

DMML, 19 Feb 2019

Back to supervised learning

Geometric approach

Assume (like clustering) data are points (x_1, x_2, \dots, x_n) in space

Suppose category depends on position

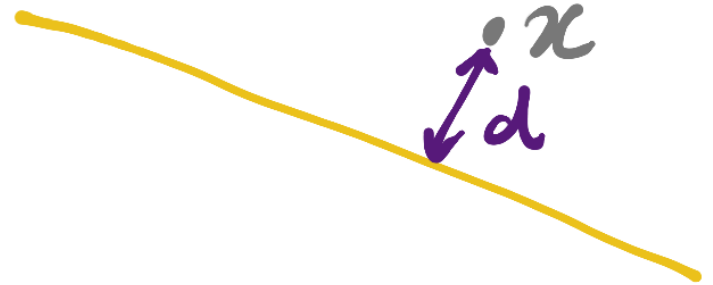
Simplest case - linearly separable

Good separator — find good \bar{w} say w^*

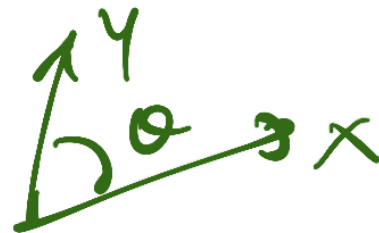
$$x^T w^* = 0$$

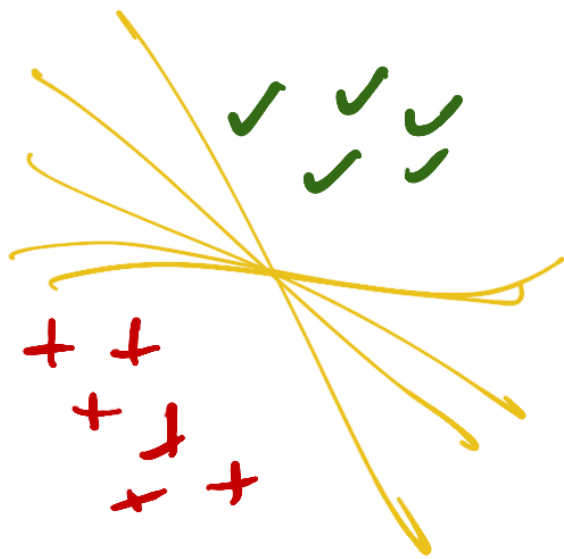
Arbitrary point x

$$d = \frac{x^T w^*}{\|w^*\|}$$



Recall: $x \cdot y = \|x\| \|y\| \cos \theta$





Bigger separation \Rightarrow Easier to find a w^*

The other quantity of interest = $\|x\|$

Aim for w^* s.t. $x^T w^* \geq 1$ for all positive points
 $x^T w^* \leq -1$ for all neg pts

Min distance is $\frac{1}{\|w\|} = \gamma$

Perceptron Algorithm

- Initialize w to 0
- Pick x in training set s.t $x^T w$ has wrong sign
 - If x is a positive example, $w \leftarrow w + x$
 - If x is a negative example, $w \leftarrow w - x$

Thm If there is w^* s.t. $x^T w^* \geq 1$ (≤ -1)
for all positive (negative) examples,
then the Perceptron algorithm makes
at most R^2 / w^{*2} updates
 $(R/\gamma)^2$

R is $\max_{x \in X} \|x\|$
└ training set

Proof

Ideal w^*

Current w

- Each update increases $w^T w^*$ by at least 1

+ve example

$$(w+x)^T w^* = w^T w^* + \underbrace{x^T w^*}_{\geq 1}$$

-ve example

$$(w-x)^T w^* = w^T w^* - \underbrace{x^T w^*}_{\leq -1}$$

- With each update, $|w|^2$ increases by at most R^2

+ve $(w+x)^T(w+x) = |w|^2 + \underbrace{2x^T w}_{< 0} + |x|^2$

$$\leq |w|^2 + |x|^2 = |R|^2$$
$$\leq |w|^2 + |R|^2$$

Suppose we make M updates

- $W^T W \geq M$ (M updates, add at least 1 each time)

