

# CS 412



---

MARCH 31<sup>ST</sup> – RETURN TO MACHINE LEARNING

# Welcome back from "Spring Break"

---

## HW3

- Graded by end-of-week

## HW4

- Out next Thursday after class
- Due two Thursdays after (4/16)

## Midterm

- Next Thursday
- Take home
- 12-8 pm CDT
- Should take around 2 hours

Mostly written, but there may be some small coding portion

HW2 regrades  
done by end of  
the week

(4/9)  
Submitted on gradescope

Next Tuesday  
will be exam  
review

Focus  
Midterm  
practice

# Going Forward

---

Slack

## Online classes

- Mostly blackboard collab during lectures.
- One panopto lecture over the weekend to catch up on material ← downloadable

## Graded things left

- HW4 Out next week,
- HW5
- Midterm
- Final Exam
- Final Project

Off:  
blackboard.

# Final Project

---

Official write up/description

Groups of up to 3

- Groups of 2, I expect ~50% more
- Groups of 3, ~100% more

*Solo projects  
are fine*

Every group will need to meet digitally with me (~15 minutes) between 4/13 and 4/17 to make sure the project is feasible

Project will be 5-10 pages, single-spaced, but with figures

- Problem introduction
- Previous work
- Novel attempt ↩
- Results
- Future work/Conclusion

# Final Project

---

Things to consider:

- Must involve a new dataset
- Must involve at least one element not covered in the course material
- Should also include a significant amount of the course material
- *For graduate students*, must include some ML research from the last 5 years

↳ a lot of researchers post their code

I **strongly** encourage you to find a project that aligns with your personal interests / research agenda

In general, I want to allow any reasonable machine learning project

# Tips for finding a project

---

Was there something difficult for you in a homework assignment?

Is there a particular dataset that you would find interesting?

- For example, there is a large amount of medical data

What sorts of data do you think would give a different result from what you've seen?

# Rationale for Ensemble Learning

---

No Free Lunch thm: There is no algorithm that is always the most accurate

Generate a group of base-learners which when combined have higher accuracy

Different learners use different ensemble or higher precision

- Algorithms
- Parameters
- Representations (Modalities)
- Training sets
- Subproblems

3 independent  
10%

/

ensemble  
method

# Bagging (Bootstrap aggregating)

---

Take M bootstrap samples (with replacement)

Train M different classifiers on these bootstrap samples

For a new query, let all classifiers predict and take an average (or majority vote)

If the classifiers make independent errors, then their ensemble can improve performance.

- This is reduced with “stable” learners

↙ less independence when the model is stable  
unstable learners get a big variance reduction

Stated differently: the variance in the prediction is reduced (we don't suffer from the random errors that a single classifier is bound to make).



# Boosting

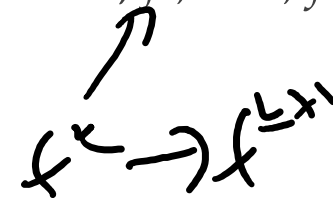
---

Boosting works differently.

1. Boosting does not involve bootstrap sampling

2. → Decision models are grown sequentially: each model is created using information from previously grown trees

3. Like bagging, boosting involves combining a large number of models,  $f^1, \dots, f^B$



# Boosting

---

Train classifiers (e.g. decision trees) in a sequence.

A new classifier should focus on those cases which were incorrectly classified in the last round.

Combine the classifiers by letting them vote on the final prediction (like bagging).

Each classifier is “weak” but the ensemble is “strong.”

AdaBoost is a specific boosting method.

10<sup>10</sup> →

$f_k$  “focus” on  
elements incorrectly  
classified by  $f' \rightarrow f^k$

# Example

Stability:  
to what extent  
does a change in the  
data set effect

a change in the model



incorrectly  
classified  
by  $f_1$

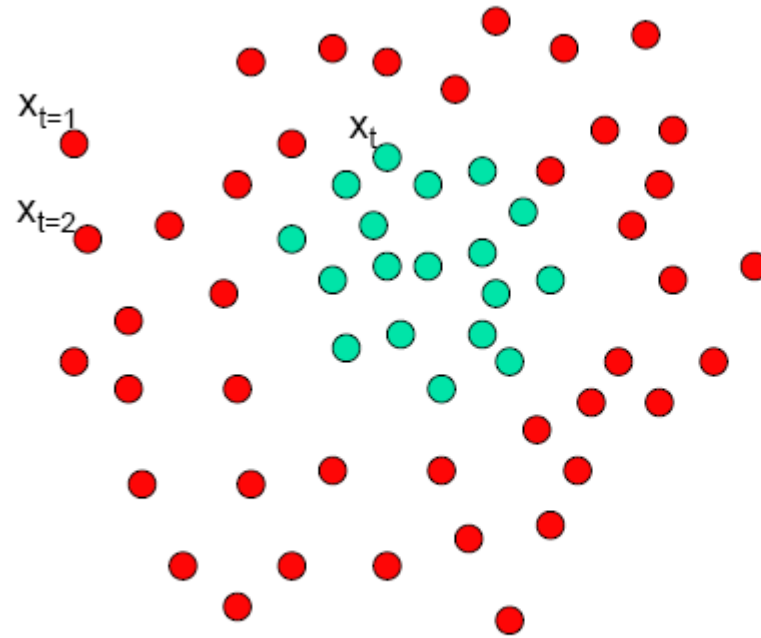
weight  
these points  
more  
when we  
build  $f_2$

This line is one simple classifier saying that everything to the left + and everything to the right is -

linear separators

# Boosting - Example

---



Each data point has  
a class label:

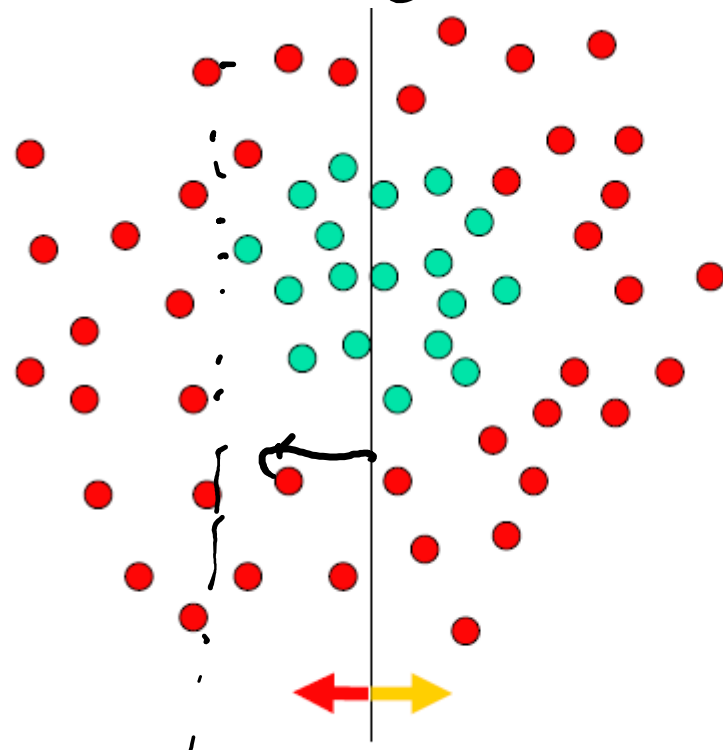
$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{cyan circle}) \end{cases}$$

and a weight:

$$w_t = 1$$

# Boosting - Example

Not a good learner



Each data point has  
a class label:

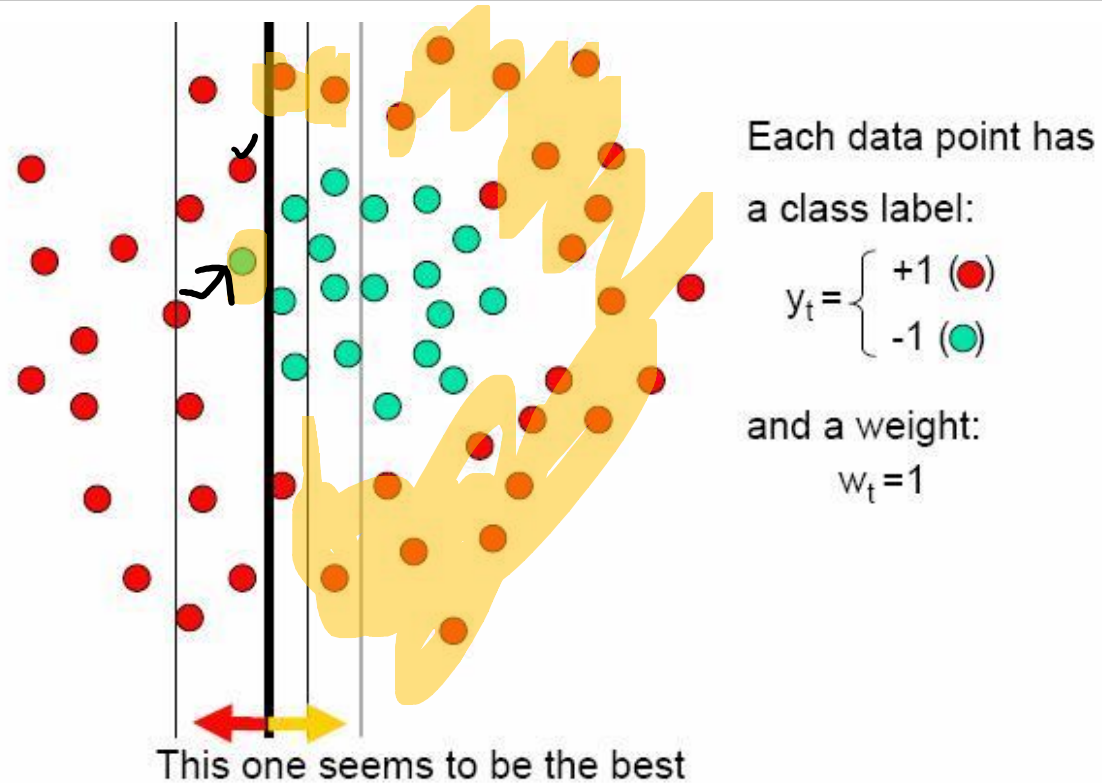
$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{cyan circle}) \end{cases}$$

and a weight:

$$w_t = 1$$

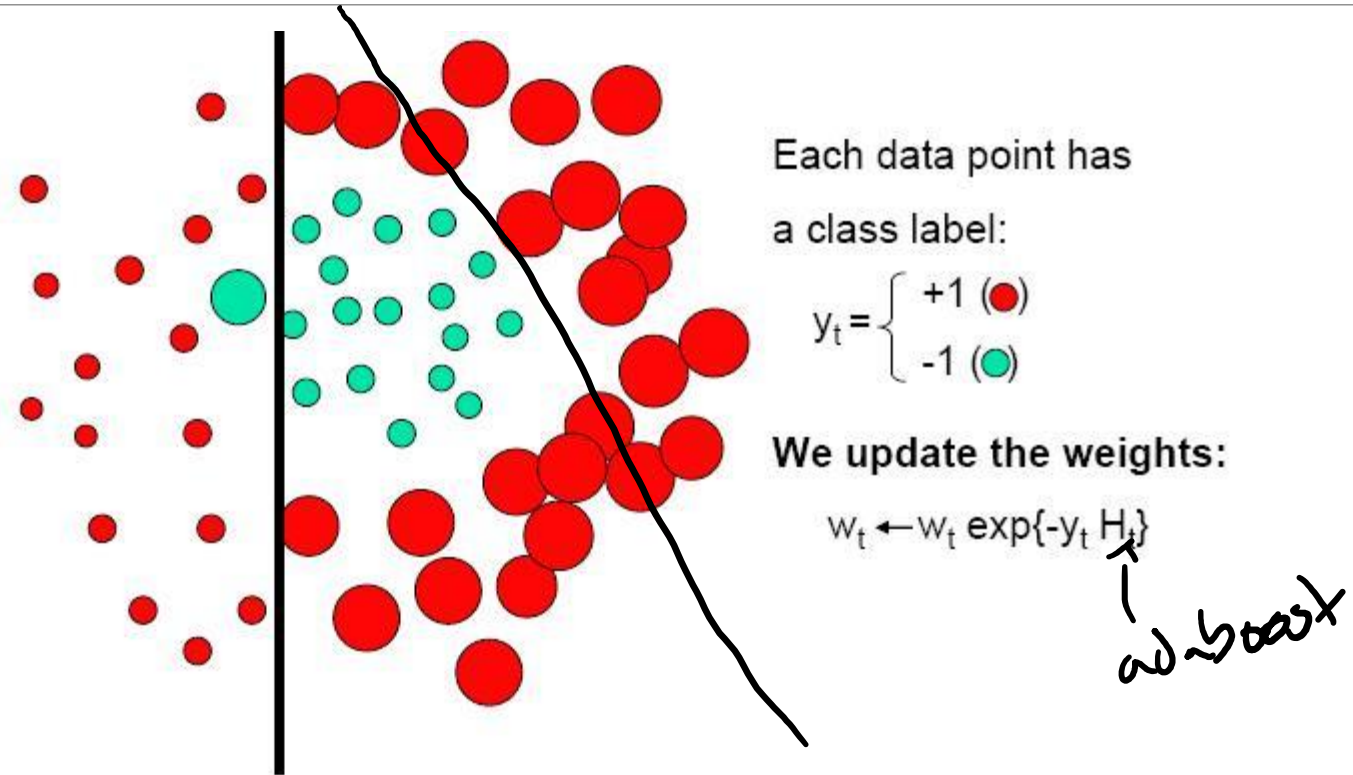
$h \Rightarrow p(\text{error}) = 0.5$  it is at chance

