

Announcements

Assignments

- HW10 (programming + “written”)
 - Due Thu 4/30, 11:59 pm

Final Exam

- Stay tuned to Piazza for details
- Date: Mon 5/11, 5:30 – 8:30 pm
- Format: “online assignment” in Gradescope
- Practice exam: Out later this week
- Recitation this Friday: Review session

Announcements

Course Evaluation Survey

- <https://cmu.smartevals.com/>
- Take time now to fill this out please

Piazza Poll 1

Did you fill out the FCE survey?

A. Yes

B. Why not?

Introduction to Machine Learning

Ensemble Methods

Instructor: Pat Virtue

Recommender Systems

Netflix Prize **COMPLETED**

Home Rules Leaderboard Update

Leaderboard

Showing Test Score. [Click here to show quiz score](#)

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries !	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

Top performing systems were ensembles

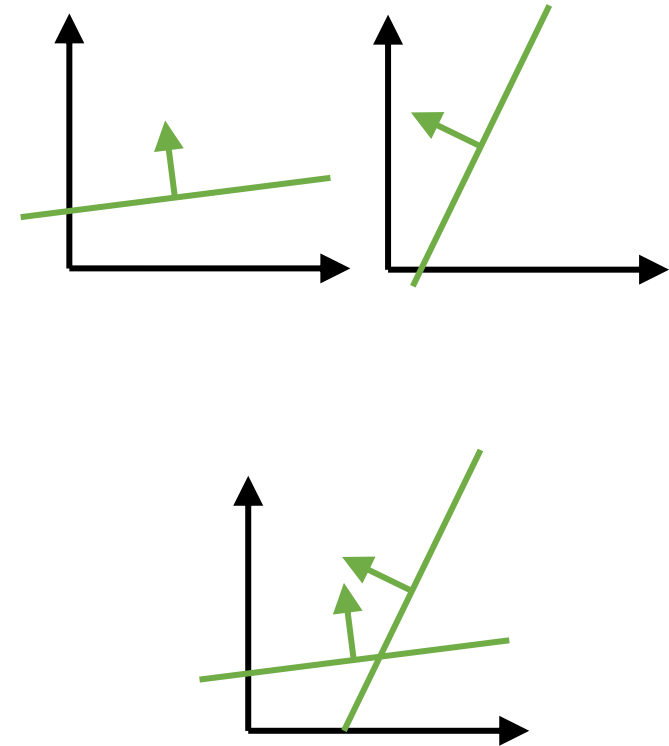
Ensemble Methods

Given: pool A of learners (that you know nothing about)

Goal: design a new learner that uses the predictions of the pool to make new predictions

Techniques

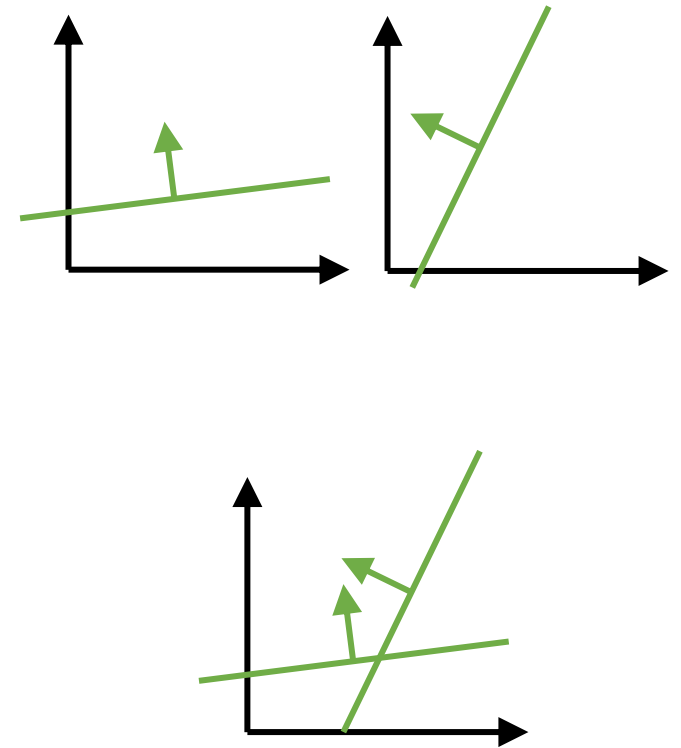
- Bagging (Bootstrapping)
 - e.g. Random Forests
- Boosting
 - e.g. Adaboost



