svm