# Boosting - Example



This is a **'weak classifier'**: It performs slightly better than chance.

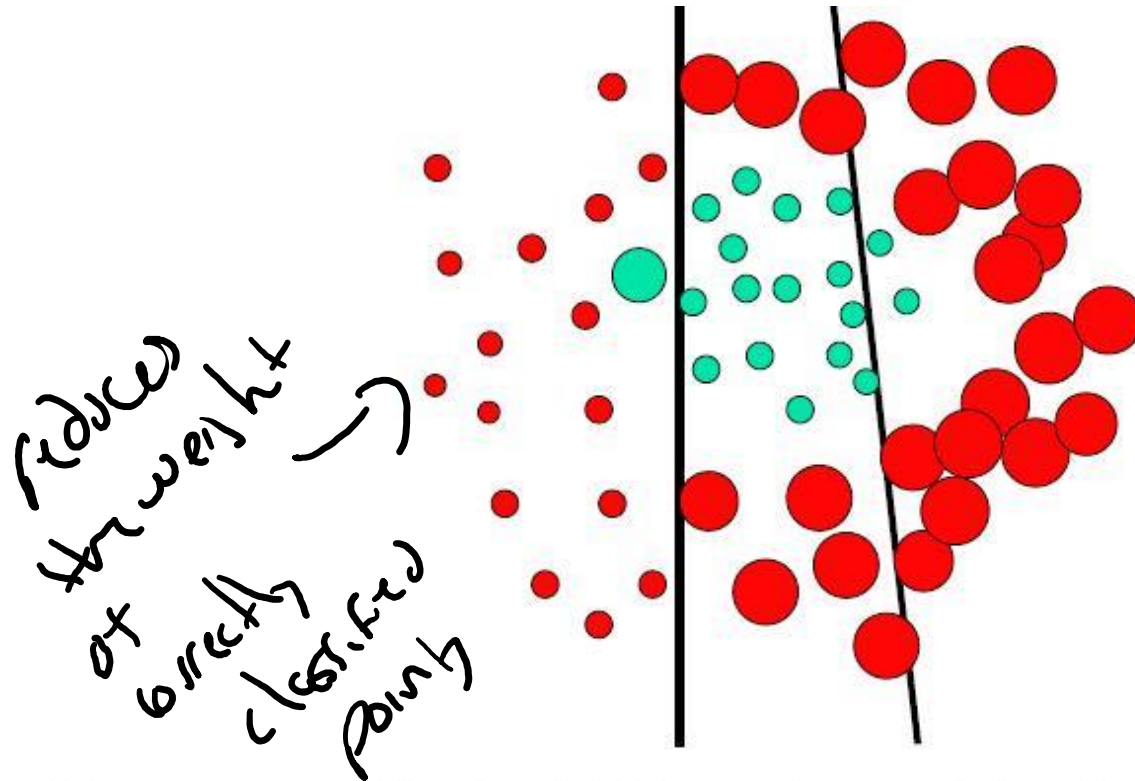
# Boosting - Example



We set a new problem for which the previous weak classifier performs at chance again

# Boosting - Example

Every point has weight 50% at 1



Each data point has a class label:

$$y_t = \begin{cases} +1 & (\text{red}) \\ -1 & (\text{teal}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

