

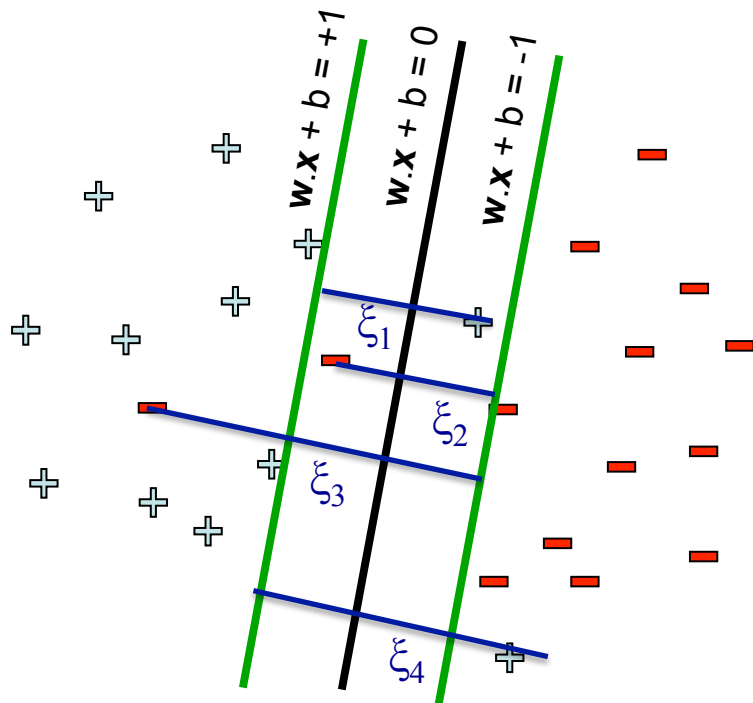
# Support vector machines (SVMs)

## Lecture 4

David Sontag  
New York University

Slides adapted from Luke Zettlemoyer, Vibhav Gogate,  
and Carlos Guestrin

# Key idea #1: Allow for *slack*



$$\text{minimize}_{\mathbf{w}, b, \xi} \quad \sum_j \xi_j$$

$$\left( \mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1 - \xi_j \quad , \quad \forall j \quad \xi_j \geq 0$$

↑  
“slack variables”

Solve for the optimal value  $\xi_j^*$  as a function of  $\mathbf{w}$  and  $b$ :

$$\text{If } (w \cdot x_j + b) y_j \geq 1, \text{ then } \xi_j^*(\mathbf{w}, b) = 0$$

$$\text{If } (w \cdot x_j + b) y_j < 1, \text{ then } \xi_j^*(\mathbf{w}, b) = 1 - (w \cdot x_j + b) y_j$$



$$\xi_j^* = \max(0, 1 - (w \cdot x_j + b) y_j)$$

## Equivalent hinge loss formulation

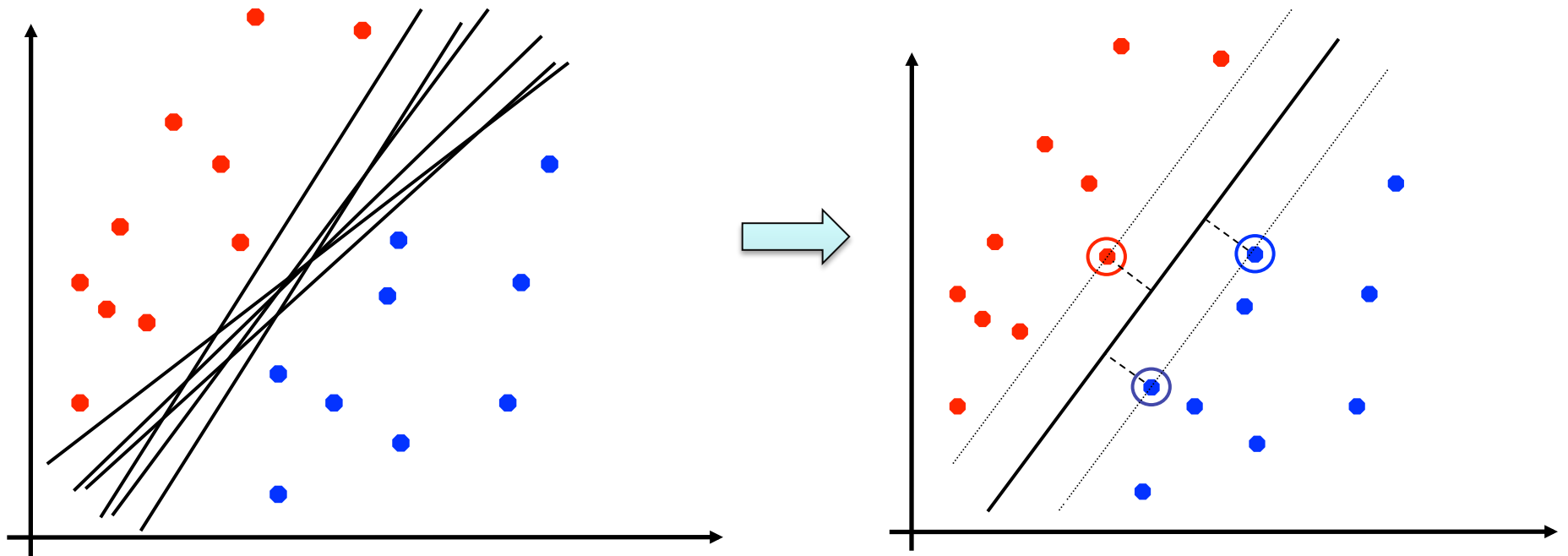
$$\begin{aligned} & \text{minimize}_{\mathbf{w}, b, \xi} \quad \sum_j \xi_j \\ & \left( \mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1 - \xi_j \quad , \forall j \quad \xi_j \geq 0 \end{aligned}$$

