

# 10-701: Introduction to Machine Learning

## Lecture 22 –Boosting

Hoda Heidari

\* Slides adopted from F24 offering of 10701 by Henry Chai.

# Bias-Variance Decomposition

- Suppose  $y=f(x)+\varepsilon$ , where  $\varepsilon$  is noise with  $\mathbb{E}[\varepsilon] = 0, \text{Var}(\varepsilon) = \sigma^2$ .
- The expected squared prediction error at a point  $x$ :

$$\begin{aligned}\mathbb{E}_{D,\varepsilon}[(y - \hat{f}(x))^2] &= \\ &= \underbrace{(f(x) - \mathbb{E}[\hat{f}(x)])^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]}_{\text{Variance}} + \sigma^2\end{aligned}$$

# Ensemble Learning

- **Key idea:** Different models make different errors. By aggregating their predictions, ensembles can reduce variance, bias, or both—leading to better accuracy and robustness than any single model alone.
- **Common ensemble methods**
  - **Bagging (Bootstrap Aggregating):**  
Trains many models independently on different random subsets of the data to reduce variance and overfitting.
  - **Boosting:**  
Trains models sequentially, with each new model focusing on the errors of the previous one to reduce bias.
  - **Stacking:**  
Combines predictions from multiple models using another model (a *meta-learner*) that learns how best to blend them.

# Decision Trees: Pros & Cons

- Pros
  - Interpretable
  - Efficient (computational cost and storage)
  - ...
- Cons
  - Learned greedily: each split only considers the immediate impact on the splitting criterion
  - Prone to overfit
  - High variance
    - Can be addressed via **bagging** → random forests
  - Limited expressivity (especially short trees, i.e., stumps)
    - Can be addressed via **boosting**

# AdaBoost

- Intuition: iteratively reweight inputs, giving more weight to inputs that are difficult-to-predict correctly
- Analogy:
  - You all have to take a test (😱) ...
  - ... but you're going to be taking it one at a time.
  - After you finish, you get to tell the next person the questions you struggled with.
  - Hopefully, they can cover for you because...
  - ... if “enough” of you get a question right, you'll all receive full credit for that problem

- Input:  $\mathcal{D} \left( y^{(n)} \in \{-1, +1\} \right), T$
- Initialize data point weights:  $\omega_0^{(1)}, \dots, \omega_0^{(N)} = \frac{1}{N}$
- For  $t = 1, \dots, T$ 
  1. Train a weak learner,  $h_t$ , by minimizing the *weighted* training error
  2. Compute the *weighted* training error of  $h_t$ :

$$\epsilon_t = \sum_{n=1}^N \omega_{t-1}^{(n)} \mathbb{1} \left( y^{(n)} \neq h_t(\mathbf{x}^{(n)}) \right)$$

3. Compute the **importance** of  $h_t$ :

$$\alpha_t = \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

4. Update the data point weights:

$$\omega_t^{(n)} = \frac{\omega_{t-1}^{(n)}}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(\mathbf{x}^{(n)}) = y^{(n)} \\ e^{\alpha_t} & \text{if } h_t(\mathbf{x}^{(n)}) \neq y^{(n)} \end{cases} = \frac{\omega_{t-1}^{(n)} e^{-\alpha_t y^{(n)} h_t(\mathbf{x}^{(n)})}}{Z_t}$$

- Output: an aggregated hypothesis

$$g_T(\mathbf{x}) = \text{sign}(H_T(\mathbf{x})) \\ = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$$

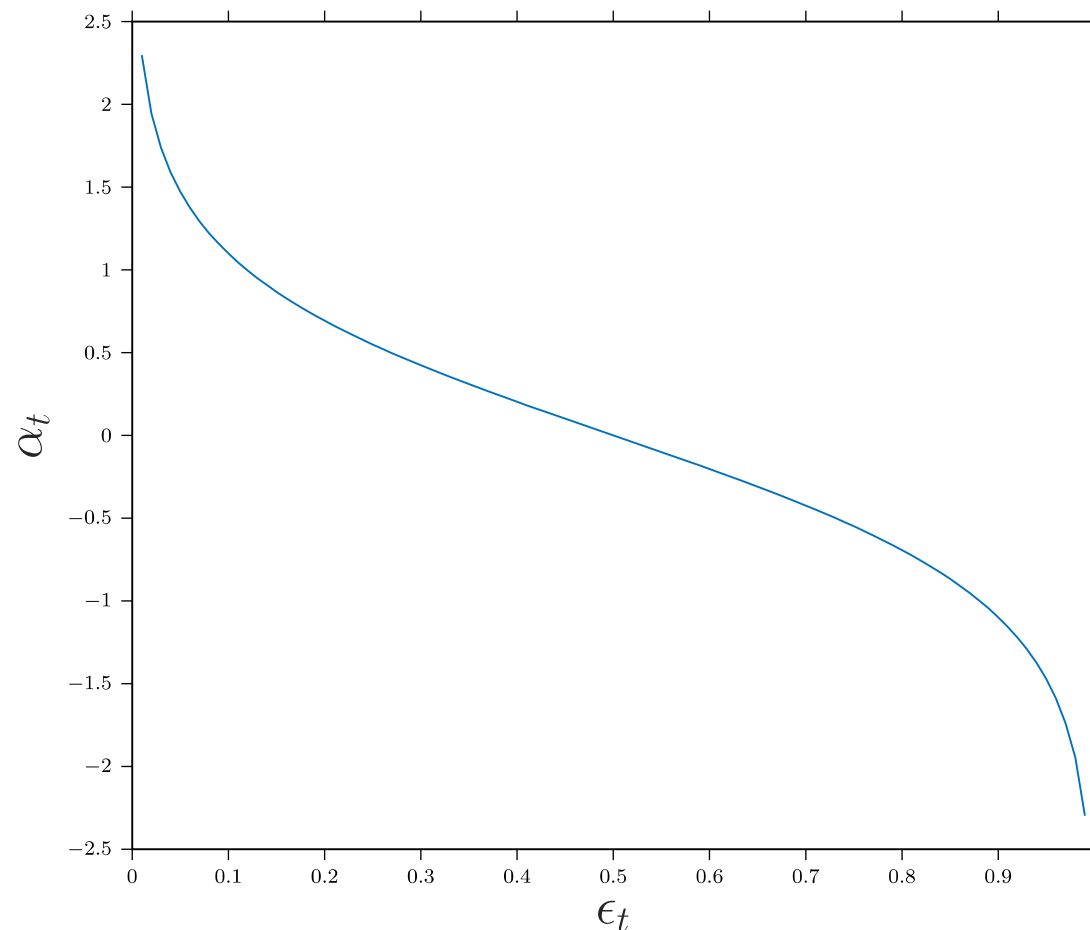
# Setting $\alpha_t$

$\alpha_t$  determines the contribution of  $h_t$  to the final, aggregated hypothesis:

$$g(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$

Intuition: we want good weak learners to have high importances

$$\alpha_t = \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$



## Updating $\omega^{(n)}$

- Intuition: we want incorrectly classified inputs to receive a higher weight in the next round

$$\omega_t^{(n)} = \frac{\omega_{t-1}^{(n)}}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(\mathbf{x}^{(n)}) = y^{(n)} \\ e^{\alpha_t} & \text{if } h_t(\mathbf{x}^{(n)}) \neq y^{(n)} \end{cases}$$

- If  $\epsilon_t < \frac{1}{2}$ , ....