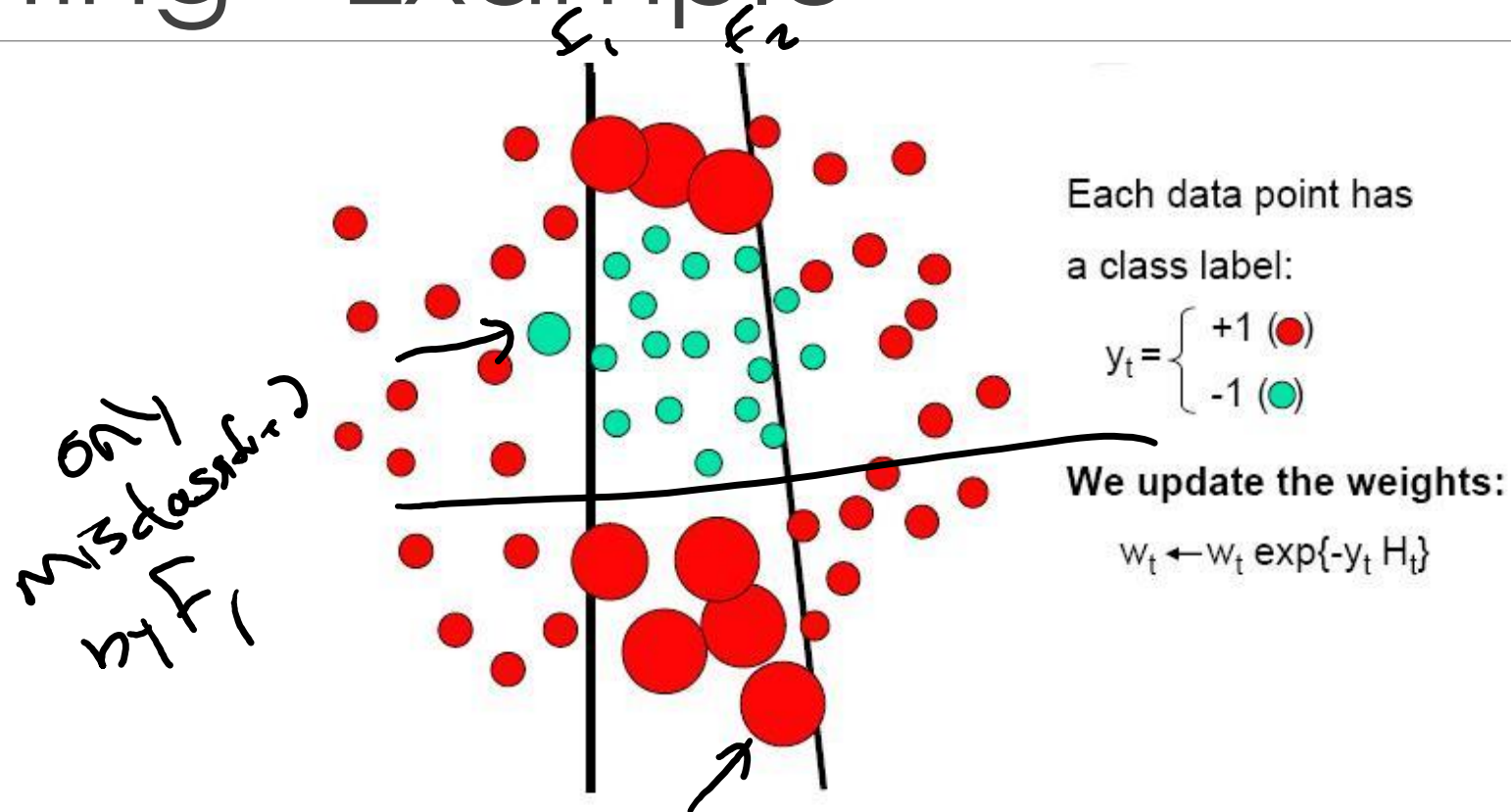
teal points

Reduce weight of correctly classified points

We set a new problem for which the previous weak classifier performs at chance again



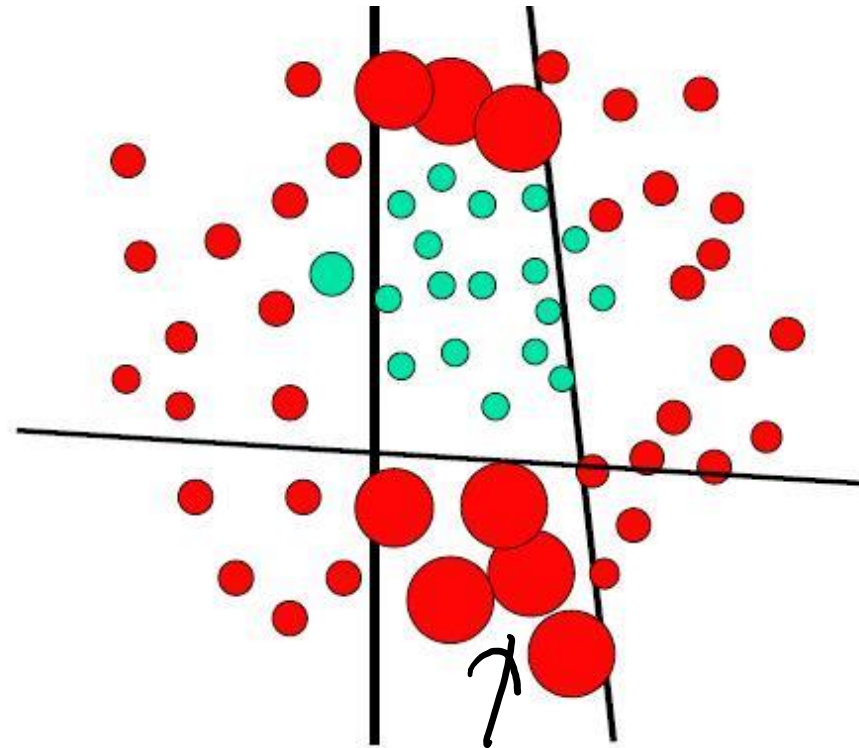
# Boosting - Example



We set a new problem for which the previous weak classifier performs at chance again

# Boosting - Example

X | 0  
X | 0  
1 | 0  
→



Each data point has  
a class label:

$$y_t = \begin{cases} +1 & (\text{red}) \\ -1 & (\text{teal}) \end{cases}$$

We update the weights:

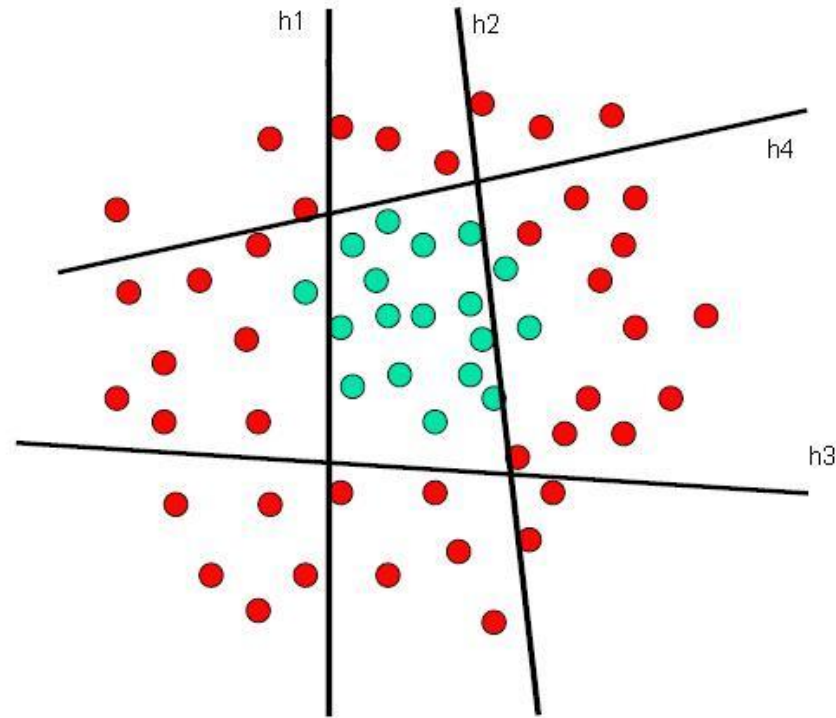
$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