Bagging

[Breiman, 1996]

1. Run independent weak learners on bootstrap replicates (sample with replacement) of the training set
2. Average/vote over weak hypotheses



Bagging

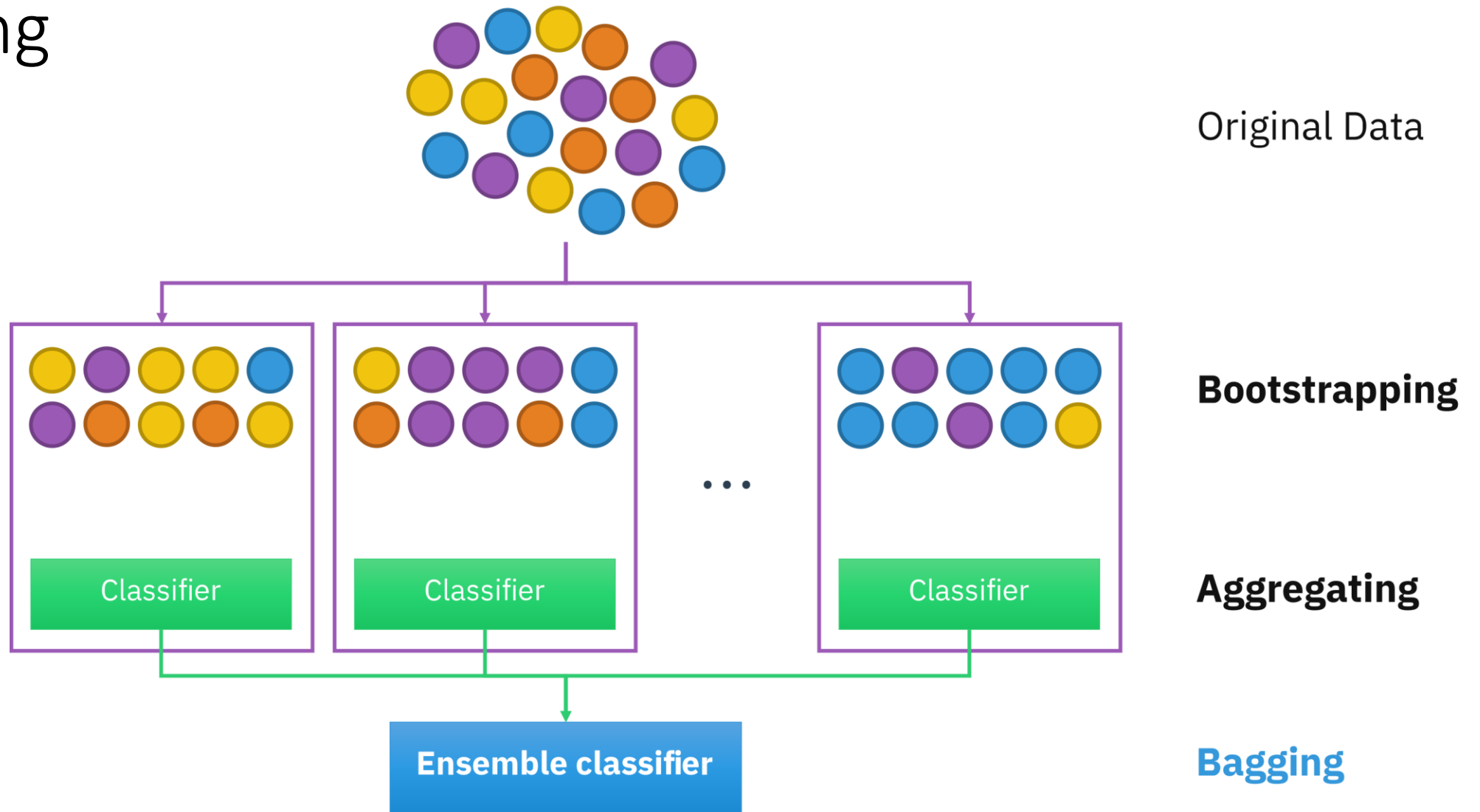
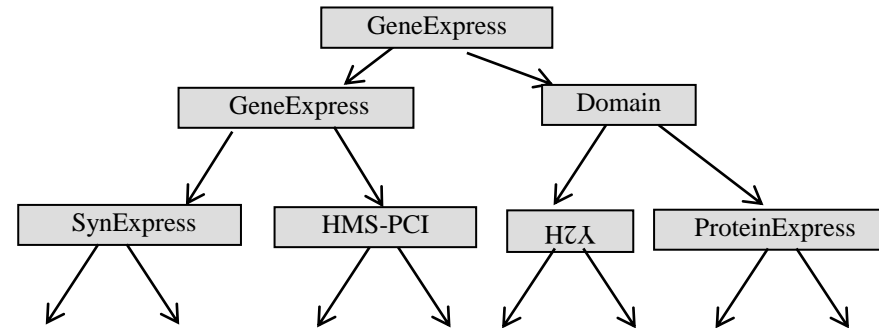
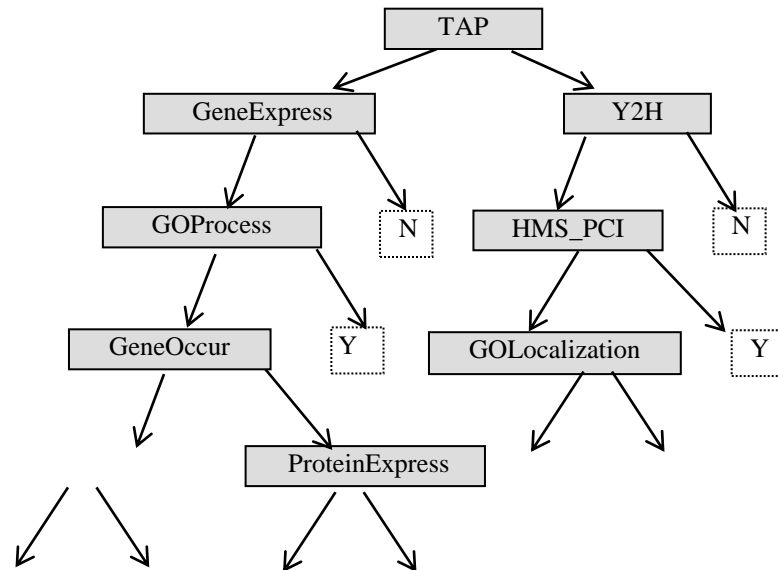