- $|W|^2 \leq M |R|^2$ (M updates, add at most $|R|^2$ each time)

$$|W| \leq \sqrt{M} |R|$$

Observe $\frac{W^T W^*}{|W^*|} \leq |W|$

Projection of W on W^* is at most $|W|$

$$|w| \leq R\sqrt{M}$$

$$w^T w^* \geq M$$

$$\frac{M}{w^*} \leq R\sqrt{M}$$

$$\sqrt{M} \leq R w^*$$

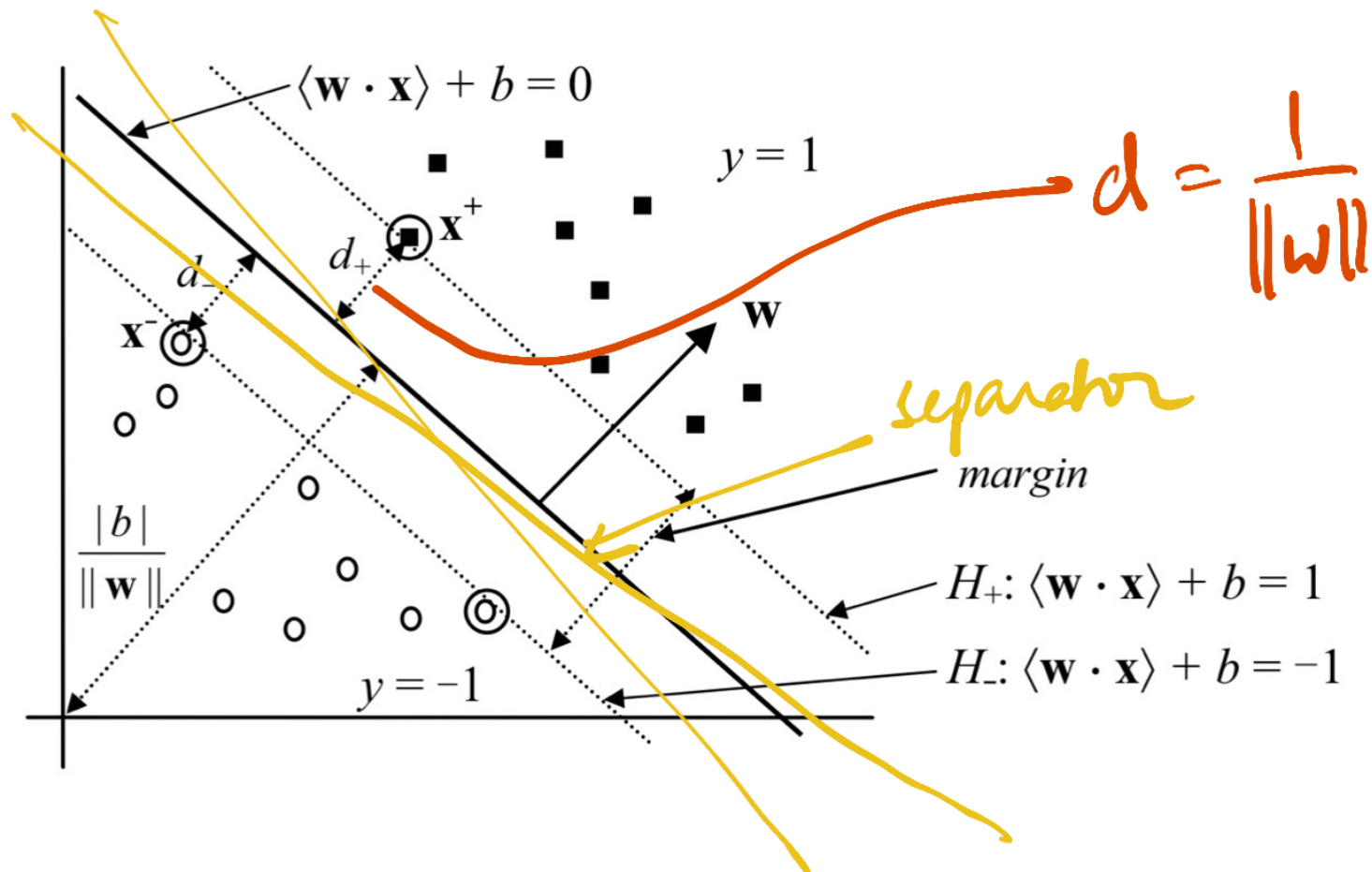
$$M \leq R^2 (w^*)^2$$

Not clear that the w we get after

M updates is w^*

Can we find w^* ?

Instead of iterative update, formulate a global optimization problem.



Total width of the "margin" is $\frac{2}{\|w\|}$

Find w that maximizes margin & separates the points as $+1/-1$

For convenience.

$$\text{Minimize } \frac{\|w\|^2}{2} \approx \frac{\langle w \cdot w \rangle}{2}$$