Substituting  $\xi_j = \max(0, 1 - (w \cdot x_j + b) y_j)$  into the objective, we get:

$$\min_{w, b} \sum_j \max \left( 0, 1 - (w \cdot x_j + b) y_j \right)$$

Now an *unconstrained* optimization problem. No longer a *linear* objective, but it is *convex*.

## Key idea #2: seek large margin

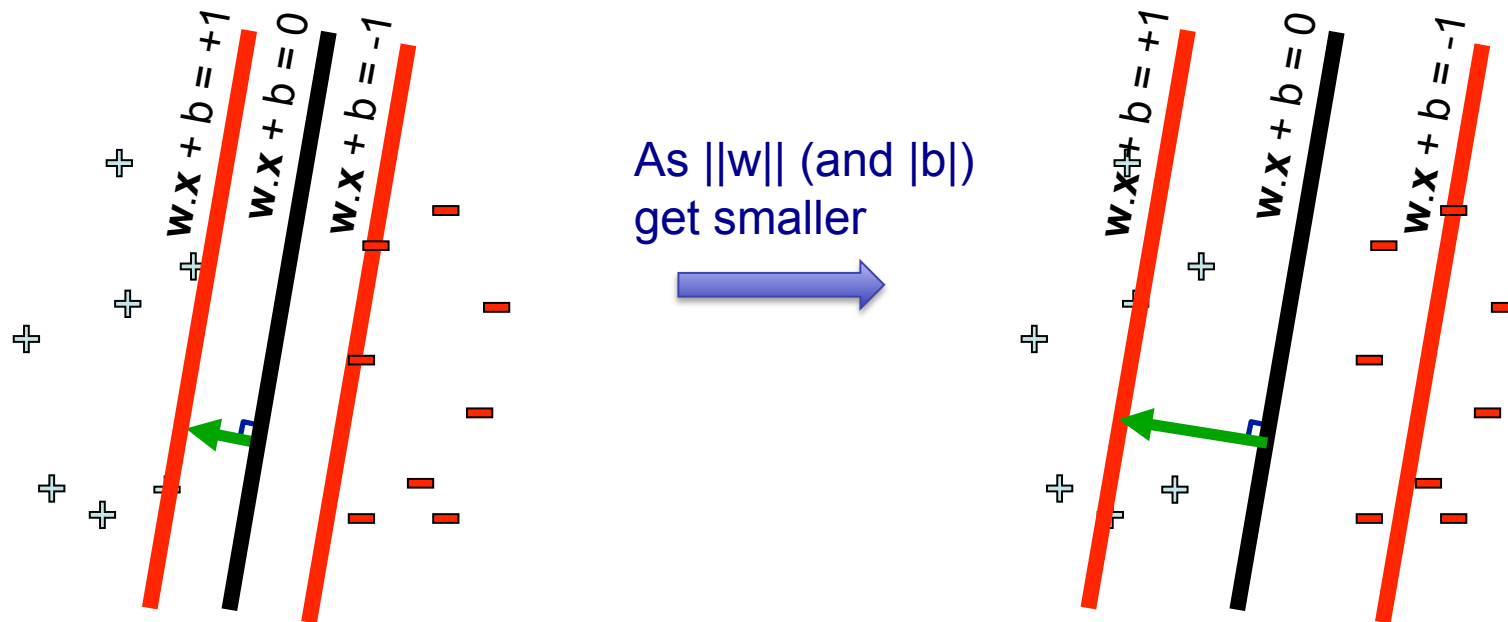


## Key idea #2: seek large margin

- Consider the constraints:

$$y_t (w \cdot x_t + b) \geq 1 \quad \forall t$$

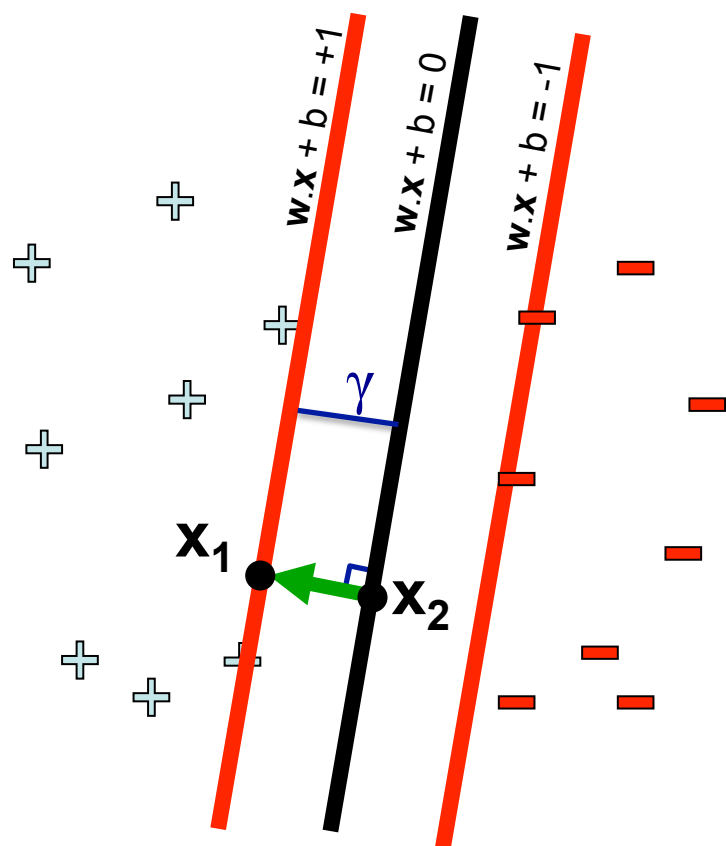
- As the norm of the weight vector  $\|w\|$  and  $b$  get **smaller**, the optimization problem becomes infeasible:



# What is $\gamma$ (geometric margin) as a function of $\mathbf{w}$ ?

$\gamma_i =$  Distance to  $i$ 'th data point

$$\gamma = \min_i \gamma_i$$



$$w \cdot x_1 + b = 1$$

$$w \cdot x_2 + b = 0$$

---


$$w \cdot (x_1 - x_2) = 1$$

Plug in

We also know that:

$$x_1 - x_2 = \gamma \frac{w}{\|w\|}$$

$$1 = w \cdot \left( \gamma \frac{w}{\|w\|} \right) = \frac{\gamma}{\|w\|} w \cdot w = \gamma \|w\|$$

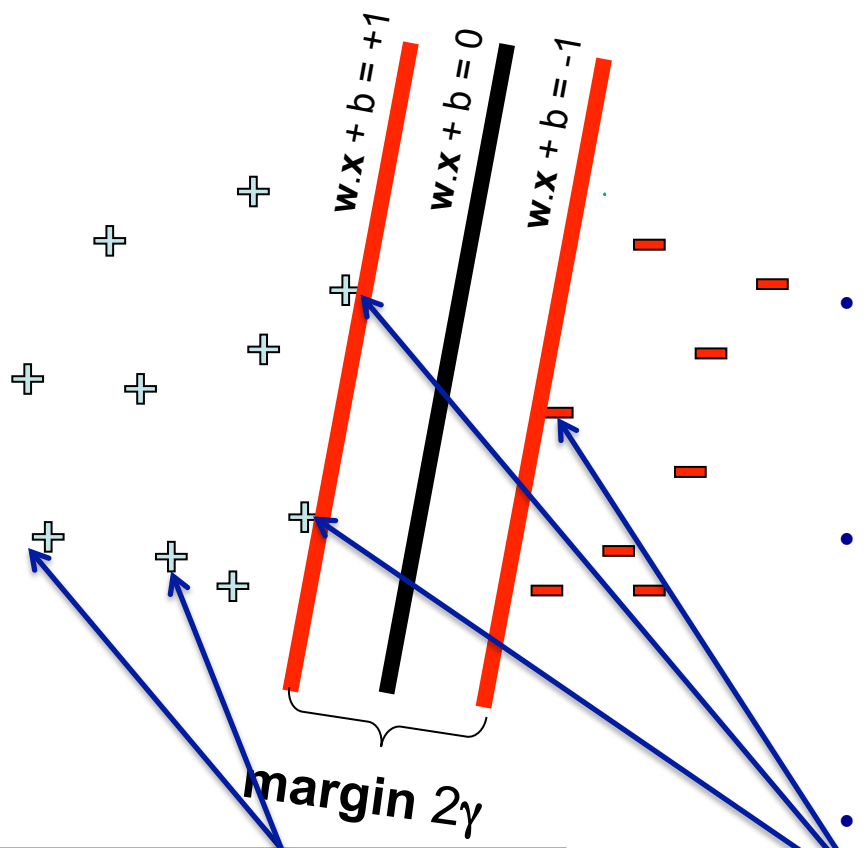
$$\text{So, } \gamma = \frac{1}{\|w\|}$$

