



# CIS 419/519 Recitation

---

Chang Liu

Nov 11 2020

# Content

---

- VC Dimension
- SVM
- Adaboost

---

# Part I: VC Dimension

# Definitions

---

- We say that a set  $S$  of examples is **shattered** by a set of functions  $H$  if for every partition of the examples in  $S$  into positive and negative examples there is a function in  $H$  that gives exactly these labels to the examples
- The **VC dimension** of hypothesis space  $H$  over instance space  $X$  is the size of the largest finite subset of  $X$  that is shattered by  $H$ .

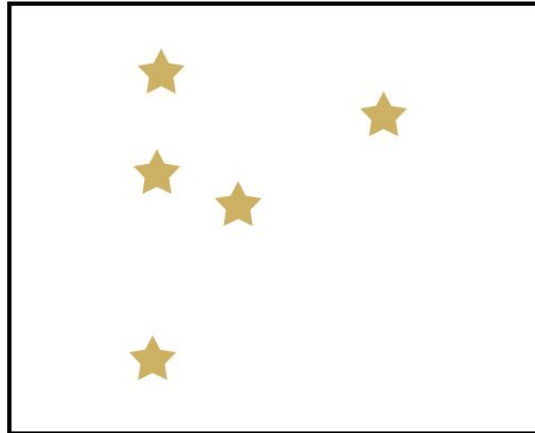
Shatter?

Size of largest finite subset of  $X$ ?

# Walkthrough Example

Setting:

- We have 5 data points scatter randomly in a 2D space
- We propose a linear separator, in Hypothesis Space (H)
- Objective: Find  $VC(H)$



# Strategy Summary:

---

We will use a different strategy:

1. Guess the VC dimension (in this case, we guess 3)
2. Find a set of size 3 that is shattered by  $H$
3. Show that no set of size 4 is shattered by  $H$

This is enough to show that the biggest set shattered by  $H$  has size 3

# Why we want it?

- Help us to figure out how **expressive a Hypothesis Space** is, especially for  $|H| = \text{infinity}$
- Using  $VC(H)$  as a measure of expressiveness, we can get an **Occam algorithm for infinite hypothesis spaces**.
- Given a sample  $D$  of  $m$  examples, find some  $h \in H$  that is **consistent** with all  $m$  examples
- If  $m > \frac{1}{\epsilon} \{8VC(H) \log \frac{13}{\epsilon} + 4 \log \left(\frac{2}{\delta}\right)\}$
- Then with probability at least  $(1 - \delta)$ ,  $h$  has error less than  $\epsilon$ . (that is, if  $m$  is polynomial we have a **PAC** learning algorithm; to be efficient, we need to produce the hypothesis  $h$  efficiently.)