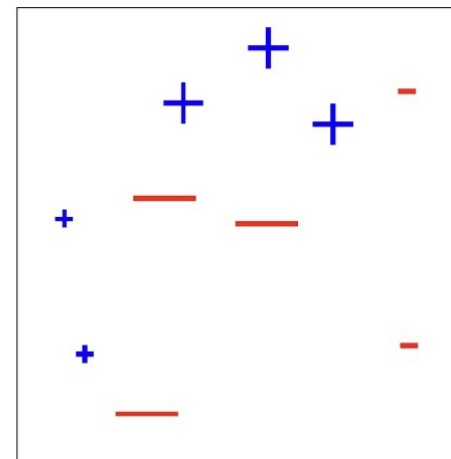
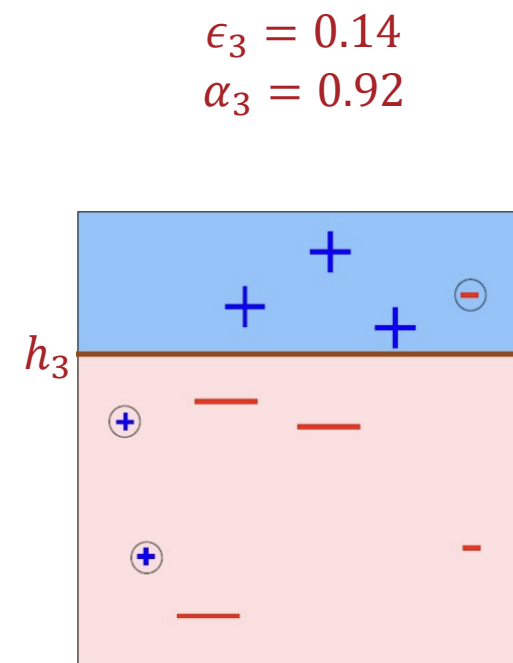
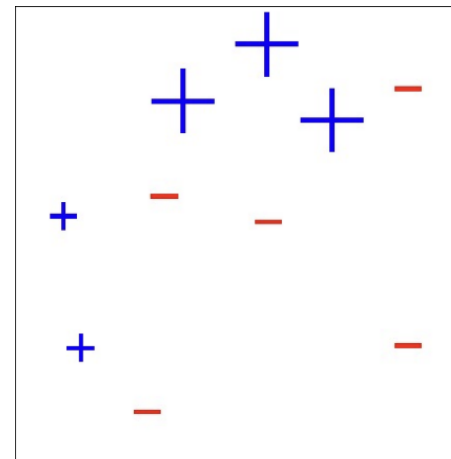
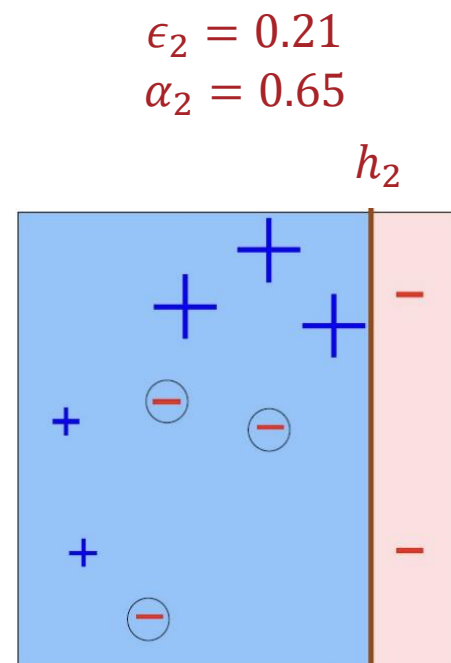
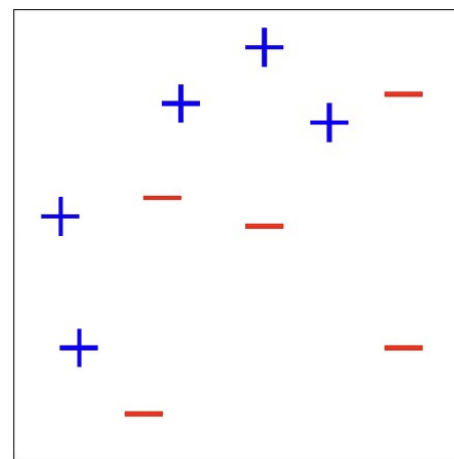
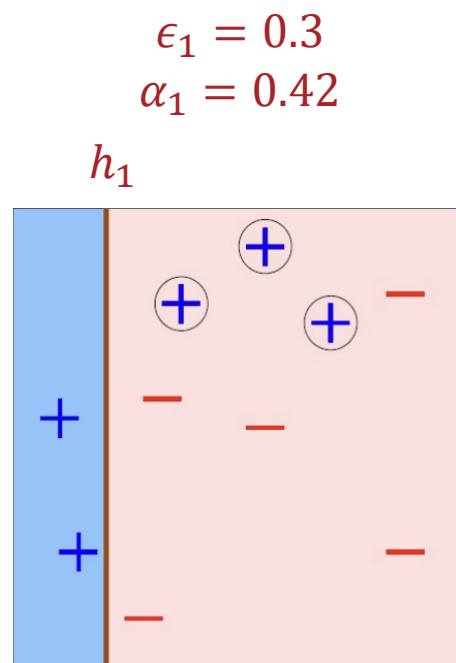
## Updating $\omega^{(n)}$

- Intuition: we want incorrectly classified inputs to receive a higher weight in the next round

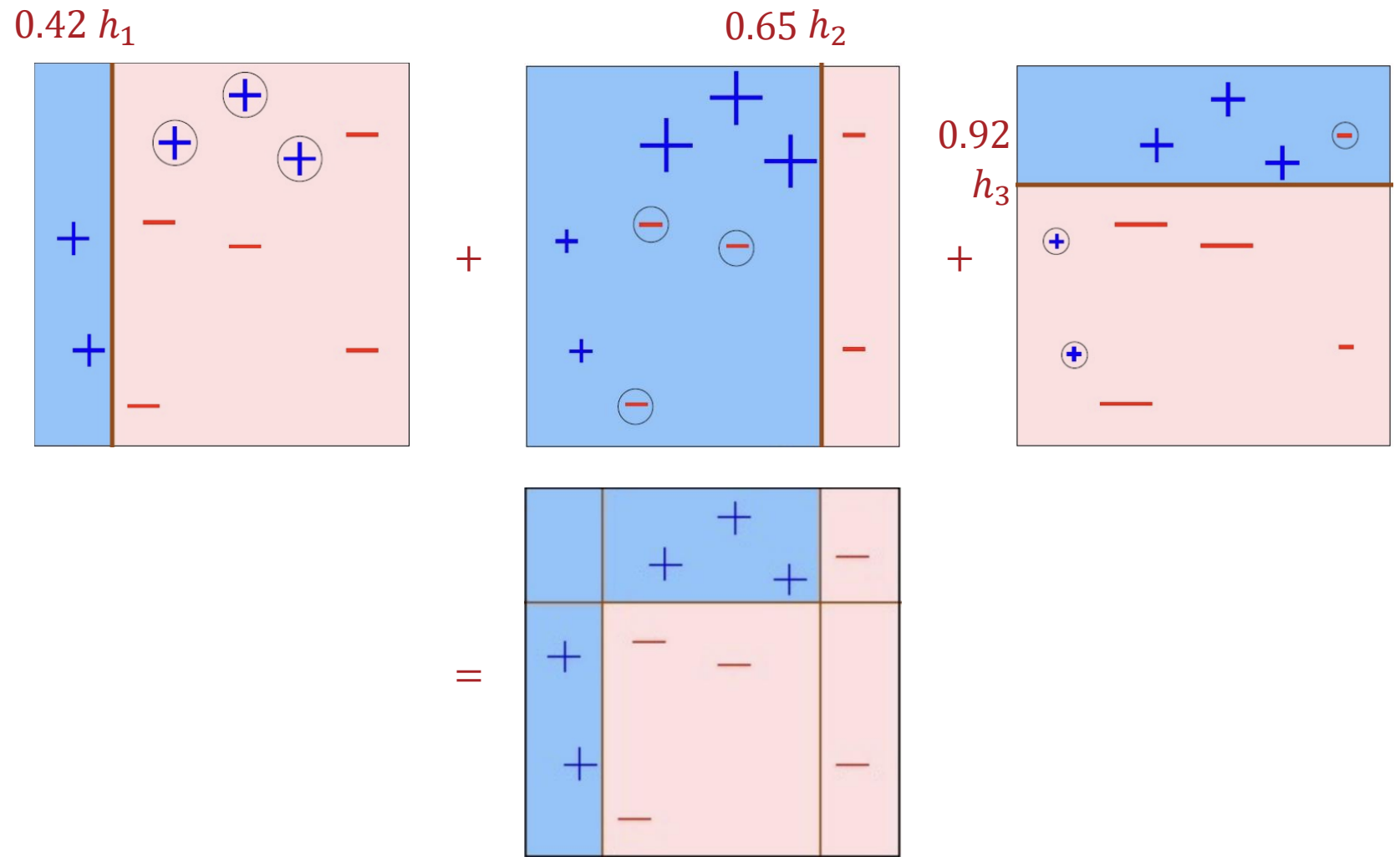
$$\omega_t^{(n)} = \frac{\omega_{t-1}^{(n)}}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(\mathbf{x}^{(n)}) = y^{(n)} \\ e^{\alpha_t} & \text{if } h_t(\mathbf{x}^{(n)}) \neq y^{(n)} \end{cases} = \frac{\omega_{t-1}^{(n)} e^{-\alpha_t y^{(n)} h_t(\mathbf{x}^{(n)})}}{Z_t}$$