Optimization problem

$$\text{Minimize } \frac{\langle W \cdot W \rangle}{2}$$

Subject to

$$\langle W \cdot x_i \rangle + b \geq 1 \text{ for positive } x_i$$

$$\langle W \cdot x_i \rangle + b \leq -1 \text{ for negative } x_i$$

Associate $y_i = +1/-1$ with
positive/negative x_i

$$y_i (\langle W \cdot x_i \rangle + b) \geq 1 \text{ for all } x_i$$

Minimize $\frac{\langle W \cdot W \rangle}{2}$ \leftarrow Quadratic in W

Subject to

$$y_i (\langle W \cdot x_i \rangle + b) \geq 1 \text{ for all } i$$

Unknown : W

\nwarrow Linear in W

constrained

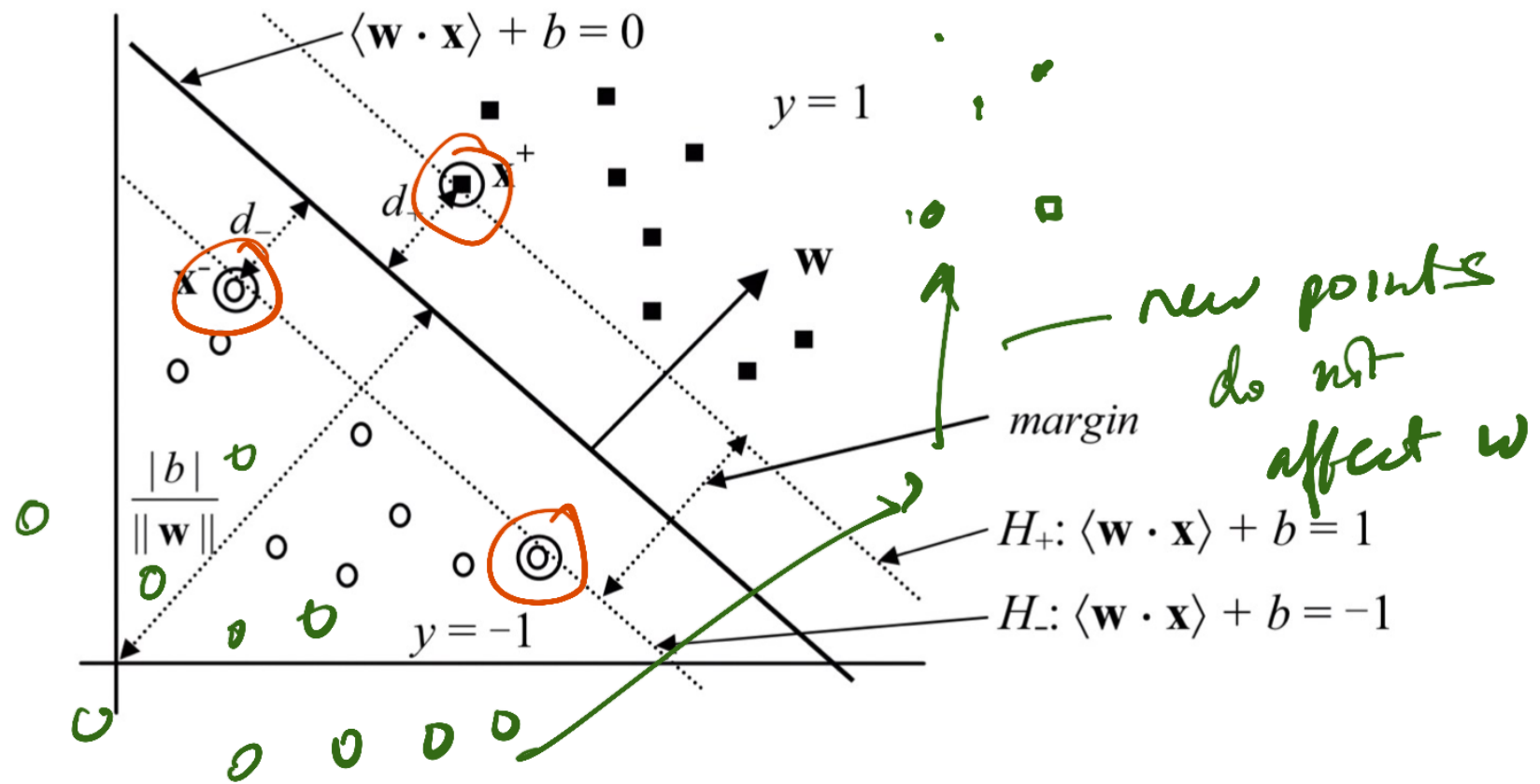
Quadratic optimization problem
 \downarrow

Under certain (syntactic) conditions [KKT]

This can be solved

Solves same problem as Perceptron Algo,
but provides max-margin separator

Support Vector Machine (SVM)



Solution will be in terms of the x_i 's
 that lie on the margin = support
vectors

Optimization problem

→ Dual

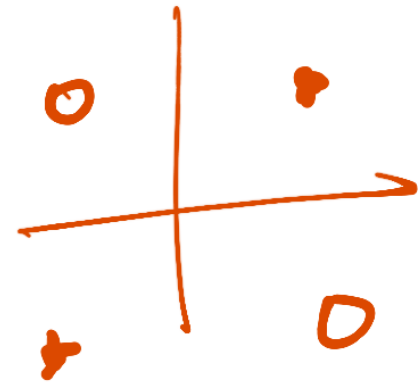
Solve dual problem - the solution can
be expressed entirely in terms of
dot products $x_i \cdot x_j$

