

BOOSTING



Outline

- Boosting – General idea
- AdaBoost
 - Visual demo
 - Detailed example
 - From a loss perspective

Boosting – General idea

- Create a **strong** classifier by combining multiple **weak** classifiers
- A **weak classifier** guarantees a slightly better-than-random error rate



Initialize a uniform distribution $D^{(1)}$ over trainset S

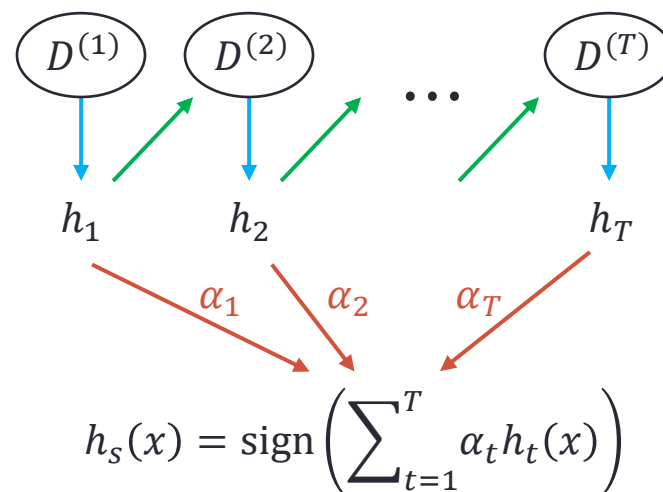
For $t=1, \dots, T$:

Learn (weak) model $h_t = \mathcal{A}(S, D^{(t)})$

Compute error

Update the distribution $D^{(t+1)}$

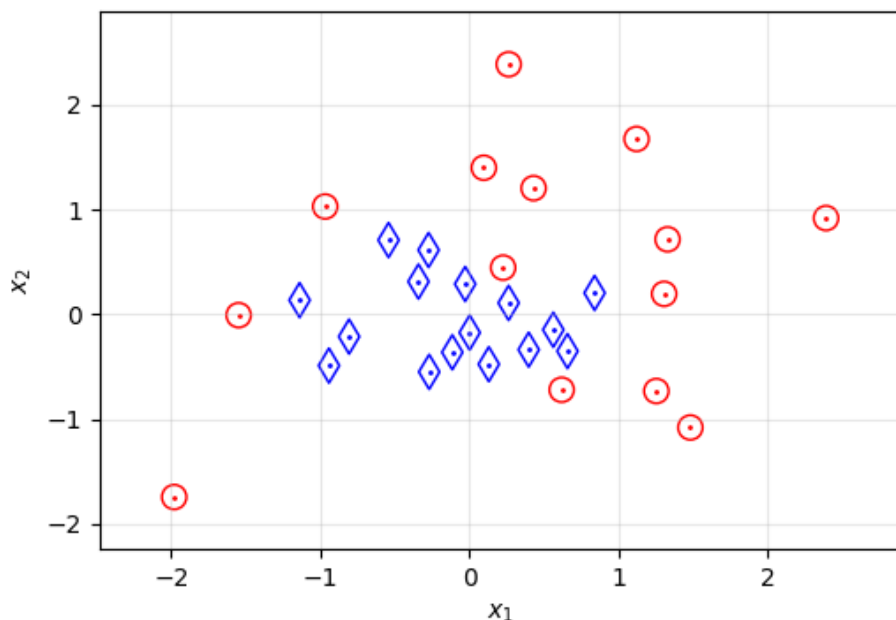
Return final hypothesis



AdaBoost – Visualization

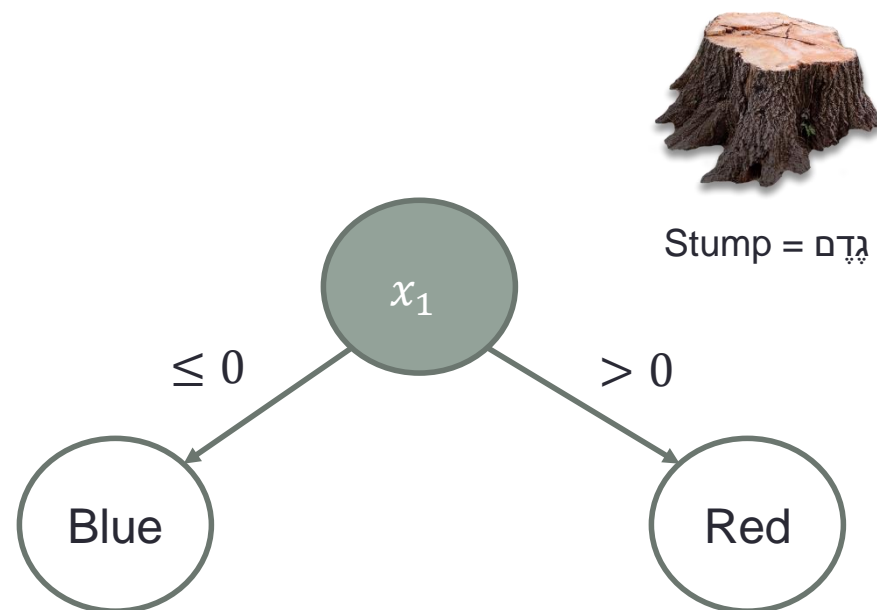
Dataset

- 2 features x_1, x_2
- 2 classes (red/blue)