(assuming there is a data point on the  $\mathbf{w} \cdot \mathbf{x} + b = +1$  or  $-1$  line)

Final result: can maximize  $\gamma$  by minimizing  $\|w\|_2$ !!!

# (Hard margin) support vector machines

$$\text{minimize}_{w,b} \quad w \cdot w$$
$$\left( w \cdot x_j + b \right) y_j \geq 1, \quad \forall j$$



- Example of a **convex optimization** problem
  - A quadratic program
  - Polynomial-time algorithms to solve!
- Hyperplane defined by **support vectors**
  - Could use them as a lower-dimension basis to write down line, although we haven't seen how yet
- More on these later

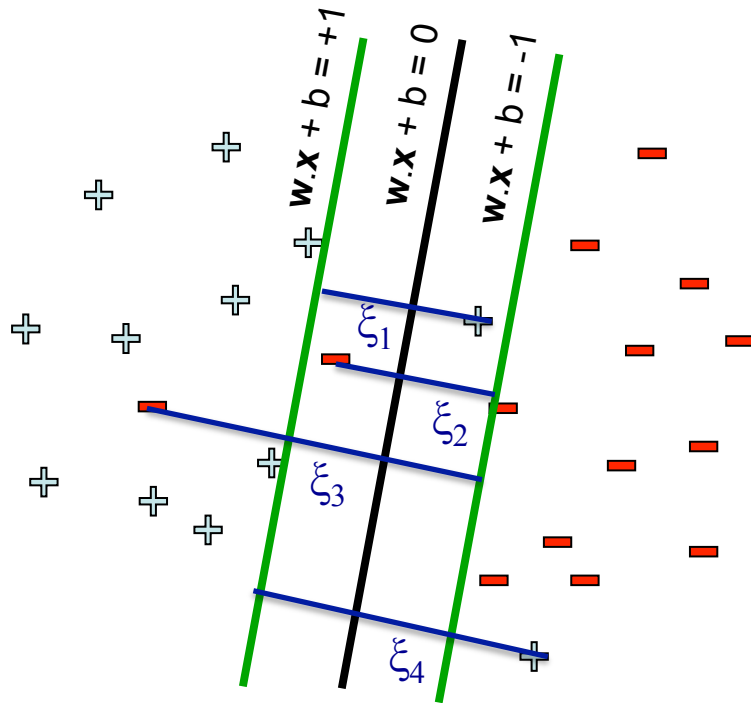
## Non-support Vectors:

- everything else
- moving them will not change  $w$

## Support Vectors:

- data points on the canonical lines

# Allowing for slack: “Soft margin SVM”



$$\text{minimize}_{w,b} \quad w \cdot w + C \sum_j \xi_j$$
$$\left( w \cdot x_j + b \right) y_j \geq 1 - \xi_j, \quad \forall j \quad \xi_j \geq 0$$

↑  
“slack variables”

## Slack penalty $C > 0$ :

- $C = \infty \rightarrow$  have to separate the data!
- $C = 0 \rightarrow$  ignores the data entirely!
- **Select using cross-validation**

For each data point:

- If margin  $\geq 1$ , don't care
- If margin  $< 1$ , pay linear penalty



# Equivalent formulation using hinge loss

$$\begin{aligned} & \text{minimize}_{\mathbf{w}, b} \quad \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ & \left( \mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1 - \xi_j, \quad \forall j \quad \xi_j \geq 0 \end{aligned}$$

Substituting  $\xi_j = \max(0, 1 - (w \cdot x_j + b) y_j)$  into the objective, we get:

$$\min \|w\|^2 + C \sum_j \max(0, 1 - (w \cdot x_j + b) y_j)$$

Recall, the **hinge loss** is  $\ell_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - \hat{y}y)$

$$\min_{\mathbf{w}, b} \|w\|_2^2 + C \sum_j \ell_{\text{hinge}}(y_j, w \cdot x_j + b)$$

This is called **regularization**;  
used to prevent overfitting!

This part is empirical risk minimization,  
using the hinge loss