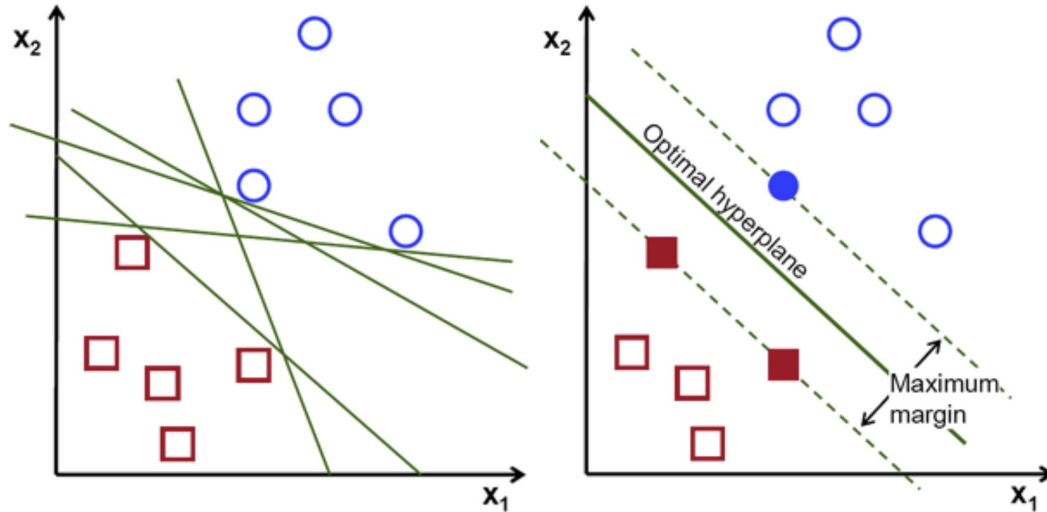
---

# Part 2: SVM

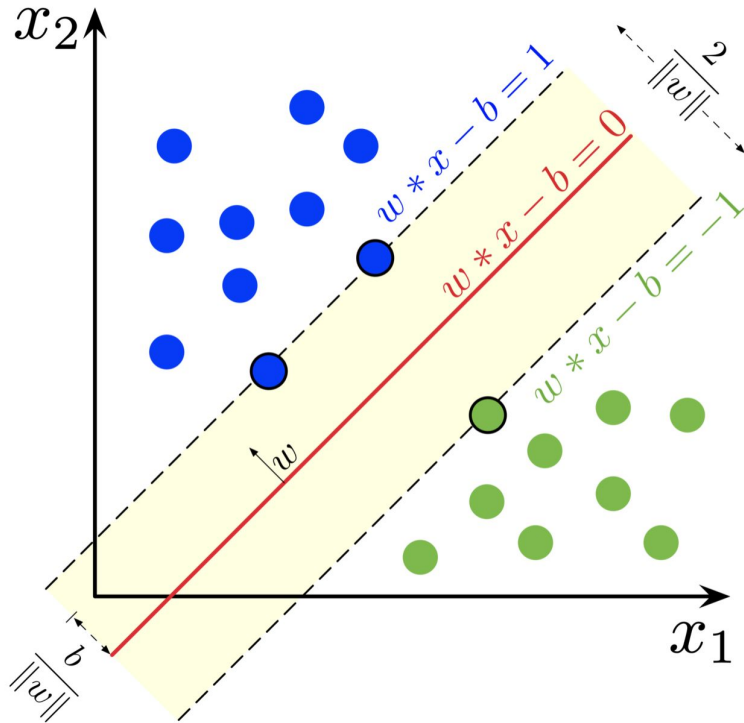


# What is Support Vector Machine?

- Finding the separator that maximize the margin of 2 sets of data



# Margin

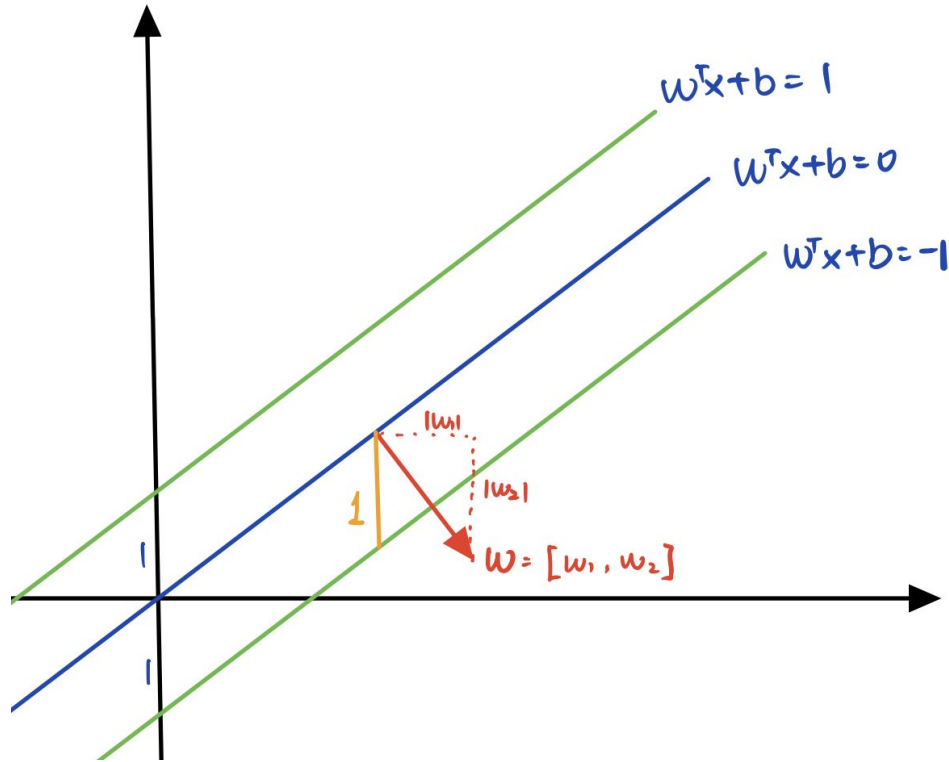


$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

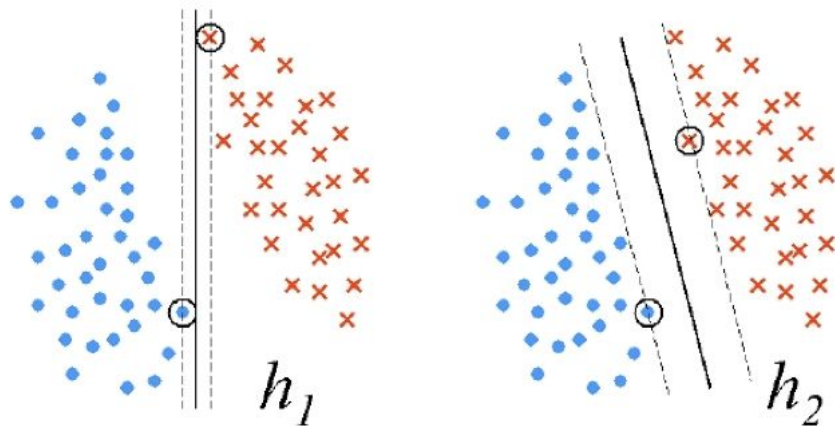
$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall (\mathbf{x}_i, y_i) \in S$$

- Why Margin =  $1/\|\mathbf{w}\|$ ?