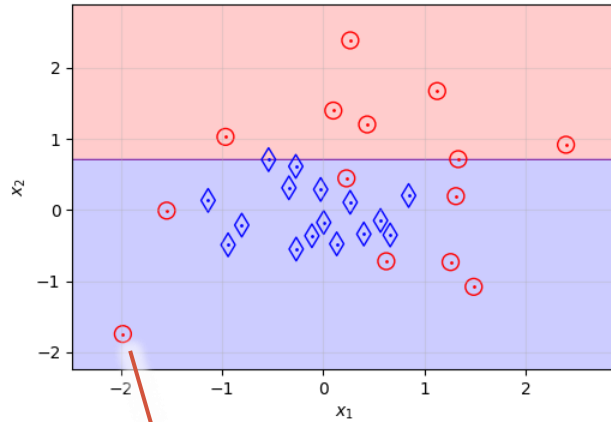
Weak learner choice

Decision stumps,
i.e., decision trees with depth 1

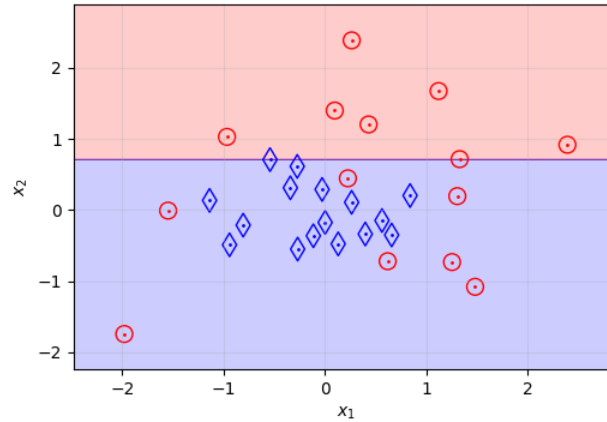


AdaBoost – Visualization

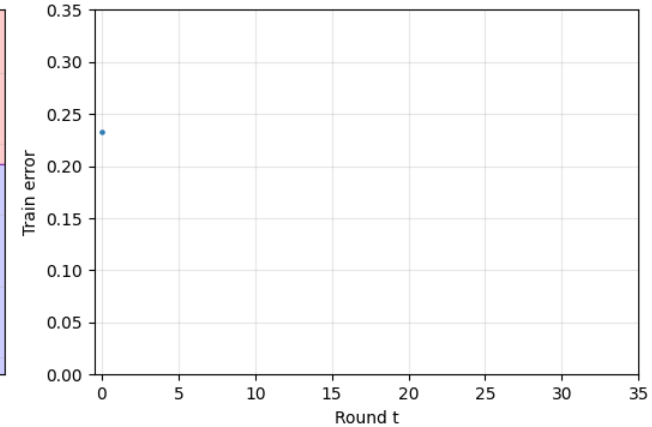
Weak learner at $t=1$



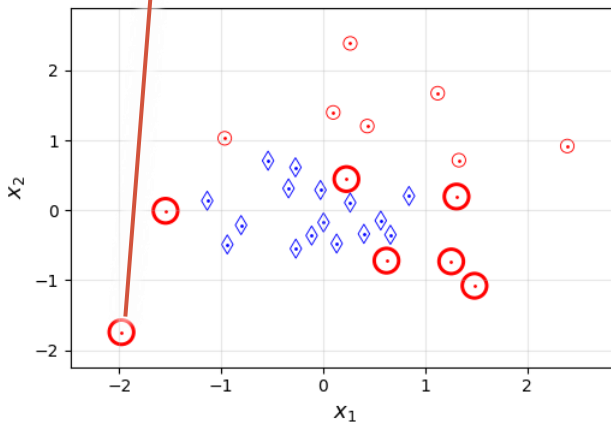
Strong learner at $t=1$



Train error

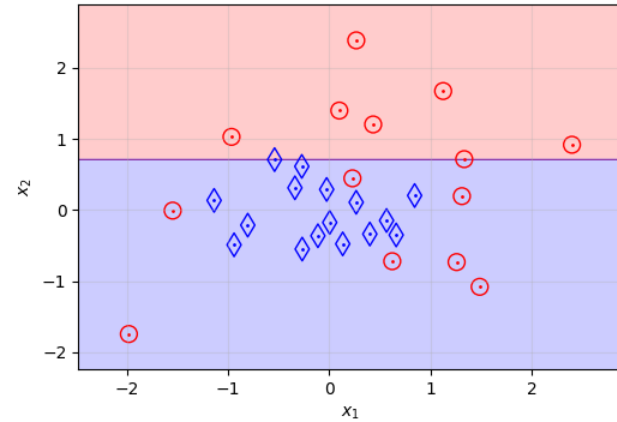


Misclassified points
get higher probability

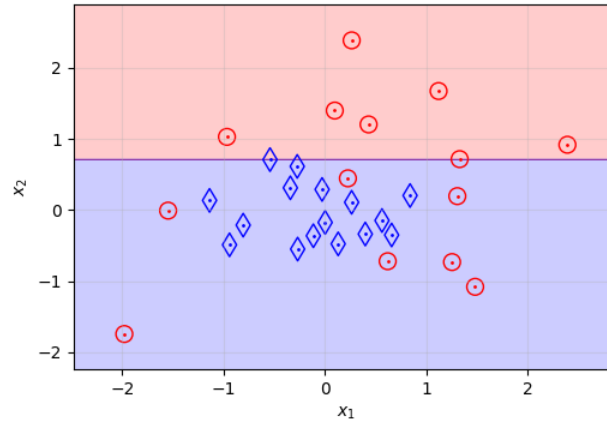


AdaBoost – Visualization

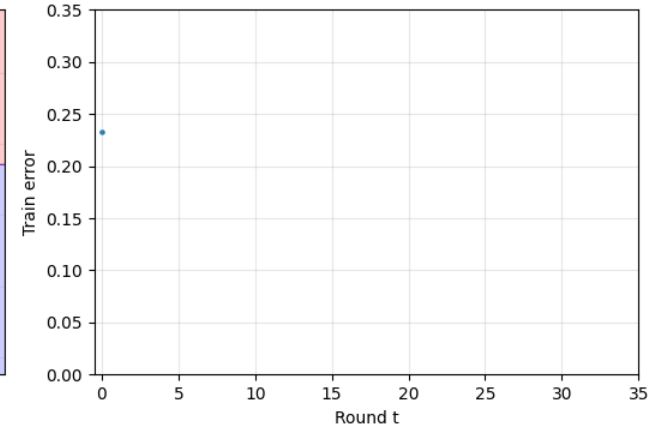
Weak learner at $t=1$



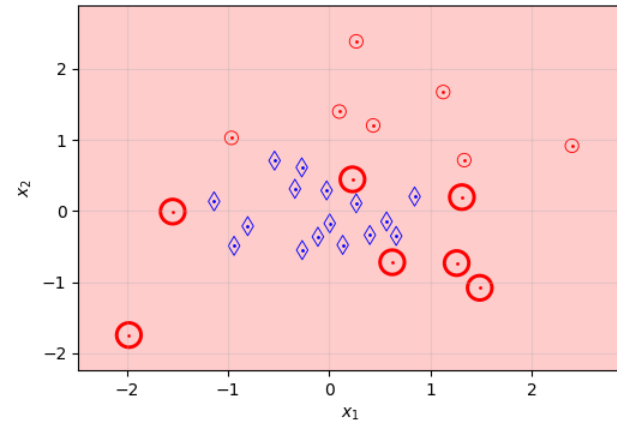
Strong learner at $t=1$



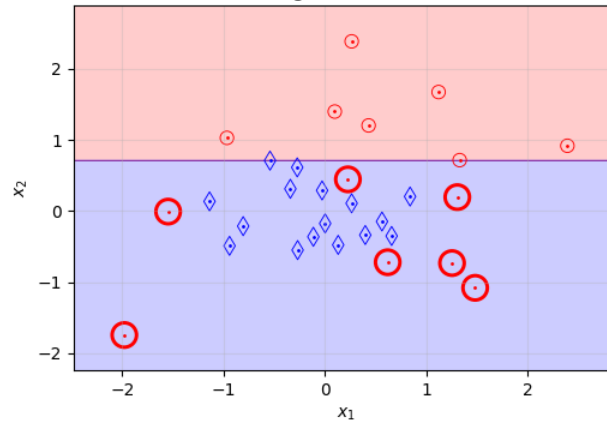