↑  
cross-entropy  
of the classifier

We set a new problem for which the previous weak classifier performs at chance again

# Boosting - Example

---



The strong (non- linear) classifier is built as the combination of all the weak (linear) classifiers.

# Boosting Intuition

We adaptively weigh each data case.

Weight = 2  
= 2 points in  
the feature space

Data cases which are wrongly classified get high weight (the algorithm will focus on them)

Each boosting round learns a new (simple) classifier on the weighed dataset.

These classifiers are weighed to combine them into a single powerful classifier.

$$F_1 \rightarrow \alpha_1 F_1$$

↳ not strictly voting

Classifiers that obtain low training error rate have high weight.

We stop by using monitoring a hold out set (cross-validation).

# Boosting Intuition

---

Why do we want to have

1. Weak learners?
2. Learners trained on different data?

↳ Bagging: bootstrap

↳ Boosting: weighted points

# Boosting Intuition

---

Why do we want to have

1. Weak learners?
2. Learners trained on different data?

We want the trainers to be as independent as possible

even with boosted data

learners can be similar w/ strong learners

boosted data causes a bigger difference

in models if the models are weak/unstable

# Boosting

or weak model

Train a set of weak hypotheses:  $h_1, \dots, h_T$ .

# hypotheses = T

The combined hypothesis  $H$  is a weighted majority vote of the  $T$  weak hypotheses.

- Each hypothesis  $h_t$  has a weight  $\alpha_t$ .

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

$\alpha_t$  = each weight for each model  $t$

During the training, focus on the examples that are misclassified.

- At round  $t$ , example  $x_i$  has the weight  $D_t(i)$ .

$x_i$   $D_t(i)$

# Boosting

$$D_1(i) = \frac{1}{m}$$

Binary classification problem

Training data:

$(x_1, y_1), \dots, (x_m, y_m)$ , where  $x_i \in X, y_i \in Y = \{-1, 1\}$

$m = \text{number of data points}$

$D_t(i)$ : the weight of  $x_i$  at round  $t$ .  $D_1(i) = 1/m$

A learner  $L$  that finds a weak hypothesis  $h_t: X \rightarrow Y$  given the training set and  $D_t$

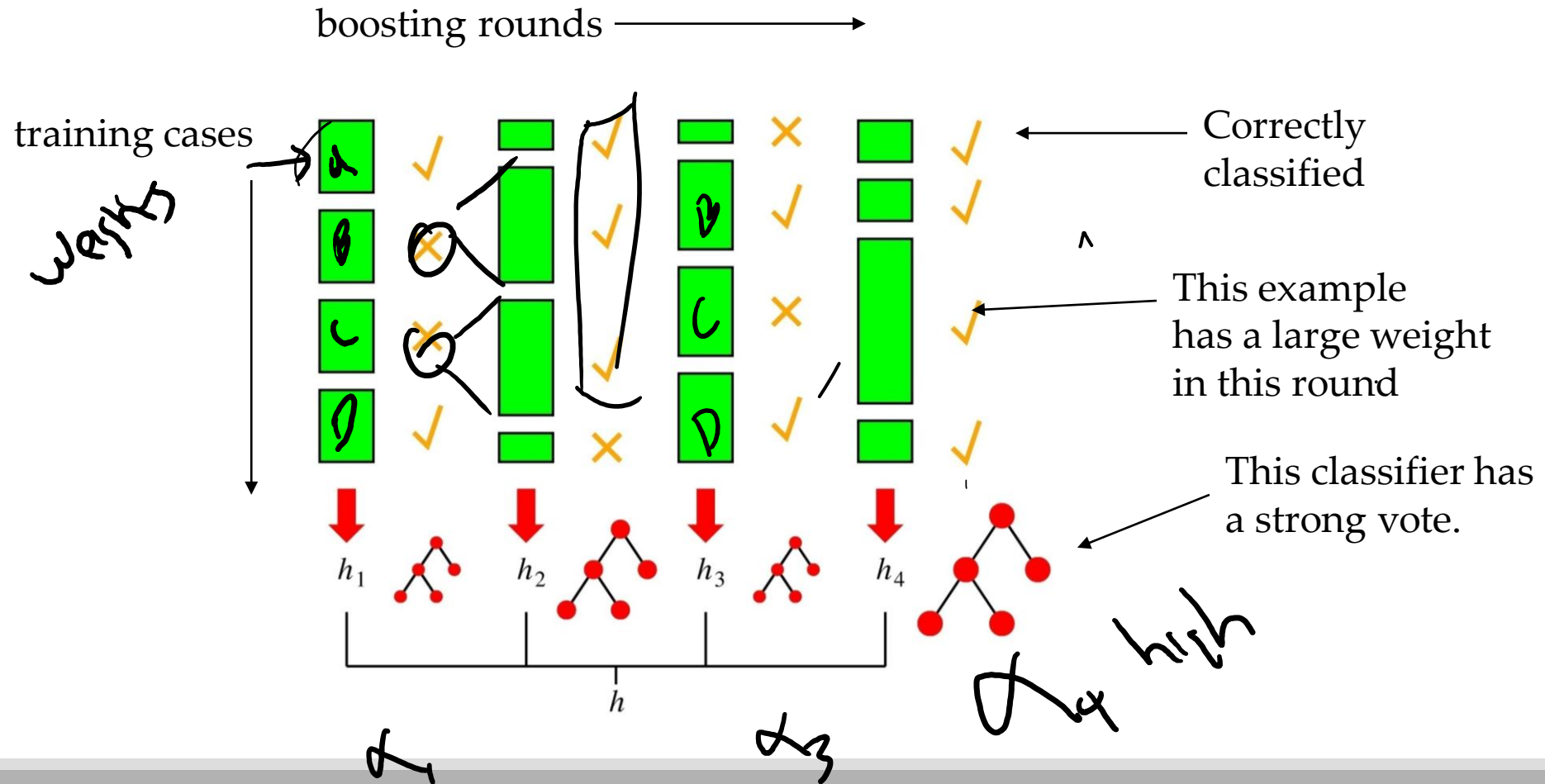
The error of a weak hypothesis  $h_t$ :

$$\varepsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i] = \sum_{i: h_t(x_i) \neq y_i} D_t(i)$$