- If  $\epsilon_t < \frac{1}{2}$ , then  $\frac{1-\epsilon_t}{\epsilon_t} > 1$
- If  $\frac{1-\epsilon_t}{\epsilon_t} > 1$ , then  $\alpha_t = \frac{1}{2} \log \left( \frac{1-\epsilon_t}{\epsilon_t} \right) > 0$
- If  $\alpha_t > 0$ , then  $e^{-\alpha_t} < 1$  and  $e^{\alpha_t} > 1$

# AdaBoost: Example



# AdaBoost: Example



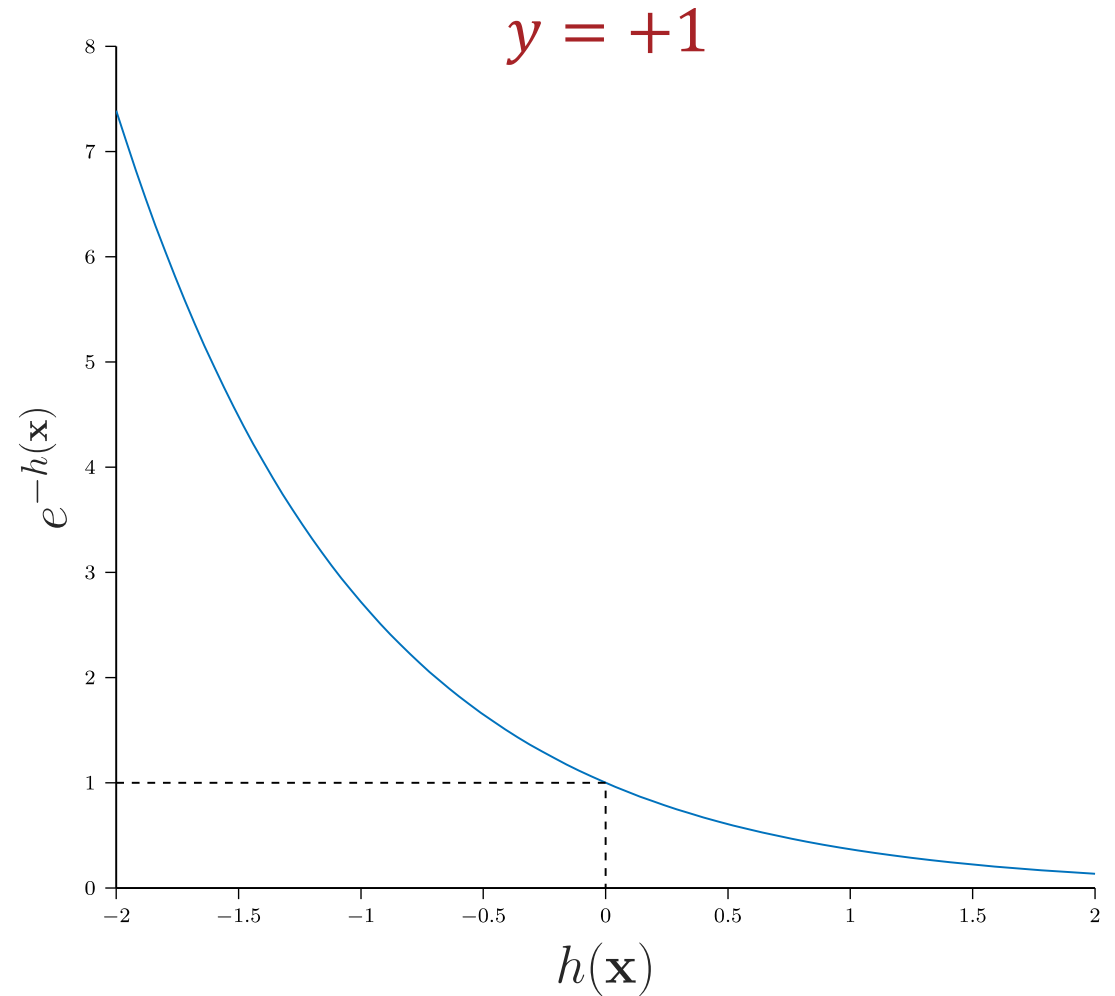
# Why AdaBoost?

1. If you want to use weak learners ...
  2. ... and want your final hypothesis to be a weighted combination of weak learners, ...
  3. ... then Adaboost greedily minimizes the exponential loss:  
$$e(h(\mathbf{x}), y) = e^{-yh(\mathbf{x})}$$
1. Because they're low variance / computational constraints
  2. Because weak learners are not great on their own
  3. Because the exponential loss upper bounds binary error!

# Exponential Loss

$$e(h(\mathbf{x}), y) = e^{-yh(\mathbf{x})}$$

The more  $h(\mathbf{x})$  “agrees with”  $y$ , the smaller the loss and the more  $h(\mathbf{x})$  “disagrees with”  $y$ , the greater the loss



# Exponential Loss

- Claim:

$$\frac{1}{N} \sum_{n=1}^N e^{(-y^{(n)} h(\mathbf{x}^{(n)}))} \geq \frac{1}{N} \sum_{n=1}^N \mathbb{1} \left( \text{sign} \left( h(\mathbf{x}^{(n)}) \right) \neq y^{(n)} \right)$$

- Consequence:

$$\frac{1}{N} \sum_{n=1}^N e^{(-y^{(n)} h(\mathbf{x}^{(n)}))} \rightarrow 0$$

$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \mathbb{1} \left( \text{sign} \left( h(\mathbf{x}^{(n)}) \right) \neq y^{(n)} \right) \rightarrow 0$$

# Exponential Loss

- Claim: if  $g_T = \text{sign}(H_T)$  is the Adaboost hypothesis, then

$$\frac{1}{N} \sum_{n=1}^N e^{(-y^{(n)} H_T(x^{(n)}))} = \prod_{t=1}^T Z_t$$

- Proof:

# Exponential Loss

- Claim: if  $g_T = \text{sign}(H_T)$  is the Adaboost hypothesis, then

$$\frac{1}{N} \sum_{n=1}^N e^{(-y^{(n)} H_T(x^{(n)}))} = \prod_{t=1}^T Z_t$$

- Proof:

$$\omega_0^{(n)} = \frac{1}{N}, \omega_1^{(n)} = \frac{e^{-\alpha_1 y^{(n)} h_1(x^{(n)})}}{N Z_1}, \omega_2^{(n)} = \frac{e^{-\alpha_1 y^{(n)} h_1(x^{(n)})} e^{-\alpha_2 y^{(n)} h_2(x^{(n)})}}{N Z_1 Z_2}$$

$$\omega_T^{(n)} = \frac{\prod_{t=1}^T e^{-\alpha_t y^{(n)} h_t(x^{(n)})}}{N \prod_{t=1}^T Z_t} = \frac{e^{-y^{(n)} \sum_{t=1}^T \alpha_t h_t(x^{(n)})}}{N \prod_{t=1}^T Z_t} = \frac{e^{-y^{(n)} H_T(x^{(n)})}}{N \prod_{t=1}^T Z_t}$$

$$\sum_{n=1}^N \omega_T^{(n)} = \sum_{n=1}^N \frac{e^{-y^{(n)} H_T(x^{(n)})}}{N \prod_{t=1}^T Z_t} = 1 \Rightarrow \frac{1}{N} \sum_{n=1}^N e^{-y^{(n)} H_T(x^{(n)})} = \prod_{t=1}^T Z_t \blacksquare$$



# Exponential Loss

- Claim: if  $g_T = \text{sign}(H_T)$  is the Adaboost hypothesis, then

$$\frac{1}{N} \sum_{n=1}^N e^{(-y^{(n)} H_T(x^{(n)}))} = \prod_{t=1}^T Z_t$$

- Consequence: one way to minimize the exponential training loss is to greedily minimize  $Z_t$ , i.e., in each iteration, make the normalization constant as small as possible by tuning  $\alpha_t$ .