Train error



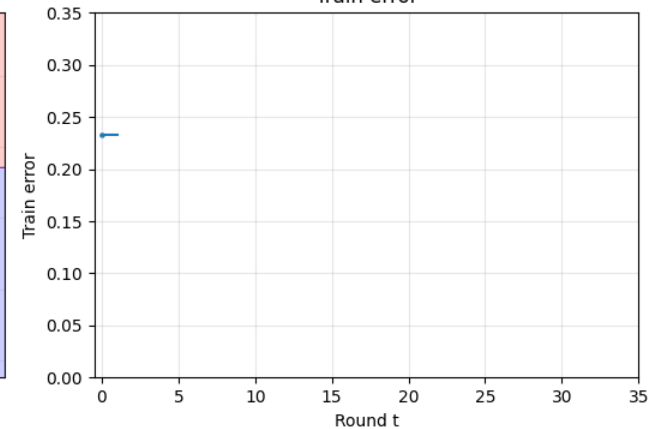
Weak learner at $t=2$



Strong learner at $t=2$

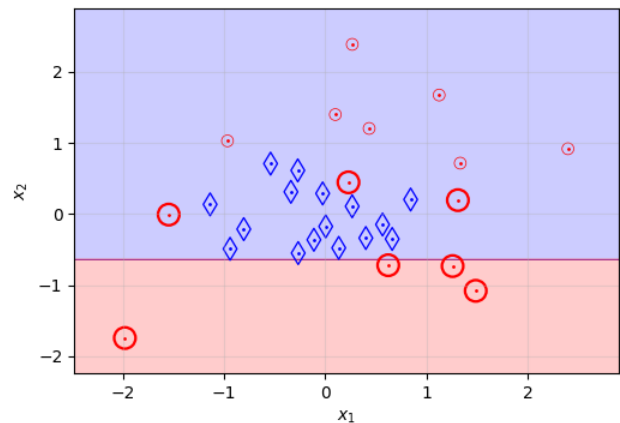


Train error

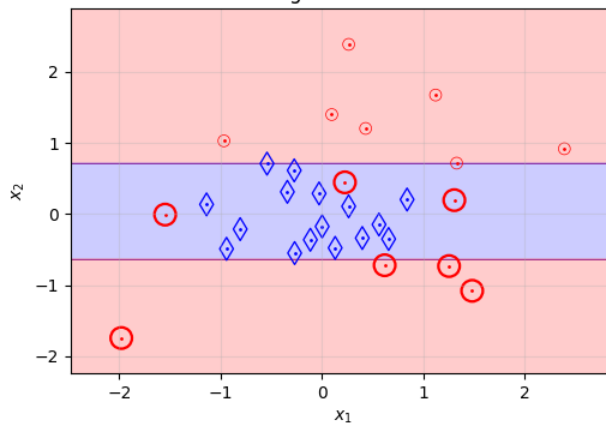


AdaBoost – Visualization

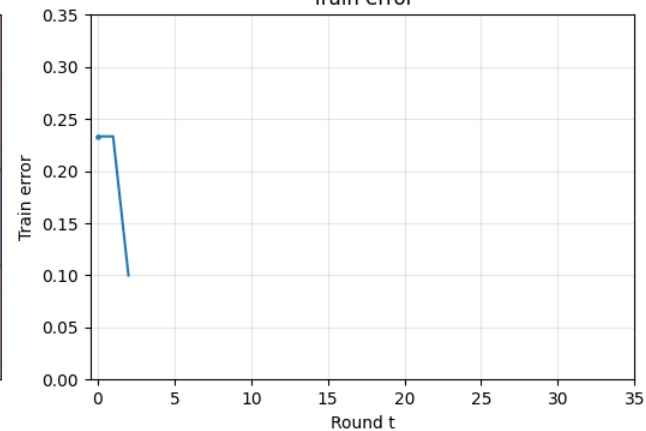
Weak learner at $t=3$



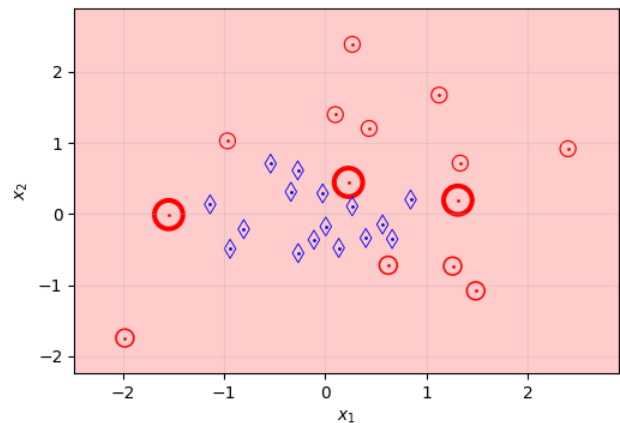
Strong learner at $t=3$



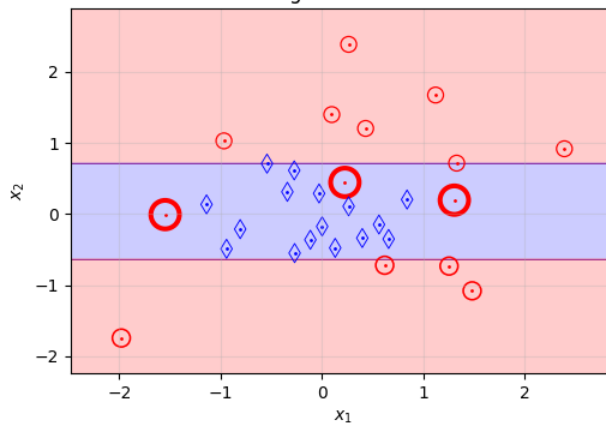
Train error



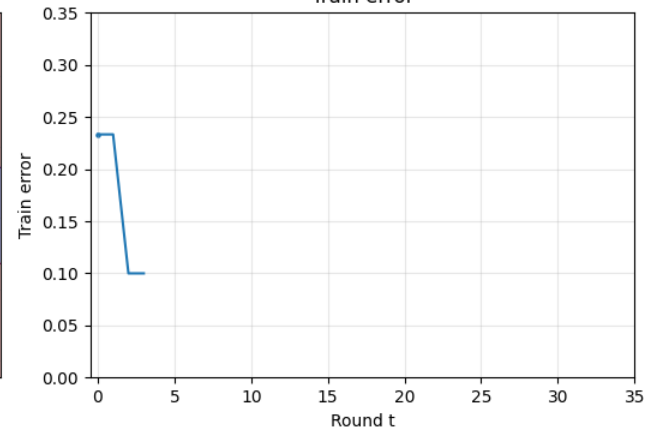
Weak learner at $t=4$



Strong learner at $t=4$

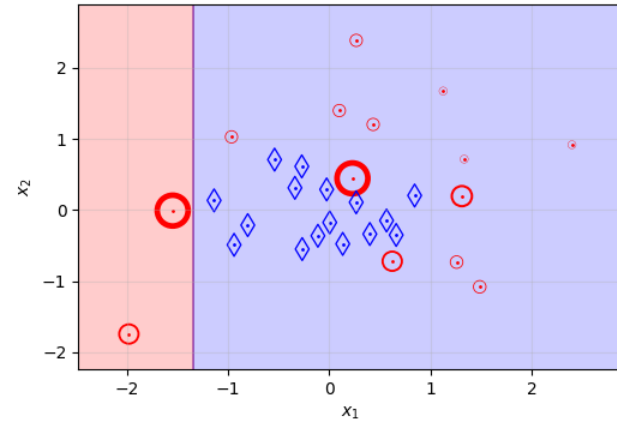


Train error

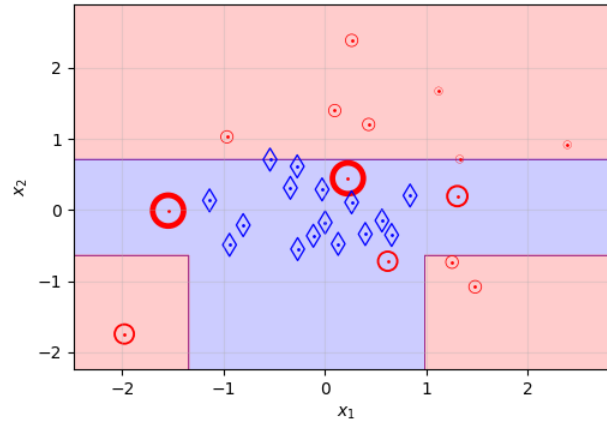


AdaBoost – Visualization

