

# Hard-margin SVM

Levent Sagun

New York University

February 11, 2016

## Problem setup

Given a set of linearly separable training data, how can one find a good separator? What do we expect from a good separator?

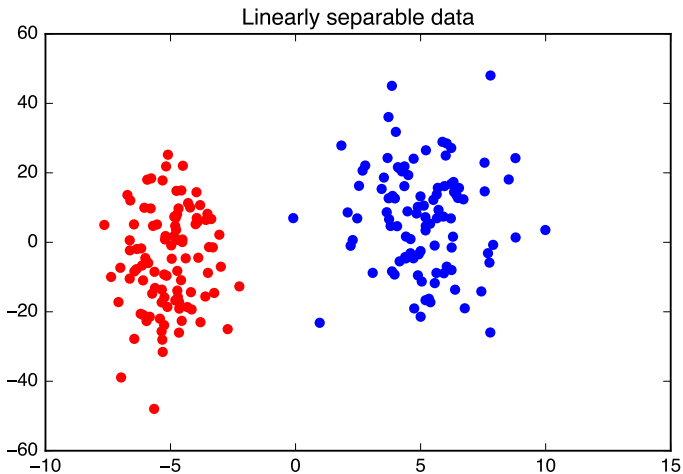
- ... that it actually separates the training points
- ... that it generalizes well

Let  $\{x^i, y^i\}_{i=1}^N \in \mathcal{D}$  be the training data, where  $x^i \in \mathbb{R}$  and  $y^i$  is either  $+1$  or  $-1$ . What does it mean that the data is linearly separable?

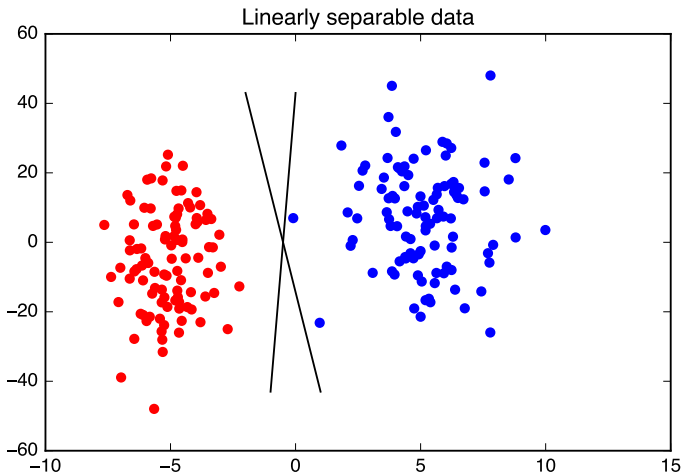
- ... that there is a hyperplane that separates the two clusters
- ... that there is possibly *a lot* of such hyperplanes

How to choose the best one?

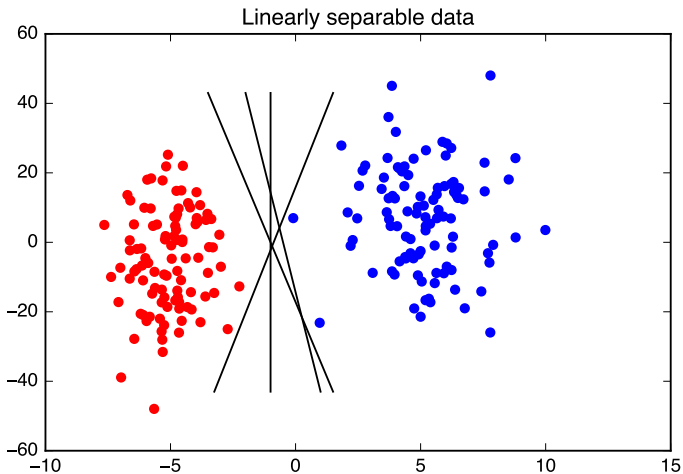
# Example



# Example



# Example



## Hyperplane parametrization

Simplest case of real variables,  $y = mx + b$  draws a line with slope  $m$  that intersects  $y$ -axis at the point  $b$ :

- Rewrite the above equation:  $(m, -1) \cdot (x, y) + b = 0$
- A better notation can be:  $(w_1, w_2) \cdot (x_1, x_2) + b = 0$
- $-w_2/w_1 = m$  captures the connection between the two

## Hyperplane parametrization

Simplest case of real variables,  $y = mx + b$  draws a line with slope  $m$  that intersects  $y$ -axis at the point  $b$ :

- Rewrite the above equation:  $(m, -1) \cdot (x, y) + b = 0$
- A better notation can be:  $(w_1, w_2) \cdot (x_1, x_2) + b = 0$
- $-w_2/w_1 = m$  captures the connection between the two

Generalize this to higher dimensions, for  $w, x \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ :

- $\ell(x) = w \cdot x + b$  where  $\ell(x) = 0$  describes a line.
- $w$  is orthogonal to  $\ell$  (check  $w \cdot v = 0$  for  $v$  along  $\ell$ .)
- What should  $\ell$  assign to the two clusters?