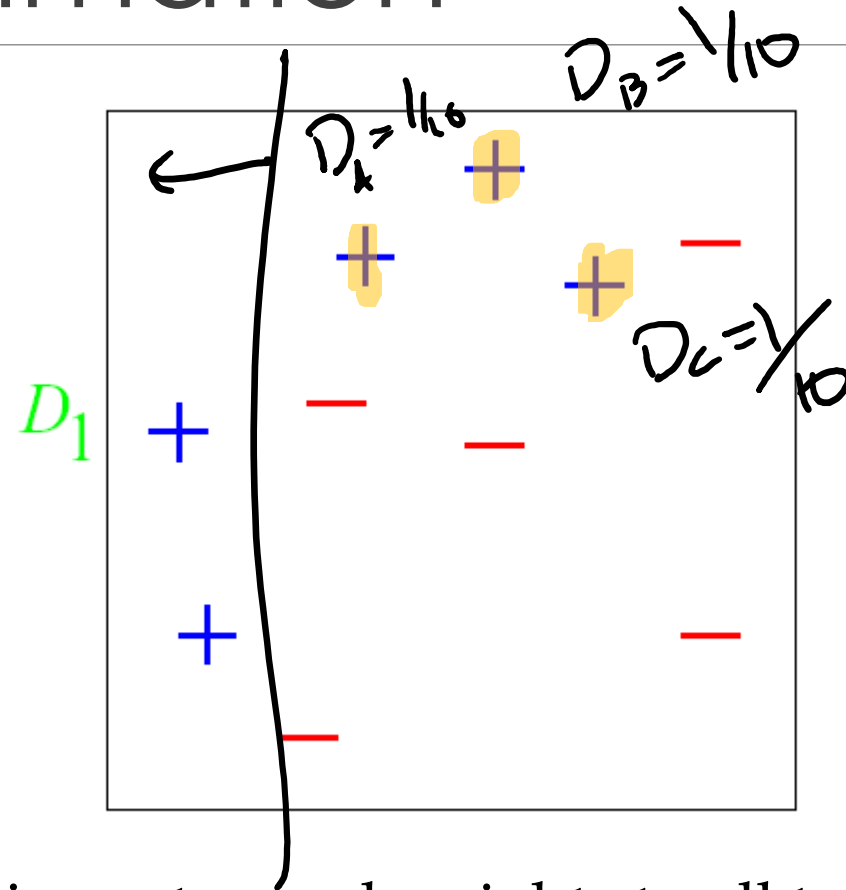
$D = \text{weight}$   
Example  $m = 100$   
 $D_1 = 1/100$   
 $\sum_{i=1}^{100} 1/100 = 100/100$



# Boosting in a Picture



# And in animation



$$\xi_e = \sum_{\text{nodes}} D$$

$$\alpha_e = \frac{\xi_e}{(1 - \xi_e)}$$

$$\xi_1 = 0.30$$

$$\alpha = \frac{0.30}{0.7} = 0.42$$

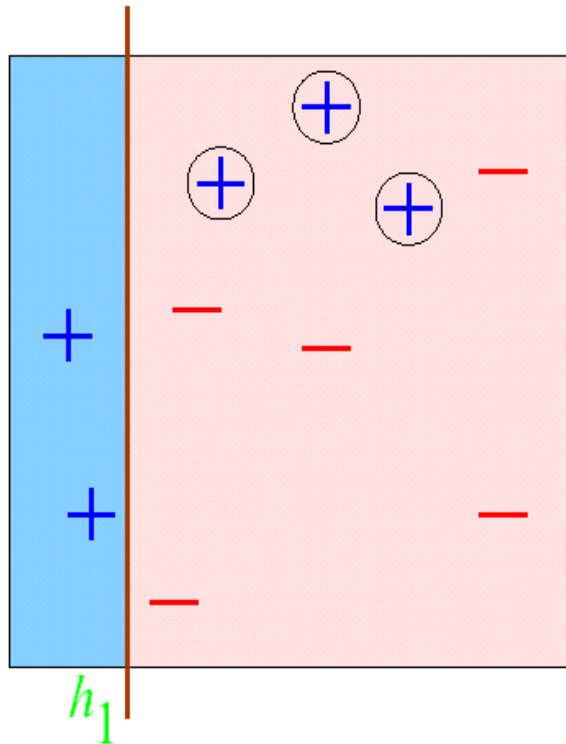
Original training set: equal weights to all training samples