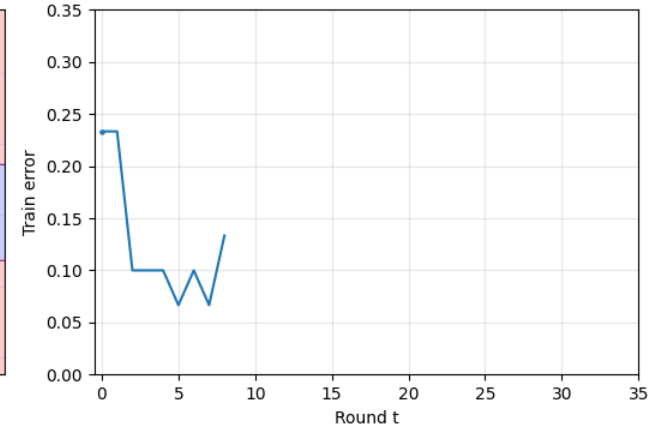
Weak learner at $t=9$



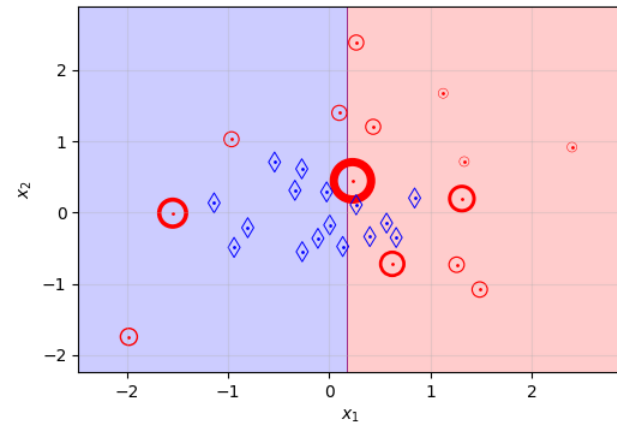
Strong learner at $t=9$



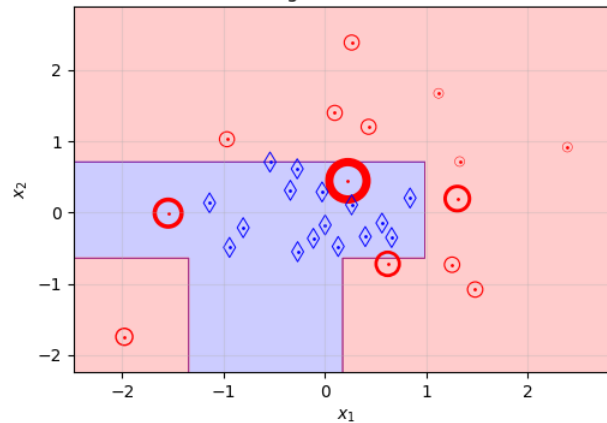
Train error



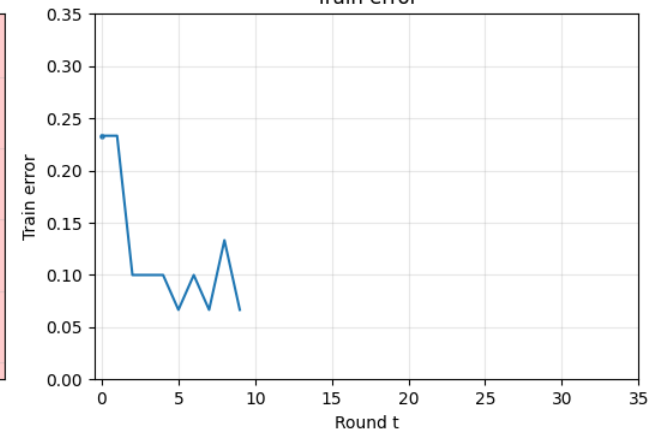
Weak learner at $t=10$



Strong learner at $t=10$

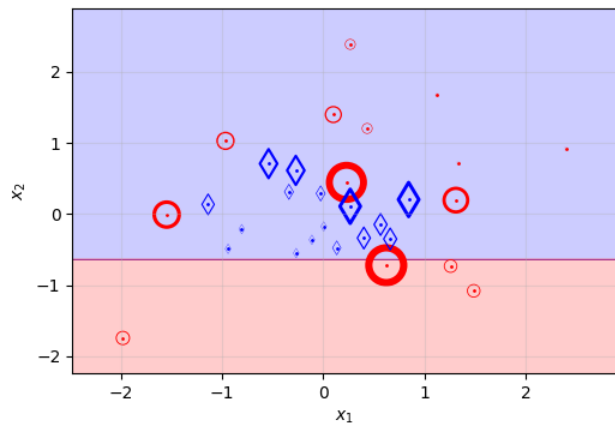


Train error

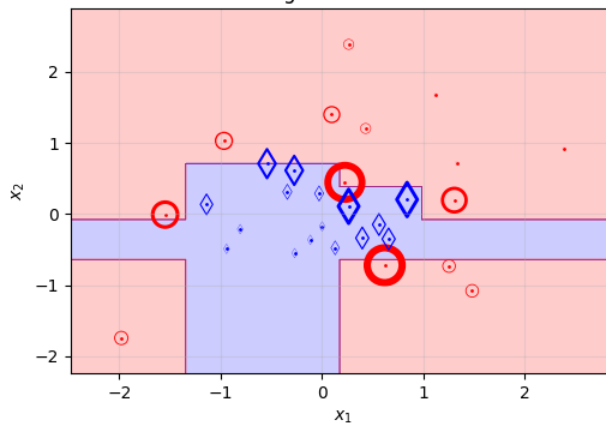


AdaBoost – Visualization

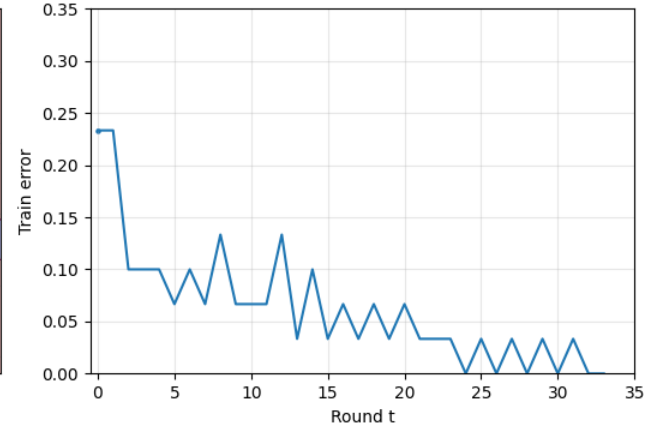
Weak learner at $t=34$



Strong learner at $t=34$

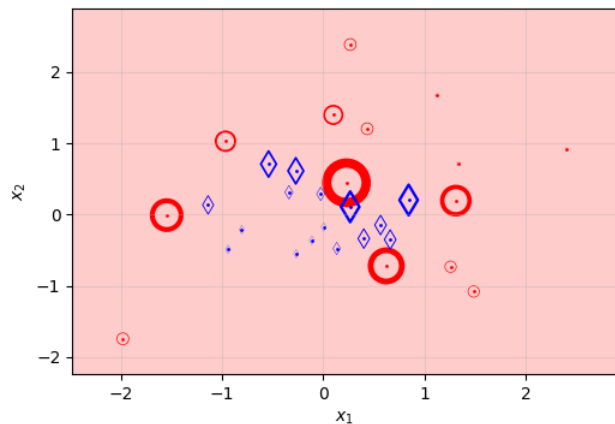


Train error

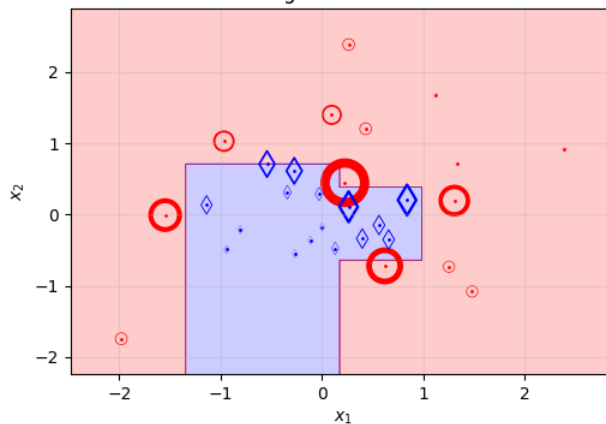


Where do you think we'll get better generalization?