- Where does it come from?

# Margin-Geometry Perspective



# Hard SVM



The margin of a linear separator

$$\mathbf{w}^T \mathbf{x} + b = 0 \text{ is } \frac{1}{\|\mathbf{w}\|}$$

$$\begin{aligned} \max \frac{1}{\|\mathbf{w}\|} &= \min \|\mathbf{w}\| \\ &= \min \frac{1}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

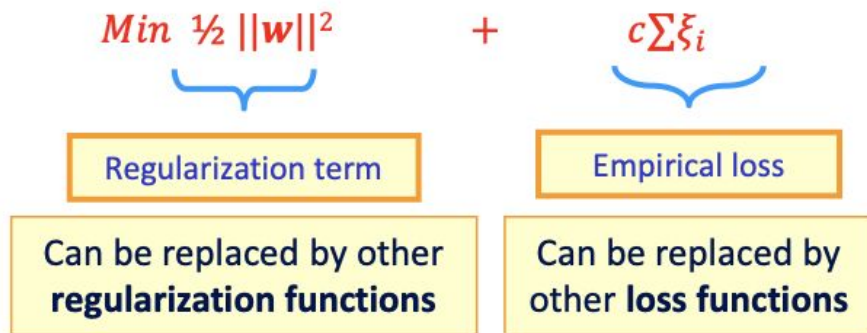
$$\text{s.t. } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall (\mathbf{x}_i, y_i) \in S$$

# Soft SVM

- The problem we solved is:

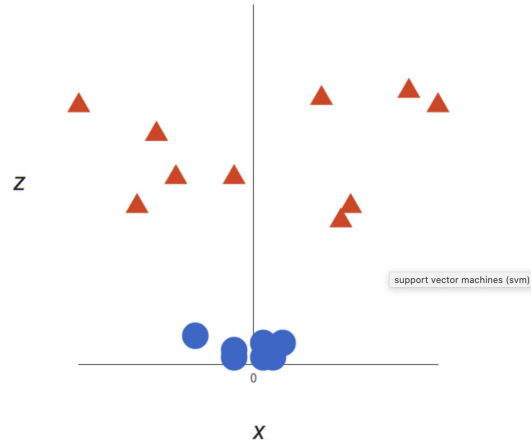
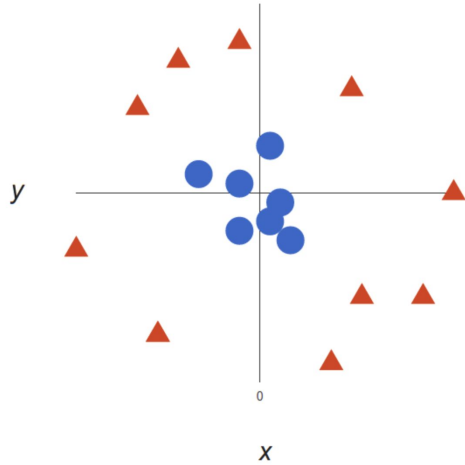
$$\text{Min } \frac{1}{2} \|\mathbf{w}\|^2 + c \sum \xi_i$$

- Where  $\xi_i > 0$  is called a **slack variable**, and is defined by:
  - $\xi_i = \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$
  - Equivalently, we can say that:  $y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i; \xi_i \geq 0$
- And this can be written as:

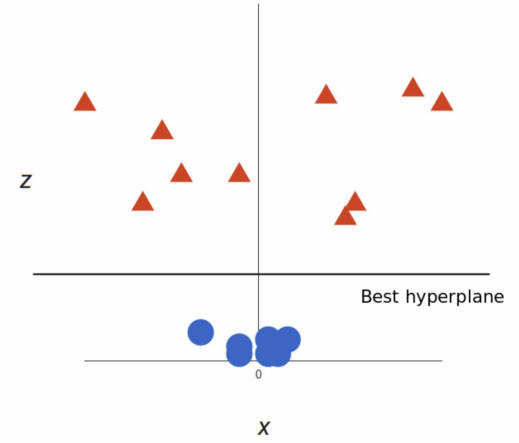


- General Form of a learning algorithm:
  - Minimize empirical loss, and Regularize (to avoid over fitting)
  - Theoretically motivated improvement over the original algorithm we've seen at the beginning of the semester.

# Nonlinear SVM & Kernels



support vector machines (svm)