# Greedy Exponential Loss Minimization

$$Z_t = \sum_{n=1}^N \omega_{t-1}^{(n)} e^{-(a) y^{(n)} h_t(x^{(n)})}$$

# Greedy Exponential Loss Minimization

$$\begin{aligned} Z_t &= \sum_{n=1}^N \omega_{t-1}^{(n)} e^{-(a)y^{(n)}h_t(x^{(n)})} \\ &= \sum_{y^{(n)}=h_t(x^{(n)})} \omega_{t-1}^{(n)} e^{-(a)} + \sum_{y^{(n)} \neq h_t(x^{(n)})} \omega_{t-1}^{(n)} e^{(a)} \\ &= e^{-(a)} \sum_{y^{(n)}=h_t(x^{(n)})} \omega_{t-1}^{(n)} + e^{(a)} \sum_{y^{(n)} \neq h_t(x^{(n)})} \omega_{t-1}^{(n)} \\ &= e^{-a}(1 - \epsilon_t) + e^a \epsilon_t \end{aligned}$$

# Greedy Exponential Loss Minimization

$$Z_t = e^{-a}(1 - \epsilon_t) + e^a \epsilon_t$$

# Greedy Exponential Loss Minimization

$$Z_t = e^{-a}(1 - \epsilon_t) + e^a \epsilon_t$$

$$\frac{\partial Z_t}{\partial a} = -e^{-a}(1 - \epsilon_t) + e^a \epsilon_t \Rightarrow -e^{-\hat{a}}(1 - \epsilon_t) + e^{\hat{a}} \epsilon_t = 0$$

$$\Rightarrow e^{\hat{a}} \epsilon_t = e^{-\hat{a}}(1 - \epsilon_t)$$

$$\Rightarrow e^{2\hat{a}} = \frac{1 - \epsilon_t}{\epsilon_t}$$

$$\Rightarrow \hat{a} = \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) = \alpha_t$$

$$\frac{\partial^2 Z_t}{\partial a^2} = e^{-a}(1 - \epsilon_t) + e^a \epsilon_t > 0$$

# Normalizing $\omega^{(n)}$

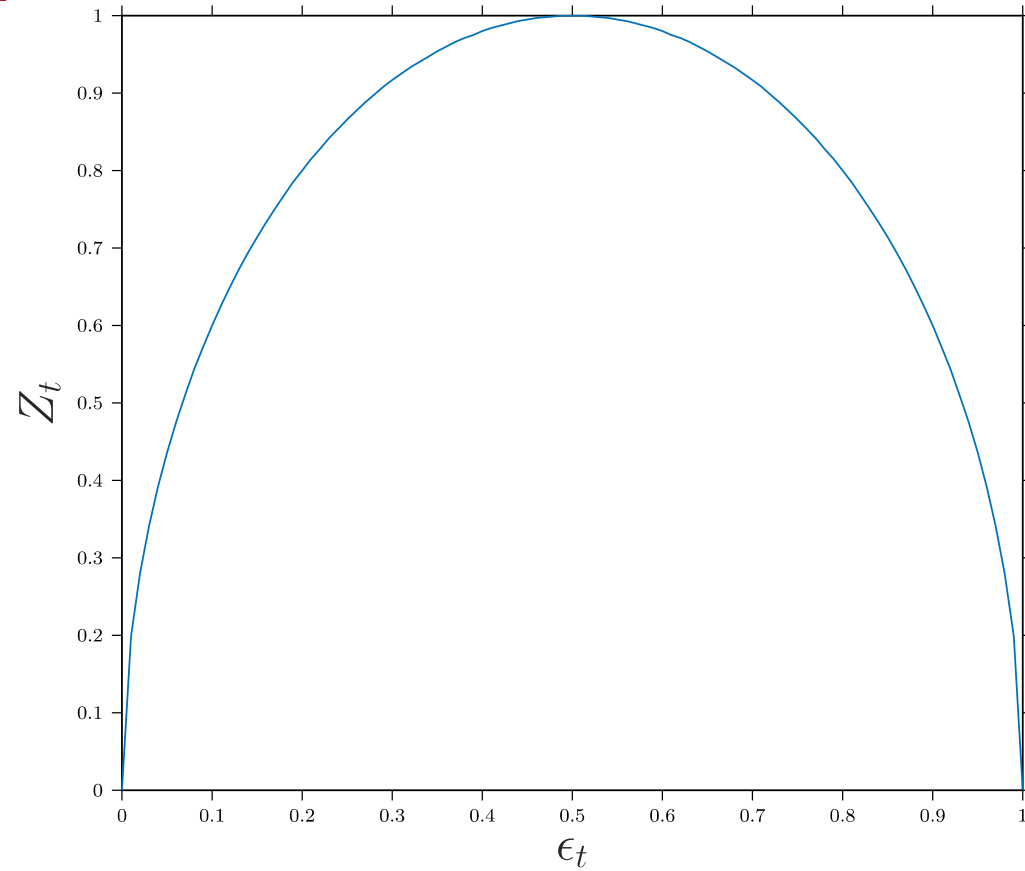
$$\begin{aligned} Z_t &= \sum_{n=1}^N \omega_{t-1}^{(n)} e^{-\alpha_t y^{(n)} h_t(x^{(n)})} \\ &= \sum_{y^{(n)}=h_t(x^{(n)})} \omega_{t-1}^{(n)} e^{-\alpha_t} + \sum_{y^{(n)} \neq h_t(x^{(n)})} \omega_{t-1}^{(n)} e^{\alpha_t} \end{aligned}$$

# Normalizing $\omega^{(n)}$

$$\begin{aligned} Z_t &= \sum_{n=1}^N \omega_{t-1}^{(n)} e^{-\alpha_t y^{(n)} h_t(x^{(n)})} \\ &= \sum_{y^{(n)}=h_t(x^{(n)})} \omega_{t-1}^{(n)} e^{-\alpha_t} + \sum_{y^{(n)} \neq h_t(x^{(n)})} \omega_{t-1}^{(n)} e^{\alpha_t} \\ &= e^{-\alpha_t} \sum_{y^{(n)}=h_t(x^{(n)})} \omega_{t-1}^{(n)} + e^{\alpha_t} \sum_{y^{(n)} \neq h_t(x^{(n)})} \omega_{t-1}^{(n)} \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \\ &= e^{-\frac{1}{2} \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)} (1 - \epsilon_t) + e^{\frac{1}{2} \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)} \epsilon_t \\ &= \sqrt{\epsilon_t (1 - \epsilon_t)} + \sqrt{\epsilon_t (1 - \epsilon_t)} = 2\sqrt{\epsilon_t (1 - \epsilon_t)} \end{aligned}$$

$Z_t$

$$Z_t = \sum_{n=1}^N \omega_{t-1}^{(n)} e^{-\alpha_t y^{(n)} h_t(x^{(n)})} = 2\sqrt{\epsilon_t(1-\epsilon_t)} < 1 \text{ if } \epsilon_t < \frac{1}{2}$$





# Training Error

$$\begin{aligned}\frac{1}{N} \sum_{n=1}^N \mathbb{1} \left( y^{(n)} \neq g_T(\mathbf{x}^{(n)}) \right) &\leq \frac{1}{N} \sum_{n=1}^N e^{(-y^{(n)} H_T(\mathbf{x}^{(n)}))} \\ &= \prod_{t=1}^T Z_t \\ &= \prod_{t=1}^T 2\sqrt{\epsilon_t(1 - \epsilon_t)} \rightarrow 0 \text{ as } T \rightarrow \infty \\ &\left( \text{as long as } \epsilon_t < \frac{1}{2} \ \forall t \right)\end{aligned}$$