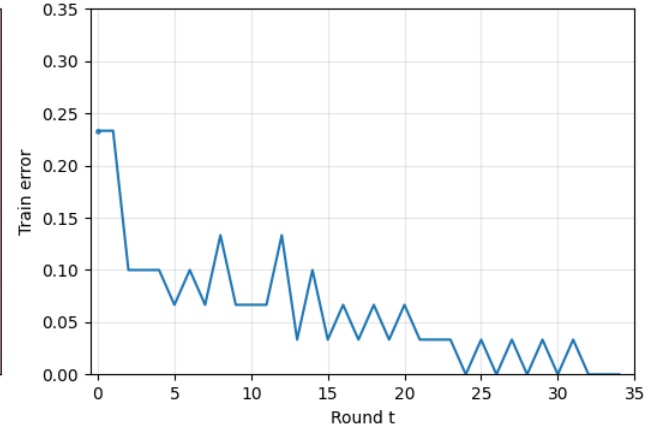
Weak learner at $t=35$



Strong learner at $t=35$



Train error



Detailed Example

Some students want to avoid “hard” courses.

They want to classify courses as “easy” or “hard”.

- First, they collect (very little) data.

Course ID	Hard?	Final exam?	Theoretical?	Advanced?	HW Number
1	1	-1	1	1	1
2	1	-1	1	-1	3
3	1	-1	1	-1	5
4	1	1	-1	-1	5
5	1	1	-1	1	5
6	-1	-1	1	-1	1
7	-1	-1	-1	-1	3

Model choice

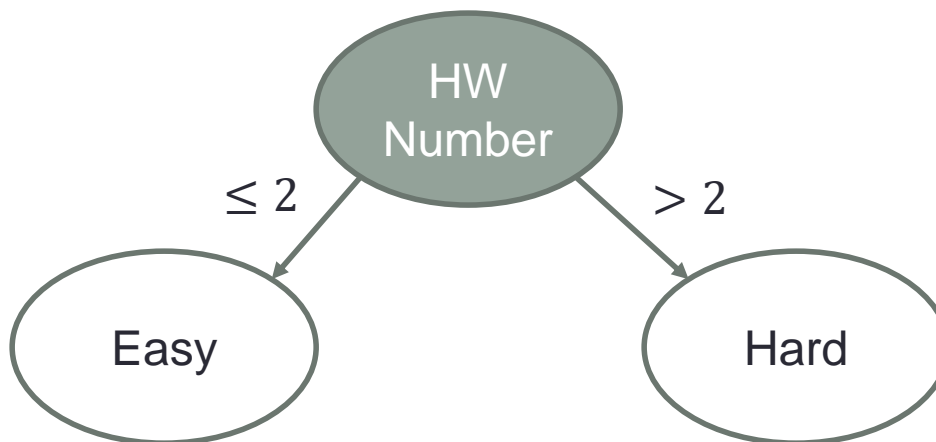
Some students want to avoid “hard” courses.

They want to classify courses as “easy” or “hard”.

- First, they collect (very little) data.

Course ID	Hard?	Final exam?	Theoretical?	Advanced?	HW Number
:	:	:	:	:	:

- Then, they train a model using AdaBoost with decision stumps.



After some preprocessing, we propose the following weak decision stumps classifiers.

Classifier	Attribute	Value	Misclassified
A	Constant “Hard”		
B	Theoretical	1	
C	Advanced	1	
D	# HW	> 2	
E	# HW	> 4	

Proposed weak classifiers

Classifier	Attribute	Value	Misclassified
A	Constant "Hard"		6, 7
B	Theoretical	1	4, 5, 6
C	Advanced	1	2, 3, 4
D	# HW	> 2	1, 7
E	# HW	> 4	1, 2

AdaBoost

- We are ready to run AdaBoost.

Initialize $D^{(1)} = \left(\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m}\right)$

For $t=1, \dots, T$:

Learn (weak) model

$$h_t = \mathcal{A}(S, D^{(t)})$$

Compute error on current distribution

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

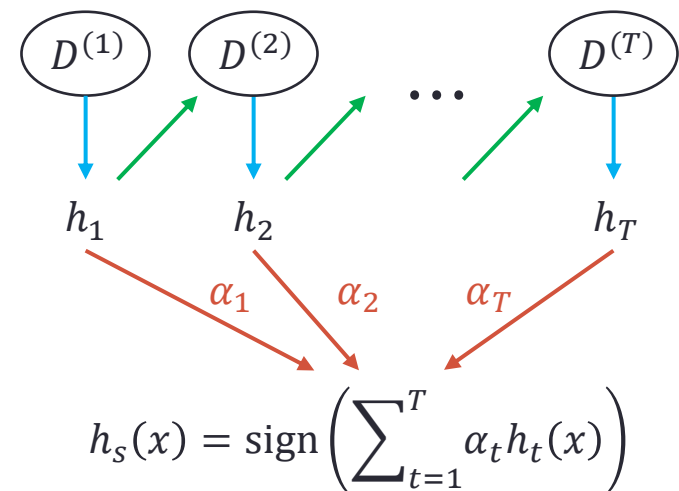
Update weights and distribution

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{\sum_j D_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))}$$

Return final hypothesis

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



- Fill in the classifiers, errors, weights, and distributions of the first 3 rounds.

$t = 1$

$D_1^{(t)}$ $1/7$

$D_2^{(t)}$ $1/7$

$D_3^{(t)}$ $1/7$

$D_4^{(t)}$ $1/7$

$D_5^{(t)}$ $1/7$

$D_6^{(t)}$ $1/7$

$D_7^{(t)}$ $1/7$

 h_t pick the classifier with the lowest weight

ϵ_t $1/7 + 1/7$

α_t $1/2 * \ln(5/2)$

Initialize $D^{(1)} = \left(\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m}\right)$ For $t=1, \dots, T$:

Learn (weak) model

$h_t = \mathcal{A}(S, D^{(t)})$

Compute error on current distribution

$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$

Update weights and distribution

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

$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$

Possible classifiers (ERM)

#	Attribute	Value	Misclassified	ϵ
A	Constant "Hard"		6, 7	
B	Theoretical	1	4, 5, 6	
C	Advanced	1	2, 3, 4	
D	# HW	> 2	1, 7	
E	# HW	> 4	1, 2	

$t = 1$ $t = 2$ $D_1^{(t)}$ $1/7$ $1/10$  $D_2^{(t)}$ $1/7$ $1/10$  $D_3^{(t)}$ $1/7$ $1/10$  $D_4^{(t)}$ $1/7$ $1/10$  $D_5^{(t)}$ $1/7$ $1/10$  $D_6^{(t)}$ $1/7$ $2.5/10$  $D_7^{(t)}$ $1/7$ $2.5/10$  h_t **A** ϵ_t $2/7$ α_t $0.5 \ln \frac{5}{2}$

Updating the distribution

- We used the following weak classifier:

#	Attribute	Value	Misclassified
A	Constant "Hard"		6, 7 => we wan't to up their weight

- We wish to compute the new distribution

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t \underbrace{y_i h_t(x_i)}_{=\{1, -1\}})}{Z_t}$$