# Boost example

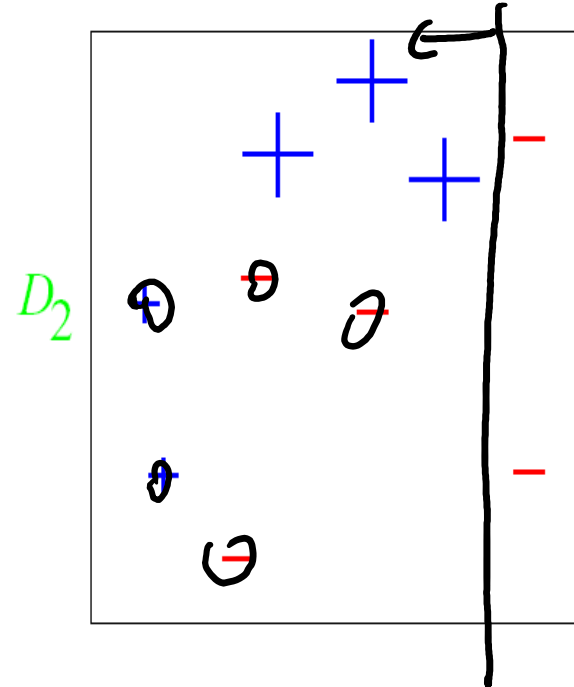
ROUND 1

$\varepsilon$  = error rate of classifier

$\alpha$  = weight of classifier e.g.  $[\varepsilon/(1-\varepsilon)]$

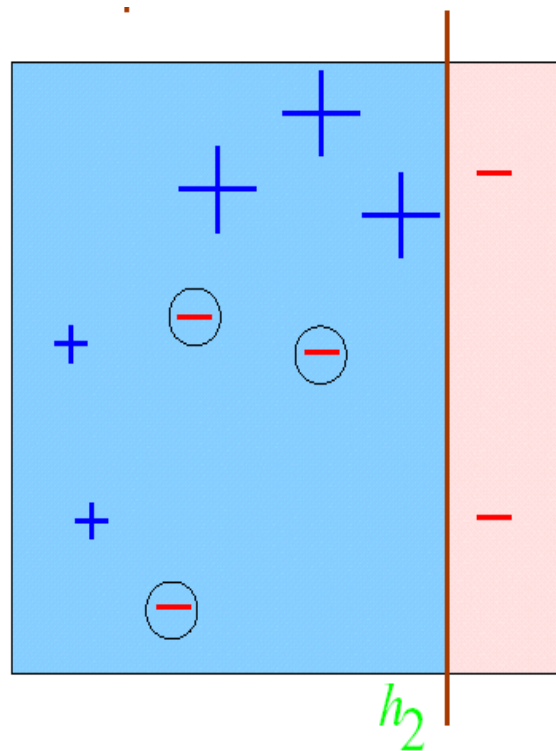


$$\begin{pmatrix} \varepsilon_1 = 0.30 \\ \alpha_1 = 0.42 \end{pmatrix}$$

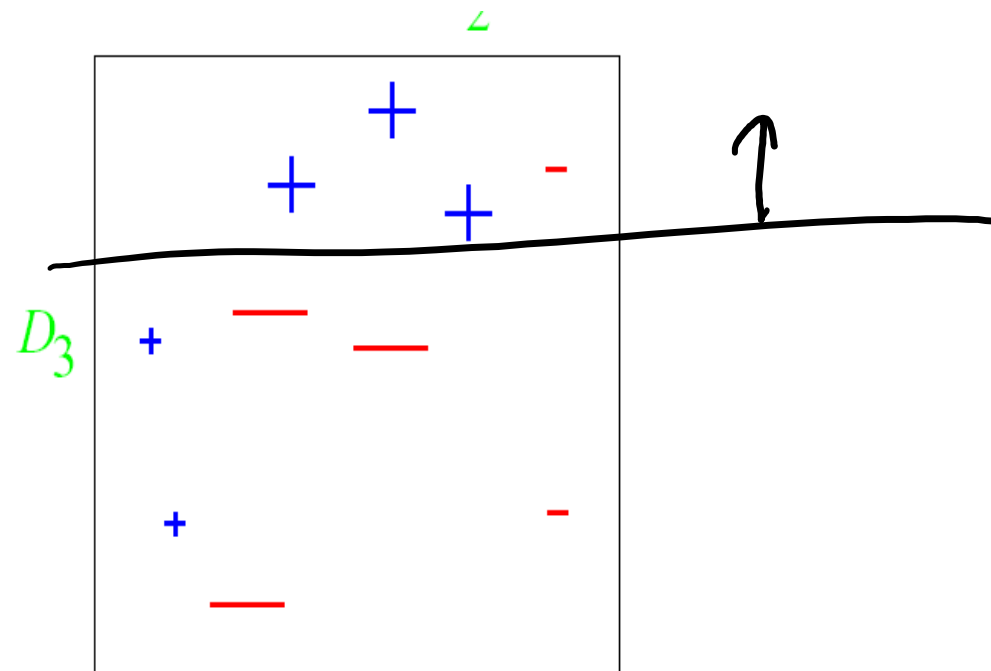


# Boost example

ROUND 2



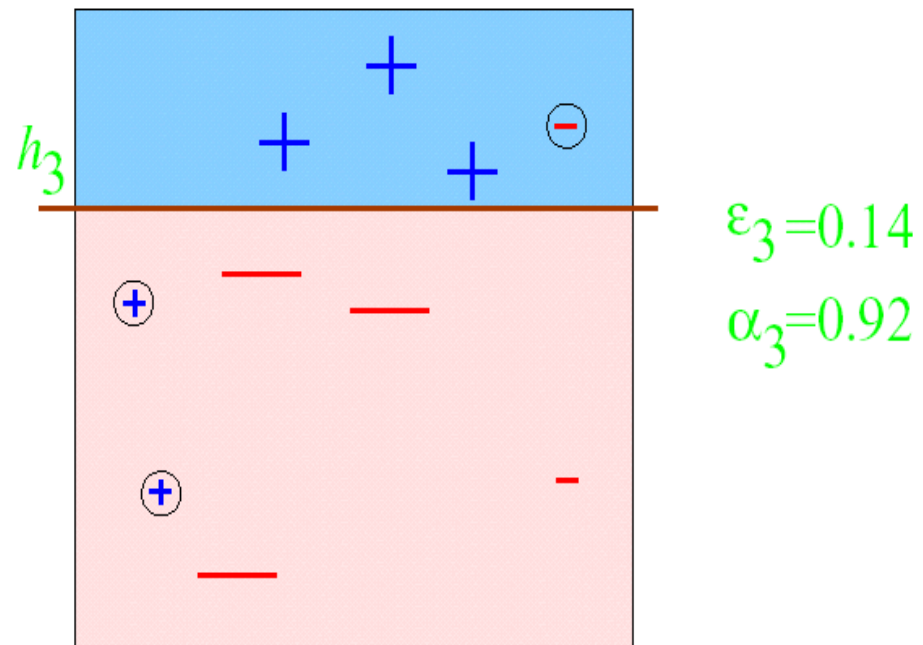
$\epsilon_2 = 0.21$   
 $\alpha_2 = 0.65$