# True Error (Freund & Schapire, 1995)

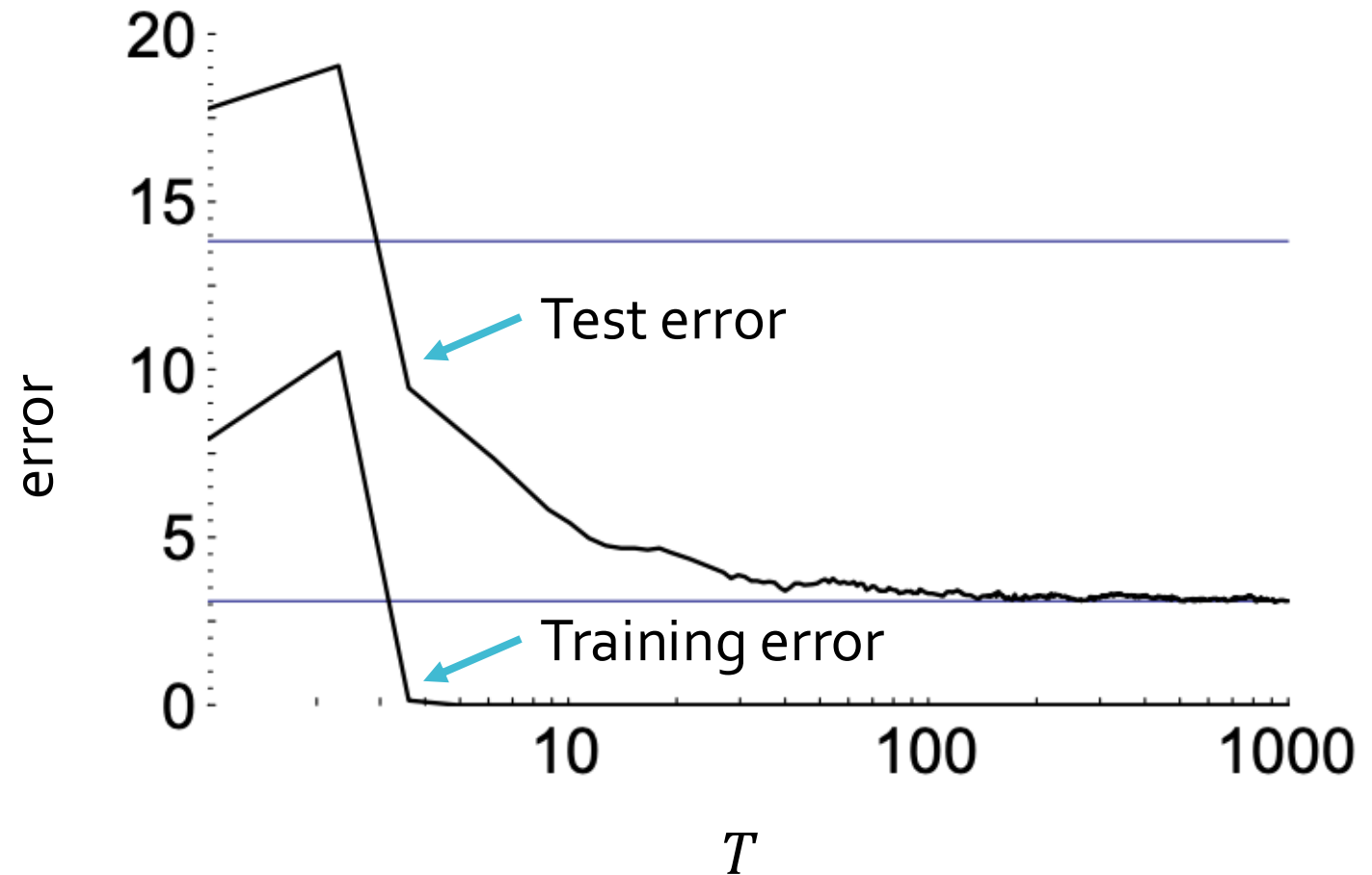
- For AdaBoost, with high probability:

$$\text{True Error} \leq \text{Training Error} + \tilde{O} \left( \sqrt{\frac{d_{vc}(\mathcal{H})T}{N}} \right)$$

where  $d_{vc}(\mathcal{H})$  is the VC-dimension of the weak learners and  $T$  is the number of weak learners.

- Empirical results indicate that increasing  $T$  does not lead to overfitting as this bound would suggest!

# Test Error (Schapire, 1989)

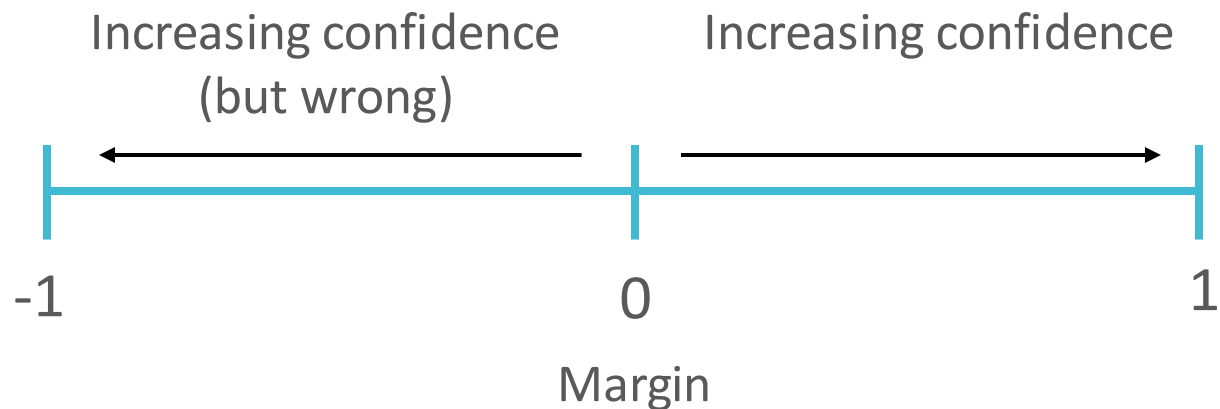


# Margins

- The *margin* of training point  $(\mathbf{x}^{(i)}, y^{(i)})$  is defined as:

$$m(\mathbf{x}^{(i)}, y^{(i)}) = \frac{y^{(i)} \sum_{t=1}^T \alpha_t h_t(\mathbf{x}^{(i)})}{\sum_{t=1}^T \alpha_t}$$

- The margin can be interpreted as how confident  $g_T$  is in its prediction: the bigger the margin, the more confident.



# True Error (Schapire, Freund et al., 1998)

- For AdaBoost, with high probability:

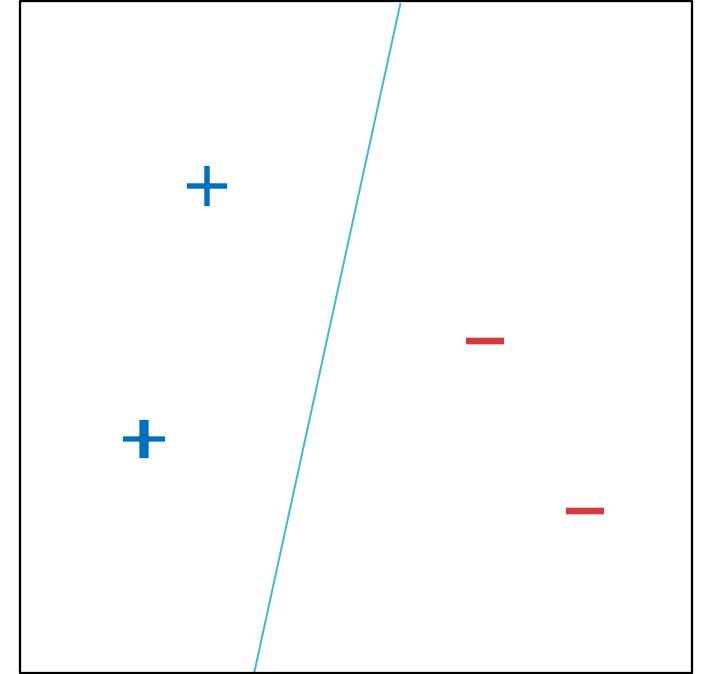
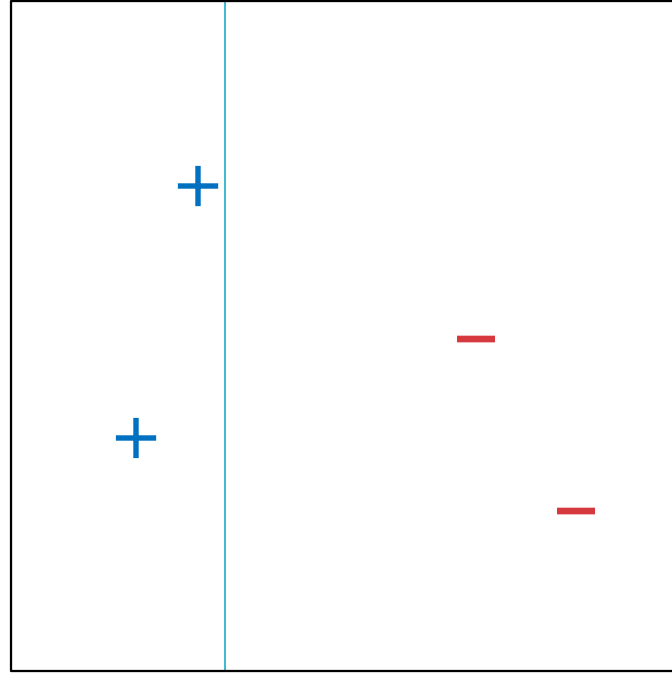
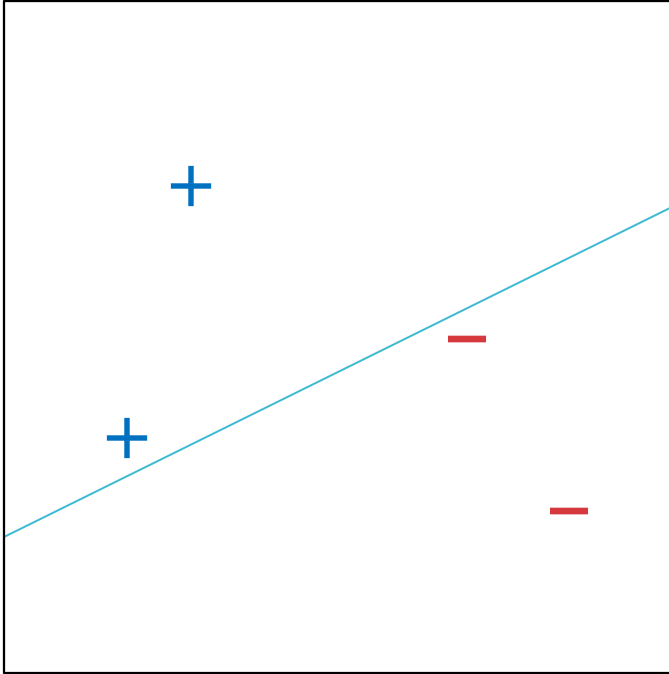
$$\text{True Error} \leq \frac{1}{N} \sum_{i=1}^N \mathbb{I}[m(\mathbf{x}^{(i)}, y^{(i)}) \leq \epsilon] + \tilde{O} \left( \sqrt{\frac{d_{vc}(\mathcal{H})}{N\epsilon^2}} \right)$$

where  $d_{vc}(\mathcal{H})$  is the VC-dimension of the weak learners and  $\epsilon > 0$  is a tolerance parameter.

- Even after AdaBoost has driven the training error to 0, it continues to target the “training margin”

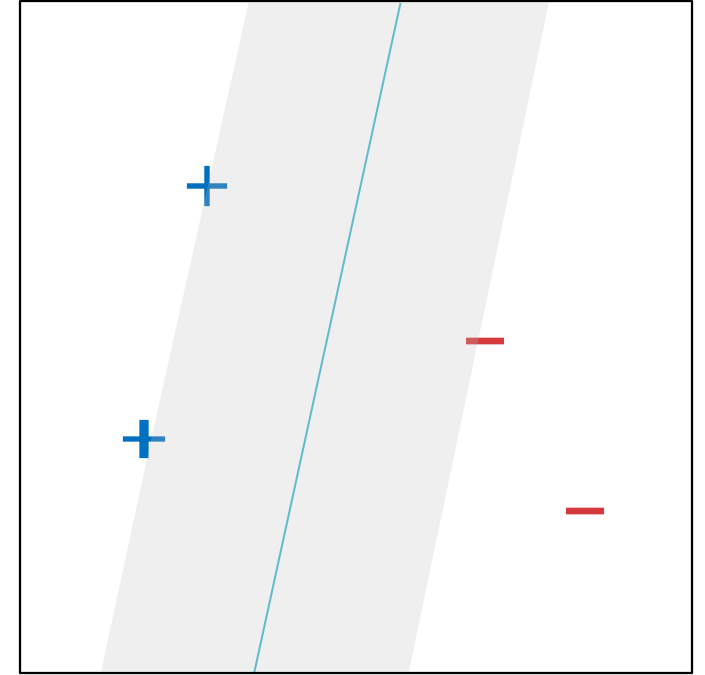
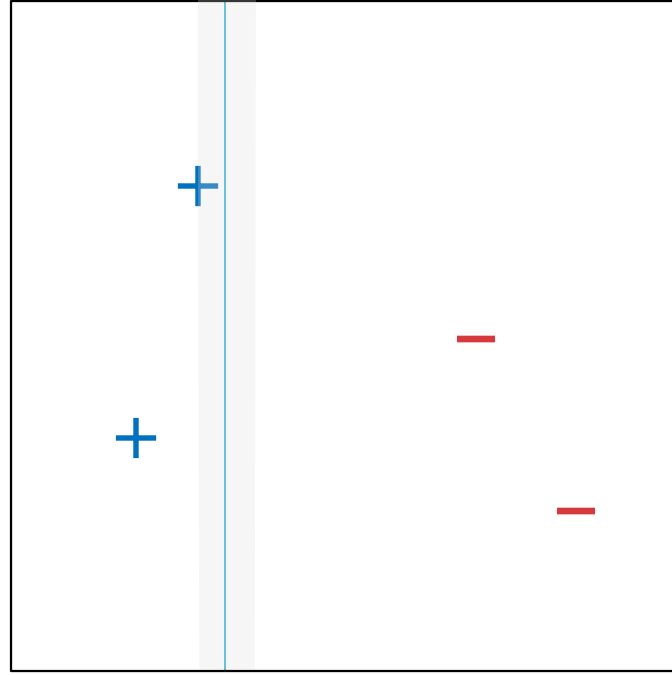
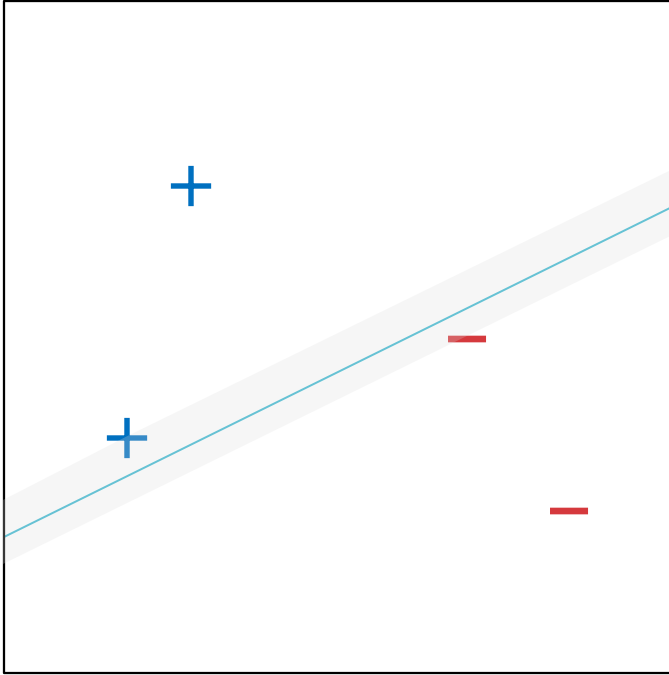
# Key Takeaways

- Boosting targets simple models, i.e., weak learners
- Greedily minimizes the exponential loss, an upper bound of the classification error
- Theoretical (and empirical) results show resilience to overfitting by targeting training margin



**In-class Poll:**  
Which linear separator is best?





# Which linear separator is best?



# Maximal Margin Linear Separators

- The margin of a linear separator is the distance between it and the nearest training data point
- Questions:
  1. How can we efficiently find a maximal-margin linear separator?
  2. Why are linear separators with larger margins better?
  3. What can we do if the data is not linearly separable?