- First, let us compute the unnormalized distribution

$\frac{1}{Z_t}$ *normalization*

$$D_i^{(2)} \propto \begin{cases} D_i^{(1)} \exp(\alpha_1) = \frac{1}{7} e^{\frac{1}{2} \ln(\frac{5}{2})} = \frac{1}{7} \cdot \sqrt{\frac{5}{2}} \propto \sqrt{\frac{5}{2}} \propto 5 & i = 6, 7 \\ D_i^{(1)} \exp(-\alpha_1) = \frac{1}{7} e^{-\frac{1}{2} \ln(\frac{5}{2})} = \frac{1}{7} \sqrt{\frac{2}{5}} \propto \sqrt{\frac{2}{5}} \propto 2 & \text{else} \end{cases}$$

*sum = 5 * 2 + 2 * 5 = 20 => $D_4^{(2)} = D_6^{(2)}: \frac{5}{20} = \frac{2.5}{10}$*

- Now, fill in the normalized distribution

	$t = 1$	$t = 2$
$D_1^{(t)}$	$1/7$	$1/10$
$D_2^{(t)}$	$1/7$	$1/10$
$D_3^{(t)}$	$1/7$	$1/10$
$D_4^{(t)}$	$1/7$	$1/10$
$D_5^{(t)}$	$1/7$	$1/10$
$D_6^{(t)}$	$1/7$	$2.5/10$
$D_7^{(t)}$	$1/7$	$2.5/10$
h_t	A	
ϵ_t	$2/7$	
α_t	$0.5 \ln \frac{5}{2}$	



Initialize $D^{(1)} = \left(\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m}\right)$

For $t=1, \dots, T$:

Learn (weak) model

$$h_t = \mathcal{A}(S, D^{(t)})$$

Compute error on current distribution

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

Update weights and distribution

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

Possible classifiers (ERM)



#	Attribute	Value	Misclassified	ϵ
A	Constant "Hard"		6, 7	$1/2$
B	Theoretical	1	4, 5, 6	$4.5/10$
C	Advanced	1	2, 3, 4	$3/10$
D	# HW	> 2	1, 7	$3.5/10$
E	# HW	> 4	1, 2	$2/10$

	$t = 1$	$t = 2$
$D_1^{(t)}$	$1/7$	$1/10$
$D_2^{(t)}$	$1/7$	$1/10$
$D_3^{(t)}$	$1/7$	$1/10$
$D_4^{(t)}$	$1/7$	$1/10$
$D_5^{(t)}$	$1/7$	$1/10$
$D_6^{(t)}$	$1/7$	$2.5/10$
$D_7^{(t)}$	$1/7$	$2.5/10$
h_t	A	E
ϵ_t	$2/7$	$1/5$
α_t	$0.5 \ln \frac{5}{2}$	$\ln 2$

Initialize $D^{(1)} = \left(\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m}\right)$

For $t=1, \dots, T$:

Learn (weak) model

$$h_t = \mathcal{A}(S, D^{(t)})$$

Compute error on current distribution

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

Update weights and distribution

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$\epsilon_1 = \frac{1}{10} + \frac{1}{10} = \frac{1}{5}$$

$$D_i^{(3)} = \frac{D_i^{(2)} e^{-\ln(2) \cdot y_i h_i(x_i)}}{Z_+}$$

$$D_2^{(3)} = D_1^{(3)} \propto D_1^{(2)} e^{\ln(2)} = \frac{2}{10} \propto 4$$

$$D_3^{(3)} = D_4^{(3)} = D_7^{(3)} \propto D_3^{(2)} e^{-\ln(2)} = \frac{1}{20} \propto 1$$

$$D_6^{(3)} = D_7^{(3)} \propto D_6^{(2)} e^{\ln(2)} = \frac{2.5}{20} \propto 2.5$$

$$\text{Sum} = 4 \cdot 2 + 3 \cdot 1 + 2 \cdot 2.5 = 16$$

$$D_2^{(3)} = D_1^{(3)} = \frac{4}{16} = \frac{1}{4}$$

$$D_3^{(3)} = D_4^{(3)} = D_7^{(3)} = \frac{1}{16} = \frac{1}{16}$$

$$D_6^{(3)} = D_7^{(3)} = \frac{2.5}{16} = \frac{5}{32}$$

Updating the distribution

- We used the following weak classifier:

#	Attribute	Value	Misclassified
E	# HW	> 4	1, 2

- We attach results for a 3rd iteration of AdaBoost without relevant computations.
- Extra:** manually run this iteration by yourselves and prove the presented results.

	$t = 1$	$t = 2$	$t = 3$
$D_1^{(t)}$	$1/7$	$1/10$	$8/32$
$D_2^{(t)}$	$1/7$	$1/10$	$8/32$
$D_3^{(t)}$	$1/7$	$1/10$	$2/32$
$D_4^{(t)}$	$1/7$	$1/10$	$2/32$
$D_5^{(t)}$	$1/7$	$1/10$	$2/32$
$D_6^{(t)}$	$1/7$	$2.5/10$	$5/32$
$D_7^{(t)}$	$1/7$	$2.5/10$	$5/32$
h_t	A	E	B
ϵ_t	$2/7$	$1/5$	$9/32$
α_t	$0.5 \ln \frac{5}{2}$	$\ln 2$	$0.5 \ln \frac{23}{9}$

Final hypothesis

Computed weights

	$t = 1$	$t = 2$	$t = 3$
h_t	A	E	B
α_t	$0.5 \ln \frac{5}{2}$	$\ln 2$	$0.5 \ln \frac{23}{9}$

Chosen classifiers

#	Attribute	Value	Misclassified
A	Constant "Hard"		6, 7
E	# HW	> 4	1, 2
B	Theoretical	1	4, 5, 6

- Final hypothesis:
$$h_S(x) = \text{sign} \left(0.5 \ln \frac{5}{2} A(x) + \ln 2 E(x) + 0.5 \ln \frac{23}{9} B(x) \right)$$

$$\approx \text{sign}(0.46 + 0.69E(x) + 0.47B(x))$$

where: A, B, E return $+1$ or -1

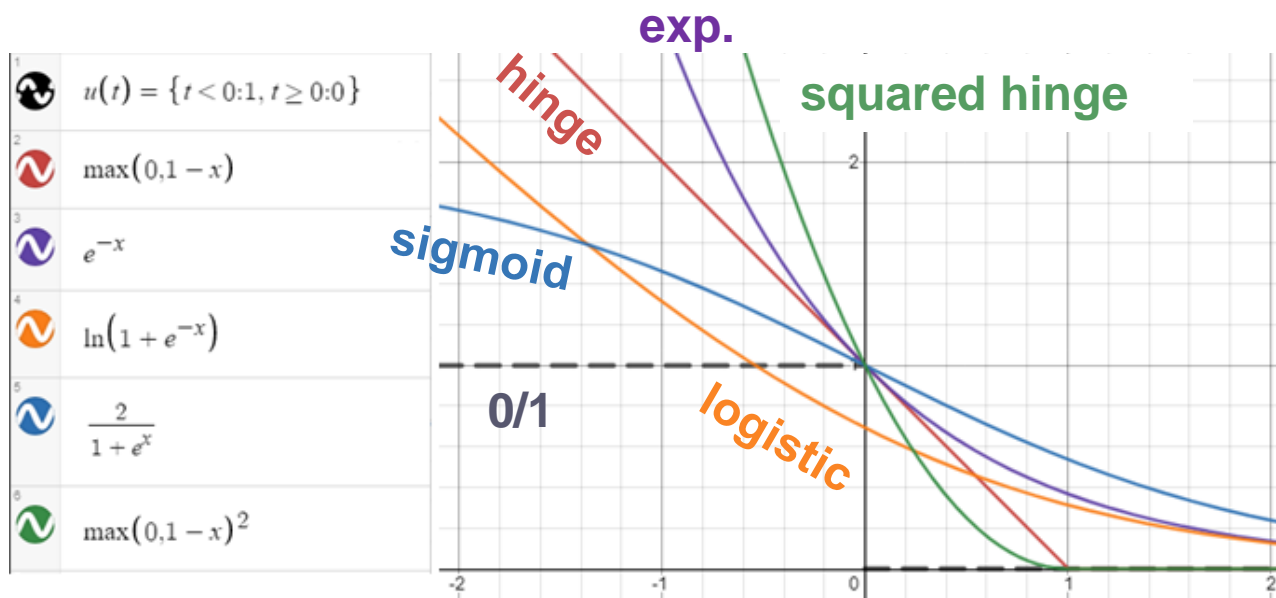
- How many of the data points are classified correctly?
 - Notice: each two weights α 's are larger than the third
 - We get a simple majority vote
- Only one wrong prediction (example 6)