If we can compute $x_i \cdot x_j$,
we can find w efficiently

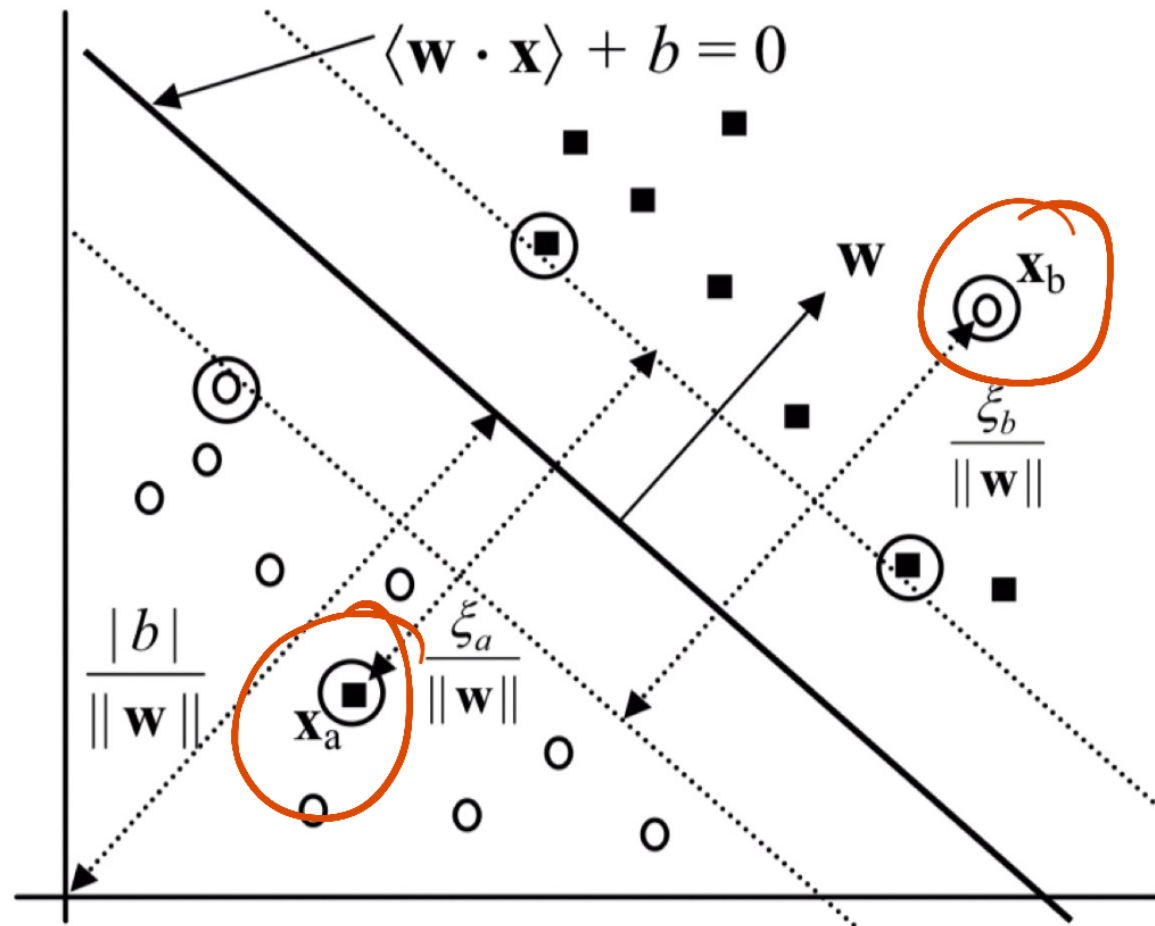
} Will play a
role
later

Not linearly separable

Extreme case \longrightarrow



Smaller case - "small errors"



Fudge factor

$$\langle W \cdot x \rangle + b \geq 1 - \xi_i \text{ for } +ve \ x$$

$$\langle W \cdot x \rangle + b \leq 1 + \xi_i \text{ for } -ve \ x$$

Find w & ξ_i 's

Penalize big ξ_i - errors

Minimize $\frac{\langle w \cdot w \rangle}{2} + \frac{1}{2} \sum_i \xi_i^2$

Subject to

$$y_i (\langle w \cdot x_i \rangle + b) \geq 1 - \xi_i, \forall i$$

Same type of quadratic constrained optimization

"Soft margin SVM"

Next time

Dealing with not linearly separable
data in a general sense

Solution - geometrically transform data
to be linearly separable

Move from (x, y) to (r, θ)