## Recall: Hyperplanes

- For linear models, decision boundaries are  $D$ -dimensional *hyperplanes* defined by a weight vector,  $[b, \mathbf{w}]$

$$\mathbf{w}^T \mathbf{x} + b = 0$$

- Problem: there are infinitely many weight vectors that describe the same hyperplane
  - $x_1 + 2x_2 + 2 = 0$  is the same line as  $2x_1 + 4x_2 + 4 = 0$ , which is the same line as  $1000000x_1 + 2000000x_2 + 2000000 = 0$
- Solution: normalize weight vectors *w.r.t. the training data*

# Normalizing Hyperplanes

- Given a dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  where  $y \in \{-1, +1\}$ ,  $\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$  is a valid **linear separator** if

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) > 0 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}$$

- For SVMs, we're *only* going to consider **linear separators** in

$$\mathcal{H} = \left\{ \hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b) : \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1 \right\}$$

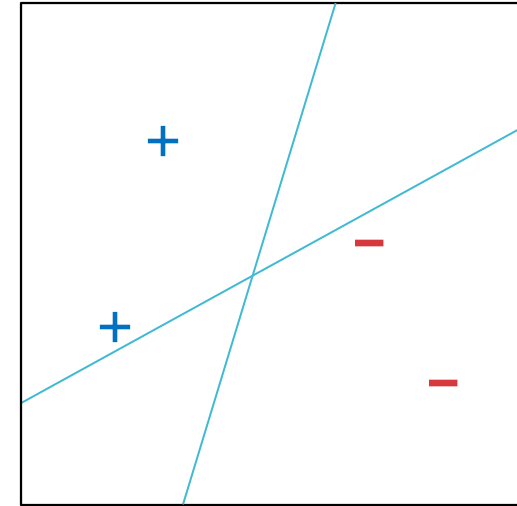
- If  $\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$  is a linear separator, then

$$\hat{y} = \text{sign}\left(\frac{\mathbf{w}^T}{\rho} \mathbf{x} + \frac{b}{\rho}\right) \in \mathcal{H} \text{ where}$$

$$\rho = \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

# Normalizing Hyperplanes: Example

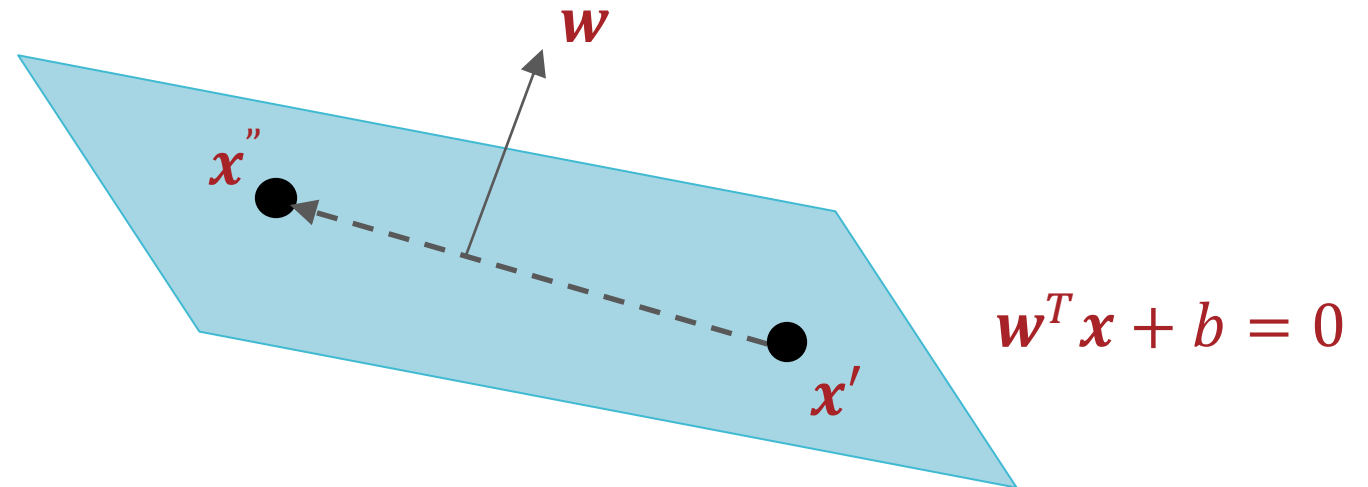
$b$	$w_1$	$w_2$	
-0.2	-0.6	1	$\notin \mathcal{H}$
-0.4	-1.2	2	$\notin \mathcal{H}$
-2	-6	10	$\notin \mathcal{H}$
-10	-30	50	$\in \mathcal{H}$
0.2	-0.6	0.2	$\notin \mathcal{H}$
0.1	-0.3	0.1	$\notin \mathcal{H}$
1	-3	1	$\notin \mathcal{H}$
2	-6	2	$\in \mathcal{H}$



$x_1$	$x_2$	$y$	$y(w^T x + b)$
0.2	0.4	+1	1.6
0.3	0.8	+1	1.8
0.7	0.6	-1	1
0.8	0.3	-1	2.2

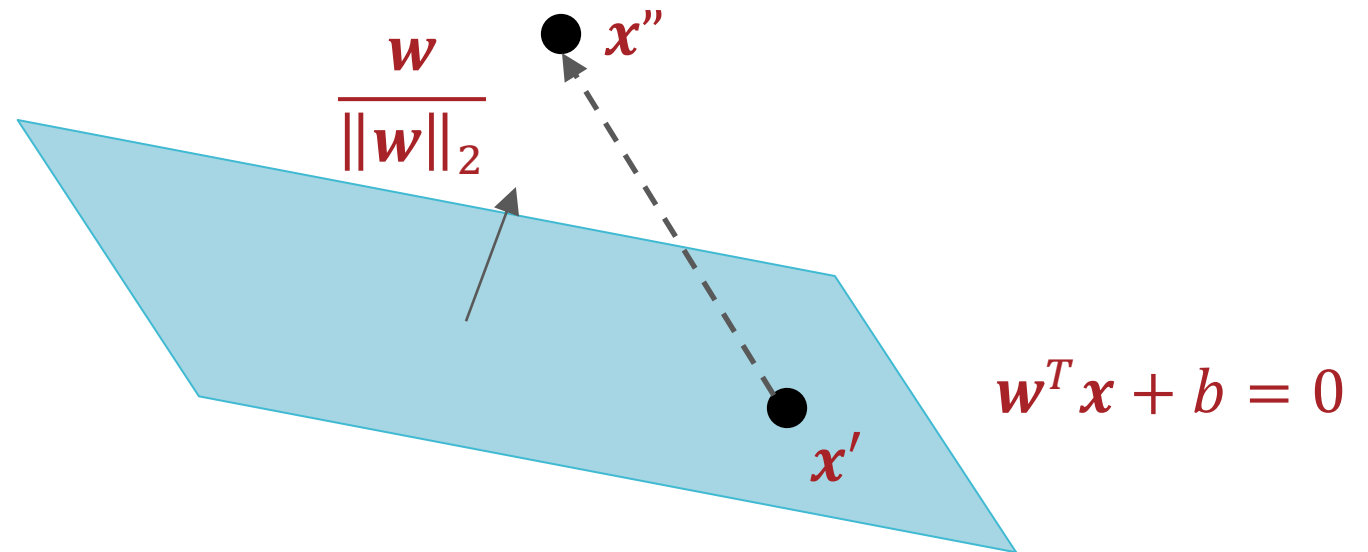
# Computing the Margin

- Claim:  $\mathbf{w}$  is orthogonal to the hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  (the decision boundary)
- A vector is orthogonal to a hyperplane if it is orthogonal to every vector in that hyperplane
- Vectors  $\alpha$  and  $\beta$  are orthogonal if  $\alpha^T \beta = 0$