AdaBoost: Guarantees

- AdaBoost constructs a strong hypothesis $h_s(x) = \text{sign}(\sum_{k=1}^t \alpha_k h_k(x))$.
- **Guarantee:** AdaBoost's training error after T iterations
is bounded by $L_S(h_s) \triangleq \frac{1}{m} \sum_i \mathbf{1}_{h_s(x_i) \neq y_i} \leq \exp\{-\gamma^2 T\}$ for some $\gamma \in (0, 1/2)$.
(under mild conditions; without proof)
- **Corollary:** AdaBoost reaches zero training error eventually. (under the same conditions)

AdaBoost optimizes the exponential loss

- AdaBoost constructs a strong hypothesis $h_s(x) = \text{sign}(\sum_{k=1}^t \alpha_k h_k(x))$.
 Focus on the “unthresholded” hypothesis
- Why does it work?
- We will show that AdaBoost greedily optimizes the exponential loss.
- Recall the exponential loss: $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \ell_{\text{exp}}(x_i, y_i) = \frac{1}{m} \sum_{i=1}^m e^{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)}$
- The exponential loss is a proxy to the zero-one loss: $\mathcal{L}_{\text{exp}}(h_s) \geq \mathcal{L}_{0/1}(h_s)$



AdaBoost optimizes the exponential loss

- **Goal:** Show AdaBoost optimizes $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \exp \{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\}$.
- **Assume:** $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ were already chosen.
- **Question:** How to choose h_t, α_t to minimize $\mathcal{L}_{\text{exp}}(h_s)$?
- **Lemma:** $\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \underbrace{\sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}}_{\triangleq \epsilon_t}$
(will be proven next)
- **Answer:** Use “two steps” greedy optimization:
 - Choose h_t minimizing the weighted error ϵ_t , e.g., with ERM w.r.t. $D^{(t)}$.
 - Choose α_t minimizing $\mathcal{L}_{\text{exp}}^{(t)}$ given $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ and h_t .
- 1. Derive $\frac{\partial}{\partial \alpha_t} \mathcal{L}_{\text{exp}}^{(t)}$ and use it to prove that the choice $\alpha_t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right)$ is optimal.

Initialize $D^{(1)} = \left(\frac{1}{m}, \dots, \frac{1}{m} \right)$

For $t=1, \dots, T$:

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

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

AdaBoost optimizes the exponential loss

- Goal: Show AdaBoost optimizes $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \exp \{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\}$.
 - Assume: $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ were already chosen.
 - Question: How to choose h_t, α_t to minimize $\mathcal{L}_{\text{exp}}(h_s)$?
 - **Lemma:** $\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \underbrace{\sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}}_{\triangleq \epsilon_t}$
 - **Proof:**
2. Extra: Show that $D_i^{(t)} = c \cdot \exp(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i))$, for $c > 0$

$$\begin{aligned}
 D_i^{(t)} &= \frac{1}{Z_{t-1}} D_i^{(t-1)} \exp(-\alpha_{t-1} y_i h_{t-1}(x_i)) \\
 &= \frac{1}{Z_{t-1} Z_{t-2}} D_i^{(t-2)} \exp(-\alpha_{t-2} y_i h_{t-2}(x_i)) \exp(-\alpha_{t-1} y_i h_{t-1}(x_i)) \\
 &= \frac{1}{Z_{t-1} Z_{t-2}} D_i^{(t-2)} \exp(-y_i (\alpha_{t-2} h_{t-2}(x_i) - \alpha_{t-1} h_{t-1}(x_i))) \\
 &= \dots = \frac{1}{\prod_{k=1}^{t-1} Z_k} \underbrace{D_i^{(0)}}_{=1/m} \exp\left(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i)\right) \quad \blacksquare
 \end{aligned}$$

Initialize $D^{(1)} = \left(\frac{1}{m}, \dots, \frac{1}{m}\right)$

For $t=1, \dots, T$:

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

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

AdaBoost optimizes the exponential loss

- Goal: Show AdaBoost optimizes $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \exp \{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\}$.
- Assume: $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ were already chosen.
- Question: How to choose h_t, α_t to minimize $\mathcal{L}_{\text{exp}}(h_s)$?
- **Lemma:** $\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \underbrace{\sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}}_{\triangleq \epsilon_t}$
- **Proof:**
 2. Extra: Show that $D_i^{(t)} = c \cdot \exp(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i))$, for $c > 0$
 3. Show that $\mathcal{L}_{\text{exp}}^{(t)} \propto \sum_i D_i^{(t)} \exp\{-y_i \alpha_t h_t(x_i)\}$

Initialize $D^{(1)} = (\frac{1}{m}, \dots, \frac{1}{m})$

For $t=1, \dots, T$:

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

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

$$\begin{aligned} \mathcal{L}_{\text{exp}}^{(t)} &= \sum_i \exp \left\{ -y_i \sum_{k=1}^t \alpha_k h_k(x_i) \right\} = \sum_i \exp \left\{ \left(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i) \right) - y_i \alpha_t h_t(x_i) \right\} \\ &= \sum_i \underbrace{\exp \left\{ -y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i) \right\}}_{=\frac{1}{c} \cdot D_i^{(t)}} \exp \{-y_i \alpha_t h_t(x_i)\} \propto \sum_i D_i^{(t)} \exp \{-y_i \alpha_t h_t(x_i)\} \end{aligned}$$

AdaBoost optimizes the exponential loss

- Goal: Show AdaBoost optimizes $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \exp \{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\}$.
- Assume: $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ were already chosen.
- Question: How to choose h_t, α_t to minimize $\mathcal{L}_{\text{exp}}(h_s)$?
- **Lemma:** $\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \underbrace{\sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}}_{\triangleq \epsilon_t}$
- **Proof:**
 2. Extra: Show that $D_i^{(t)} = c \cdot \exp(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i))$, for $c > 0$
 3. Show that $\mathcal{L}_{\text{exp}}^{(t)} \propto \sum_i D_i^{(t)} \exp\{-y_i \alpha_t h_t(x_i)\}$
 4. Extra: Prove the following form of the exp. loss after t rounds is

Initialize $D^{(1)} = (\frac{1}{m}, \dots, \frac{1}{m})$
 For $t=1, \dots, T$:

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

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

$$\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} \sum_i D_i^{(t)} + (e^{\alpha_t} - e^{-\alpha_t}) \sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}$$