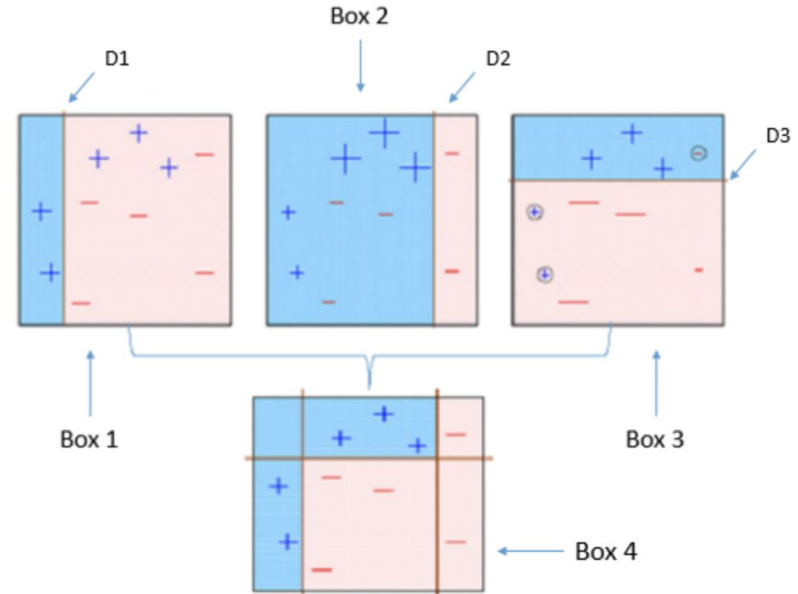
---

# Part 3: Adaboost

More Intuition

# What is Adaboost?

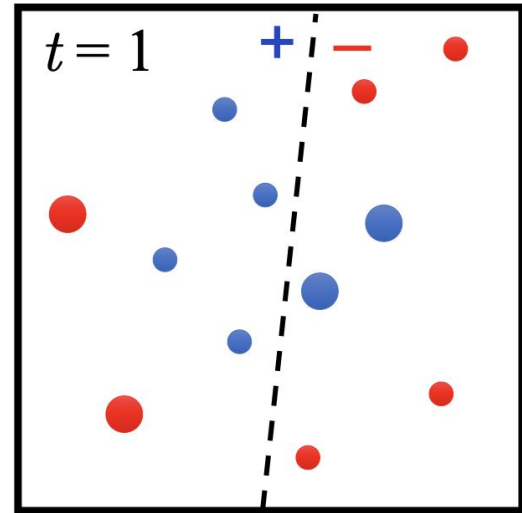
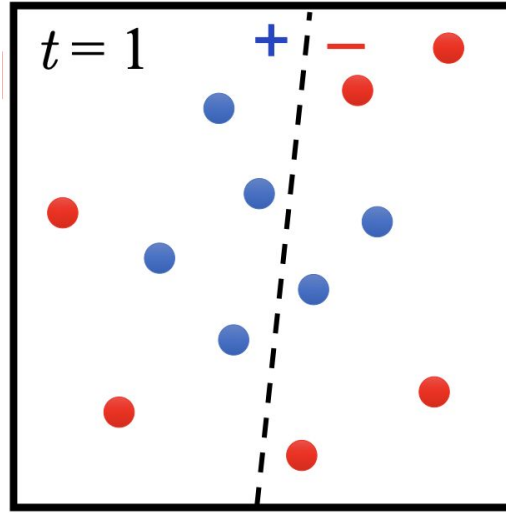
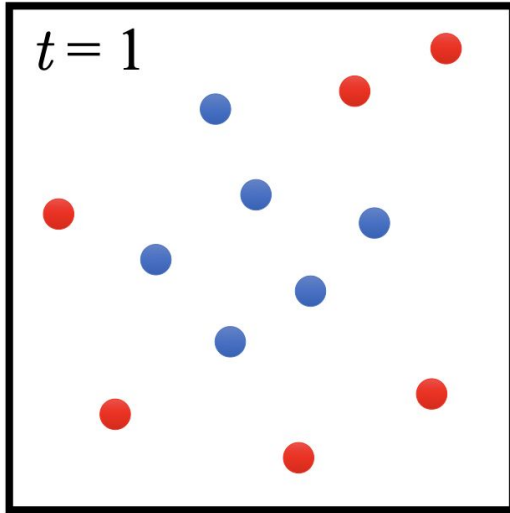
- A predictors which trains sequentially, each trying to correct its predecessor.
- **Method:** set weights to both classifiers and data points in a way that forces classifiers to concentrate on observations that are difficult to correctly classify.
- Helps combine multiple “**weak classifiers**” into a **single “strong classifier”**”



Adaboost with DT Stump



# Adaboost with Linear Separator



copyright©2020 Eric Eaton

# Adaboost with Linear Separator

- Constructing  $D_t$  on  $\{1, \dots, m\}$ :

- $D_1(i) = 1/m$

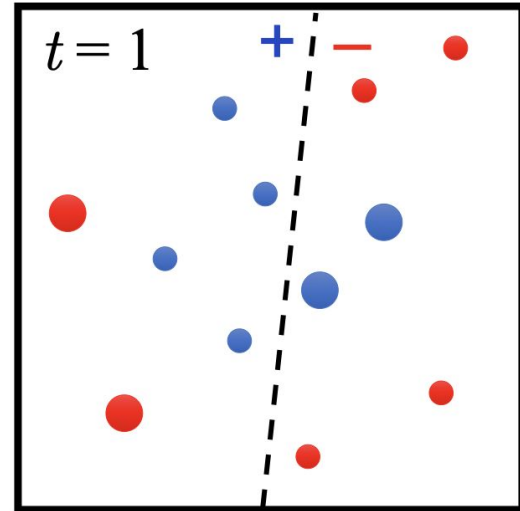
- Given  $D_t$  and  $h_t$  :