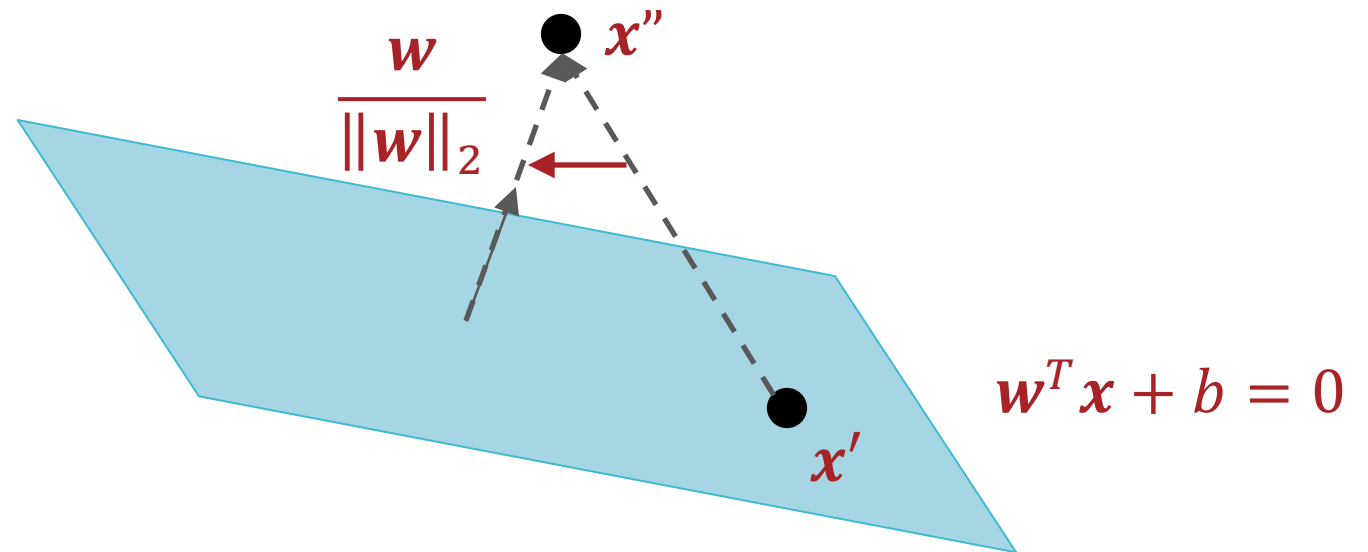
# Computing the Margin

- Let  $\mathbf{x}'$  be an arbitrary point on the hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  and let  $\mathbf{x}''$  be an arbitrary point
- The distance between  $\mathbf{x}''$  and  $\mathbf{w}^T \mathbf{x} + b = 0$  is equal to the magnitude of the projection of  $\mathbf{x}'' - \mathbf{x}'$  onto  $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , the unit vector orthogonal to the hyperplane



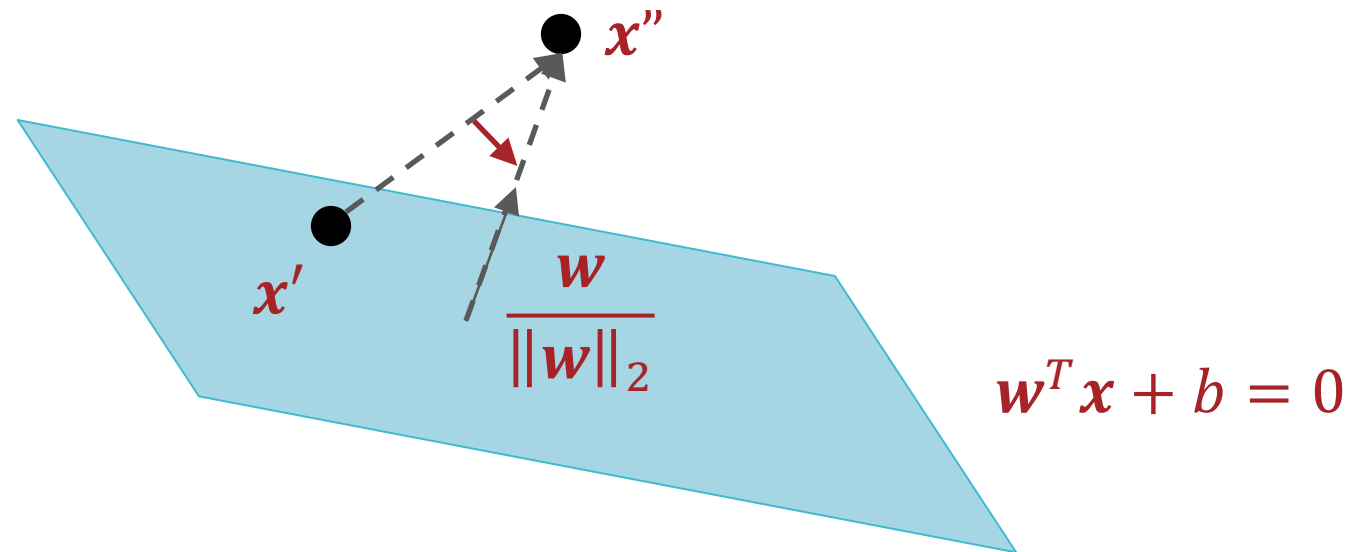
# Computing the Margin

- Let  $\mathbf{x}'$  be an arbitrary point on the hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  and let  $\mathbf{x}''$  be an arbitrary point
- The distance between  $\mathbf{x}''$  and  $\mathbf{w}^T \mathbf{x} + b = 0$  is equal to the magnitude of the projection of  $\mathbf{x}'' - \mathbf{x}'$  onto  $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , the unit vector orthogonal to the hyperplane



# Computing the Margin

- Let  $\mathbf{x}'$  be an arbitrary point on the hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  and let  $\mathbf{x}''$  be an arbitrary point
- The distance between  $\mathbf{x}''$  and  $\mathbf{w}^T \mathbf{x} + b = 0$  is equal to the magnitude of the projection of  $\mathbf{x}'' - \mathbf{x}'$  onto  $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , the unit vector orthogonal to the hyperplane





## Computing the Margin

- Let  $\mathbf{x}'$  be an arbitrary point on the hyperplane  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$  and let  $\mathbf{x}''$  be an arbitrary point
- The distance between  $\mathbf{x}''$  and  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$  is equal to the magnitude of the projection of  $\mathbf{x}'' - \mathbf{x}'$  onto  $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ , the unit vector orthogonal to the hyperplane

$$\begin{aligned} d(\mathbf{x}'', h) &= \left| \frac{\mathbf{w}^T (\mathbf{x}'' - \mathbf{x}')}{\|\mathbf{w}\|_2} \right| = \frac{|\mathbf{w}^T \mathbf{x}'' - \mathbf{w}^T \mathbf{x}'|}{\|\mathbf{w}\|_2} \\ &= \frac{|\mathbf{w}^T \mathbf{x}'' + b|}{\|\mathbf{w}\|_2} \end{aligned}$$

## Computing the Margin

- The margin of a linear separator is the distance between it and the nearest training data point

$$\begin{aligned} \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} d(\mathbf{x}^{(i)}, h) &= \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \frac{|\mathbf{w}^T \mathbf{x}^{(i)} + b|}{\|\mathbf{w}\|_2} \\ &= \frac{1}{\|\mathbf{w}\|_2} \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} |\mathbf{w}^T \mathbf{x}^{(i)} + b| \\ &= \frac{1}{\|\mathbf{w}\|_2} \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \\ &= \frac{1}{\|\mathbf{w}\|_2} \end{aligned}$$

# Maximizing the Margin

$$\begin{aligned} & \text{maximize} \quad \frac{1}{\|\mathbf{w}\|_2} \\ & \text{subject to} \quad \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1 \end{aligned}$$

$\Updownarrow$

$$\begin{aligned} & \text{minimize} \quad \|\mathbf{w}\|_2 \\ & \text{subject to} \quad \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1 \end{aligned}$$

$\Updownarrow$

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{subject to} \quad \min_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1 \end{aligned}$$

$\Updownarrow$

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{subject to} \quad y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \quad \forall (\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D} \end{aligned}$$