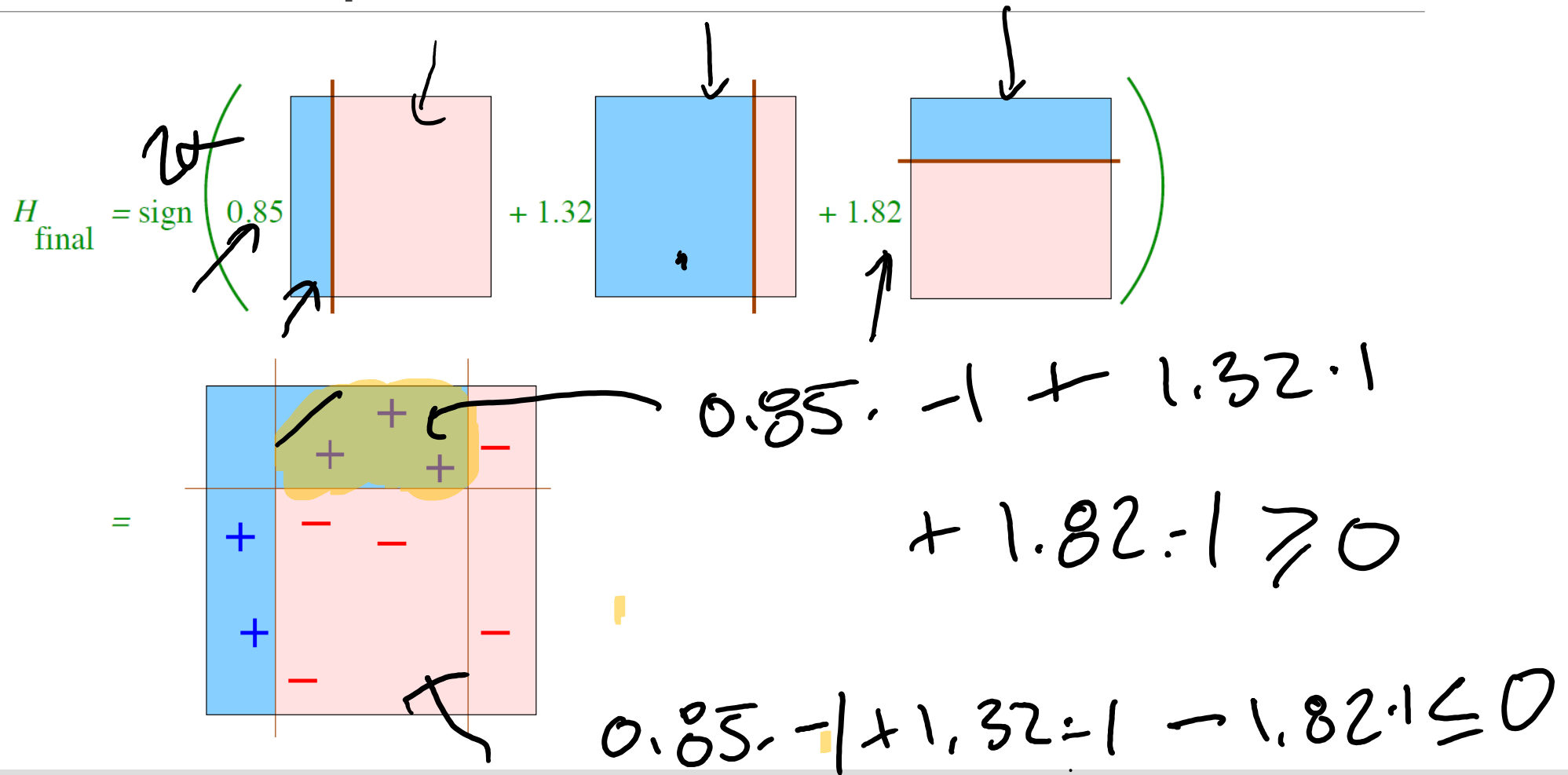
# Boost example

---

ROUND 3



# Boost example



# AdaBoost - Introduction

Linear classifier with all its desirable properties

← Decision stump -  
horizontal or vertical  
separator

Has good generalization properties

Is a feature selector with a principled strategy (minimisation of upper bound on empirical error)

Close to sequential decision making

↳ stacking  
cascading

# AdaBoost

$h_t(x)$   
returning  $(-1, 1)$

Is an algorithm for constructing a “strong” classifier as linear combination

$$f(x) = \sum_{t=1}^T \alpha_t h_t(x) \checkmark$$

of simple “weak” classifiers  $h_t(x)$ .

$h_t(x)$  - “weak” or basis classifier, hypothesis, “feature”

$H(x) = \text{sign}(f(x))$  - “strong” or final classifier/hypothesis  $\checkmark$  ensemble

$\nearrow$



# AdaBoost

---

Input – a training set:  $S = \{(x_1, y_1); \dots; (x_m, y_m)\}$

- $x_i \in X$ ,  $X$  instance space
- $y_i \in Y$ ,  $Y$  finite label space
  - in binary case  $Y = \{-1, +1\}$

Each round,  $t=1, \dots, T$ , AdaBoost calls a given weak or base learning algorithm – accepts as input a sequence of training examples ( $S$ ) and a set of weights over the training example ( $D_t(i)$ )

$D_t(i)$  = weight for point  $i$  in weak learner  $t$

# AdaBoost

---

The weak learner computes a weak classifier ( $h_t$ ), :  $h_t : X \rightarrow \mathbb{R}$

Once the weak classifier has been received, AdaBoost chooses a parameter ( $\alpha_t \in \mathbb{R}$ ) – intuitively measures the importance that it assigns to  $h_t$ .



# The main idea of AdaBoost

---

to use the weak learner to form a highly accurate prediction rule by calling the weak learner repeatedly on different distributions over the training examples.

initially, all weights are set equally, but each round the weights of incorrectly classified examples are increased so that those observations that the previously classifier poorly predicts receive greater weight on the next iteration.

# Boosting

Perform 3 iterations of the boosting algorithm on the following dataset. Double the error at each point and weight each model with  $\alpha_i = (1-\epsilon)/\epsilon$

# of learners