- $D_{t+1} = \begin{cases} D_t(i)/z_t \times e^{-\alpha_t} & \text{if } y_i = h_t(\mathbf{x}_i) \\ D_t(i)/z_t \times e^{+\alpha_t} & \text{if } y_i \neq h_t(\mathbf{x}_i) \end{cases}$

$$= \frac{D_t(i)}{z_t} \times \exp(-\alpha_t y_i h_t(\mathbf{x}_i))$$

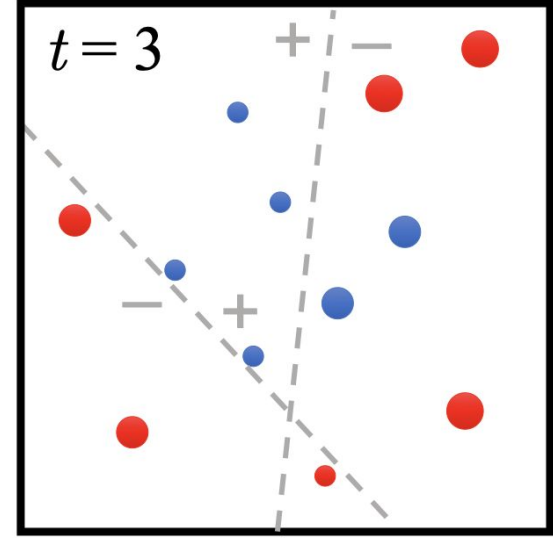
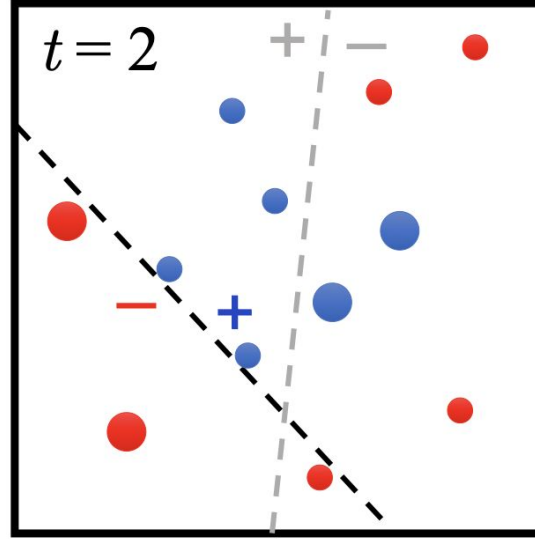
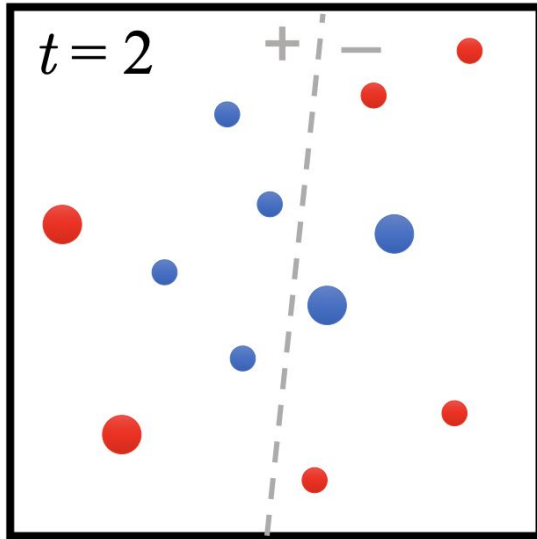
where  $z_t =$  normalization constant

and  $\alpha_t = \frac{1}{2} \ln\{ (1 - \varepsilon_t)/\varepsilon_t \}$



- Final hypothesis:  $H_{final}(\mathbf{x}) = \text{sign}(\sum_t \alpha_t h_t(\mathbf{x}))$

# Adaboost with Linear Separator



# Adaboost with Linear Separator