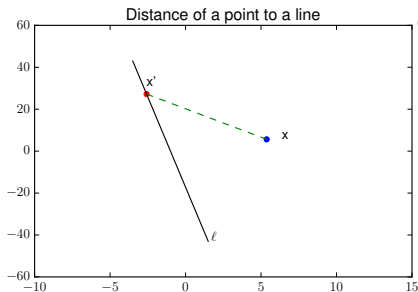
$$\ell(x) \text{ is } \begin{cases} > 0 & \text{if } x \in \text{Blue: } +1 \text{ class} \\ < 0 & \text{if } x \in \text{Red: } -1 \text{ class} \end{cases}$$

- *Observation:*  $y^i \ell(x^i)$  is always positive!

# Distance of a point to a line

For a point  $x \in \mathbb{R}^n$ , how far is  $x$  to a given line  $\ell$ ?

- Denote the distance of a point  $x$  to a line  $\ell$  by  $d(x, \ell)$ .
- Pick a point on the line, say  $x'$ , then  $d(x, \ell)$  is the projection of  $(x - x')$  onto the normal vector  $w$  of  $\ell$ .



Crash course on projections:

- Linear transformations,  $P$ , such that  $P^2 = P$ .
- Unique decomposition into image and kernel of  $P$
- Orthogonal projections:  $P = P^T$
- Vector projection:  $P_w(v) = \frac{v \cdot w}{\|w\|^2} w$



## Hard-margin SVM

Given two linearly separable clusters,  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , and a line  $\ell$  with  $\ell(x) = w \cdot x + b$  and  $\|w\| = 1$ , suppose  $x^{1,\ell} \in \mathcal{C}_1$  and  $x^{2,\ell} \in \mathcal{C}_2$  are the closest points to  $\ell$ .

- For any  $i$ ,  $y^i \ell(x^i) \geq \min\{d(x^{1,\ell}, \ell), d(x^{2,\ell}, \ell)\} > 0$
- **GOAL:** Maximize the *margin*!
- Since data is linearly separable, the maximizer will be on the set where  $d(x^{1,\ell}, \ell) = d(x^{2,\ell}, \ell)$ , let's call this  $M$ .  
(note that  $M$  depends on data points and the line)

## Hard-margin SVM

Given two linearly separable clusters,  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , and a line  $\ell$  with  $\ell(x) = w \cdot x + b$  and  $\|w\| = 1$ , suppose  $x^{1,\ell} \in \mathcal{C}_1$  and  $x^{2,\ell} \in \mathcal{C}_2$  are the closest points to  $\ell$ .

- For any  $i$ ,  $y^i \ell(x^i) \geq \min\{d(x^{1,\ell}, \ell), d(x^{2,\ell}, \ell)\} > 0$
- **GOAL:** Maximize the *margin*!
- Since data is linearly separable, the maximizer will be on the set where  $d(x^{1,\ell}, \ell) = d(x^{2,\ell}, \ell)$ , let's call this  $M$ .  
(note that  $M$  depends on data points and the line)

### Procedure:

$$\max\{M : b \in \mathbb{R}, w \in \mathbb{R}^n, \|w\| = 1\} \quad (1)$$

$$\text{subject to } y^i(w \cdot x^i + b) \geq M \quad (2)$$

## Equivalent formulation

For any pair of  $(w, b)$  we can calculate  $M$  and then considering the new pair  $(w', b') = (\frac{w}{M}, \frac{b}{M})$  we get  $y^i(\frac{w}{M} \cdot x^i + \frac{b}{M}) \geq 1$ . Therefore, maximizing  $M$  can be rephrased as minimizing  $\|w'\|$ .

### New procedure:

$$\min\{\|w'\| : b' \in \mathbb{R}, w' \in \mathbb{R}^n\} \quad (3)$$

$$\text{subject to } y^i(w' \cdot x^i + b') \geq 1 \quad (4)$$

- Note that:  $\|w'\| = \|\frac{w}{M}\| = \frac{\|w\|}{M} = \frac{1}{M}$
- This is a convex optimization problem: quadratic criterion, linear inequality constraints.
- But, what if the clusters overlap?

## Overlapping clusters

For all data points let  $t^i > 0$  be the slack variables. Let's modify equation (2) to allow each point to have a little more room:

$$\text{subject to } y^i(w \cdot x^i + b) \geq M(1 - t^i)$$

Equivalent formulation ... :

$$\min ||w||, \text{ subject to } y^i(w \cdot x^i + b) \geq 1 - t^i \text{ and } \sum t^i < C$$

Equivalent formulation ... :

$$\min \frac{1}{2} ||w||^2 + c \sum t^i$$

$$\text{subject to } y^i(w \cdot x^i + b) \geq 1 - t^i$$

# Overlapping clusters

properties, interpretation, plots...

## Exercises

- *Linear regression*; Minimizing sum of squares of errors in  $y = X\beta + \epsilon$ : Find  $\beta$  such that  $\|y - X\beta\| = f(\beta)$  is minimized.
- What's the orthogonal projection of  $y$  onto the columns of  $X$ ?
- What's the connection of the two?
- When is  $X^T X$  not invertible?
- In the overlapping case, what would happen if you modified the constraint by  $y^i(w \cdot x^i - b) \geq M - t^i$