Do so by filling in the blanks:

$$\begin{aligned} \mathcal{L}_{\text{exp}}^{(t)} &\propto \sum_i D_i^{(t)} \exp\{-y_i \alpha_t h_t(x_i)\} = \dots = e^{-\alpha_t} \sum_{i: y_i = h_t(x_i)} D_i^{(t)} + e^{\alpha_t} \sum_{i: y_i \neq h_t(x_i)} D_i^{(t)} \\ &= \dots = e^{-\alpha_t} \sum_i D_i^{(t)} + (e^{\alpha_t} - e^{-\alpha_t}) \sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i} \end{aligned}$$

AdaBoost optimizes the exponential loss

- Goal: Show AdaBoost optimizes $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \exp \{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\}$.
- Assume: $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ were already chosen.
- Question: How to choose h_t, α_t to minimize $\mathcal{L}_{\text{exp}}(h_s)$?
- **Lemma:** $\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \underbrace{\sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}}_{\triangleq \epsilon_t}$
- **Proof:**
 2. Extra: Show that $D_i^{(t)} = c \cdot \exp(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i))$, for $c > 0$
 3. Show that $\mathcal{L}_{\text{exp}}^{(t)} \propto \sum_i D_i^{(t)} \exp\{-y_i \alpha_t h_t(x_i)\}$
 4. Extra: Prove the following form of the exp. loss after t rounds is

$$\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} \underbrace{\sum_i D_i^{(t)}}_{=?} + (e^{\alpha_t} - e^{-\alpha_t}) \sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}$$

Initialize $D^{(1)} = (\frac{1}{m}, \dots, \frac{1}{m})$

For $t=1, \dots, T$:

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

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

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

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

AdaBoost optimizes the exponential loss

- **Goal:** Show AdaBoost optimizes $\mathcal{L}_{\text{exp}}(h_s) = \frac{1}{m} \sum_{i=1}^m \exp \{-y_i \sum_{k=1}^t \alpha_k h_k(x_i)\}$.
- **Assume:** $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ were already chosen.
- **Question:** How to choose h_t, α_t to minimize $\mathcal{L}_{\text{exp}}(h_s)$?
- **Lemma:** $\mathcal{L}_{\text{exp}}^{(t)} \propto e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \underbrace{\sum_i D_i^{(t)} \mathbf{1}_{h_t(x_i) \neq y_i}}_{\triangleq \epsilon_t}$
(proven)
- **Answer:** Use “two steps” greedy optimization:
 - Choose h_t minimizing the weighted error ϵ_t , e.g., with ERM w.r.t. $D^{(t)}$.
 - Choose α_t minimizing $\mathcal{L}_{\text{exp}}^{(t)}$ given $h_1, \alpha_1, \dots, h_{t-1}, \alpha_{t-1}$ and h_t .

Initialize $D^{(1)} = \left(\frac{1}{m}, \dots, \frac{1}{m}\right)$

For $t=1, \dots, T$:

$$h_t = \mathcal{A}(S, D^{(t)})$$

$$\epsilon_t = \sum_i D_i^{(t)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

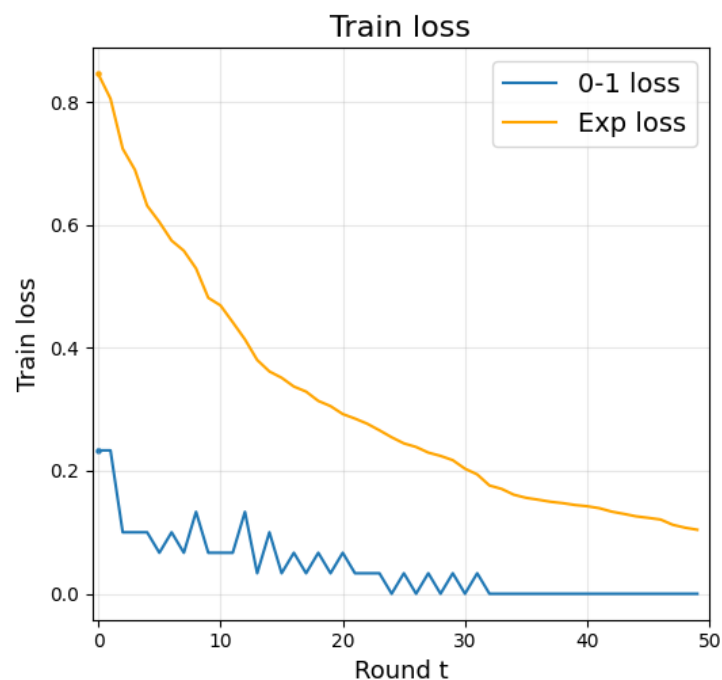
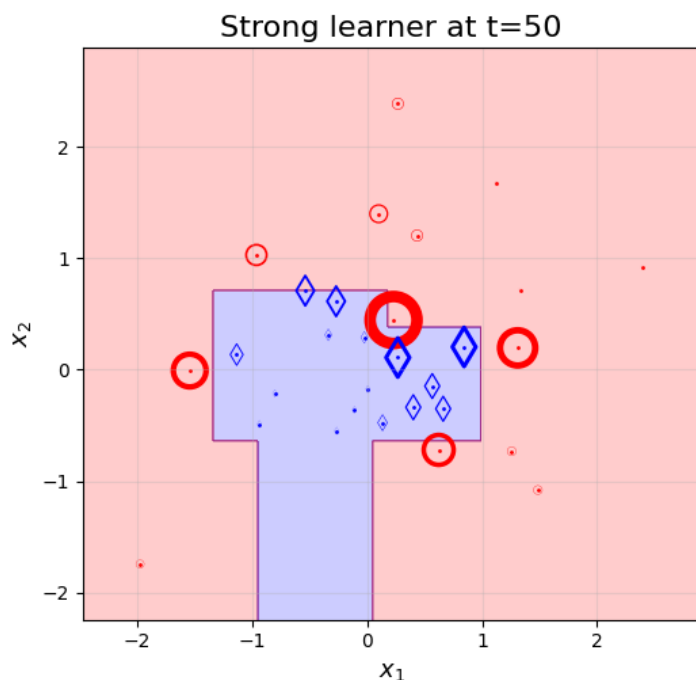
$$\alpha_t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right)$$

$$D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

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

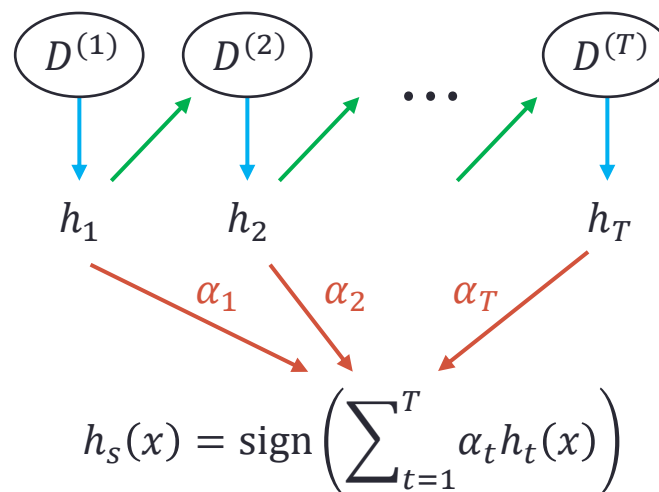
AdaBoost optimizes the exponential loss

- The simulation from earlier can also visualize the decaying exponential loss.
- After 0 classification training loss, the exponential loss keeps decreasing



Summary

- In boosting we create a **strong** classifier by combining multiple **weak** classifiers
- Each weak classifier is trained on an updated distribution where misclassified data points get larger probability mass



- AdaBoost is a boosting algorithm that greedily optimizes the **exponential loss**