

# Support vector machines I: Maximum-margin linear classifiers

Sanjoy Dasgupta

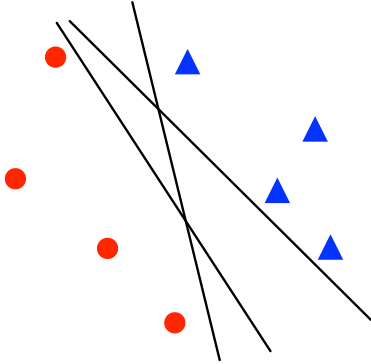
University of California, San Diego

## Topics we'll cover

- ① The margin of a linear classifier
- ② Maximizing the margin
- ③ A convex optimization problem
- ④ Support vectors

## Improving upon the Perceptron

For a linearly separable data set, there are in general many possible separating hyperplanes, and Perceptron is guaranteed to find one of them.



Is there a better, more systematic choice of separator?  
The one with the most buffer around it, for instance?

## The learning problem

Given: training data  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathbb{R}^d \times \{-1, +1\}$ .

Find:  $w \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  such that  $y^{(i)}(w \cdot x^{(i)} + b) > 0$  for all  $i$ .

By scaling  $w, b$ , can equivalently ask for

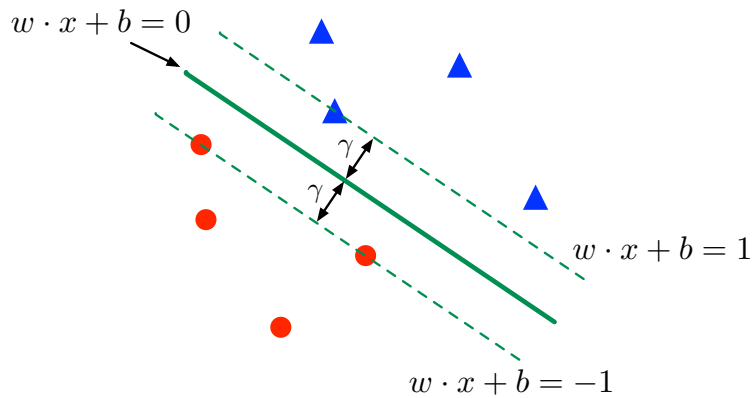
$$y^{(i)}(w \cdot x^{(i)} + b) \geq 1 \quad \text{for all } i$$

## Maximizing the margin

Given: training data  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathbb{R}^d \times \{-1, +1\}$ .

Find:  $w \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  such that

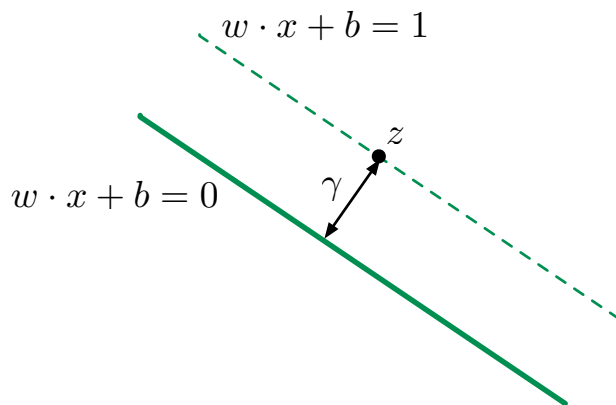
$$y^{(i)}(w \cdot x^{(i)} + b) \geq 1 \quad \text{for all } i.$$



Maximize the **margin**  $\gamma$ .

## A formula for the margin

Close-up of a point  $z$  on the positive boundary.



A quick calculation shows that  $\gamma = 1/\|w\|$ .

In short: to maximize the margin, minimize  $\|w\|$ .

## Maximum-margin linear classifier

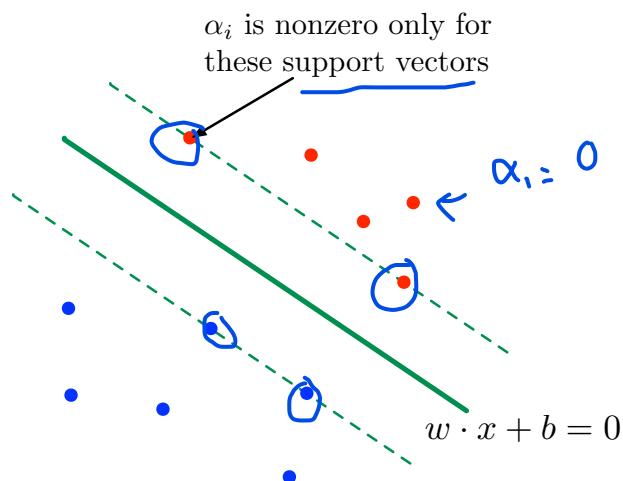
- Given  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathbb{R}^d \times \{-1, +1\}$

$$\begin{aligned} \min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \quad & \|w\|^2 \\ \text{s.t.:} \quad & y^{(i)}(w \cdot x^{(i)} + b) \geq 1 \quad \text{for all } i = 1, 2, \dots, n \end{aligned}$$

- This is a **convex optimization problem**:
  - Convex objective function
  - Linear constraints
- This means that:
  - the optimal solution can be found efficiently
  - duality** gives us information about the solution

## Support vectors

**Support vectors**: training points right on the margin, i.e.  $y^{(i)}(w \cdot x^{(i)} + b) = 1$ .



$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$  is a function of just the support vectors.

$$\alpha_i \geq 0$$

## Small example: Iris data set

Fisher's **iris** data



150 data points from three classes:

- iris setosa
- iris versicolor
- iris virginica

Four measurements: petal width/length, sepal width/length

## Small example: Iris data set

Two features: sepal width, petal width.

Two classes: setosa (red circles), versicolor (black triangles)