Image: https://en.m.wikipedia.org/wiki/Bootstrap_aggregating

Random Forests

Bagging over Features

- A collection of decision trees
- For each tree we select a subset of the attributes (recommended square root of $|A|$) and build tree using just these attributes
- An input sample is classified using majority voting



Weighted Majority Algorithm

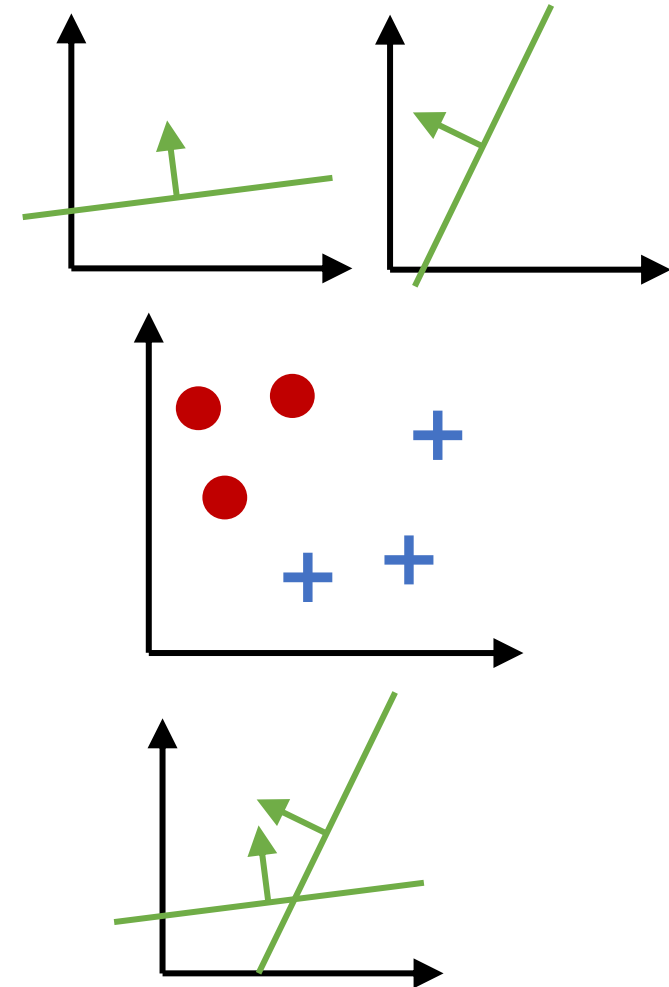
(Littlestone & Warmuth, 1994)

Given: pool A of binary classifiers (that you know nothing about)

Goal: design a new learner that uses the predictions of the pool to make new predictions

Algorithm:

- Initially weight all classifiers equally
- Receive a training example and predict the (weighted) majority vote of the classifiers in the pool
- Down-weight classifiers that contribute to a mistake by a factor of β



Weighted Majority Algorithm

(Littlestone & Warmuth, 1994)

Suppose we have a pool of T binary classifiers $\mathcal{A} = \{h_1, \dots, h_T\}$ where $h_t : \mathbb{R}^M \rightarrow \{+1, -1\}$. Let α_t be the weight for classifier h_t .

Algorithm 1 Weighted Majority Algorithm

- 1: **procedure** WEIGHTEDMAJORITY(\mathcal{A}, β)
- 2: Initialize classifier weights $\alpha_t = 1, \forall t \in \{1, \dots, T\}$
- 3: **for** each training example (\mathbf{x}, y) **do**
- 4: Predict majority vote class (splitting ties randomly)

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

- 5: **if** a mistake is made $\hat{h}(x) \neq y$ **then**
- 6: **for** each classifier $t \in \{1, \dots, T\}$ **do**
- 7: **if** $h_t(x) \neq y$, **then** $\alpha_t \leftarrow \beta \alpha_t$

Adaboost

Comparison

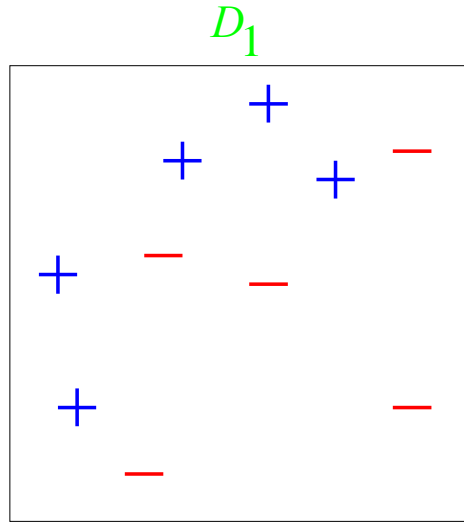
Weighted Majority Algorithm

- an example of an ensemble method
- assumes the classifiers are learned ahead of time
- only learns (majority vote) weight for each classifiers

AdaBoost

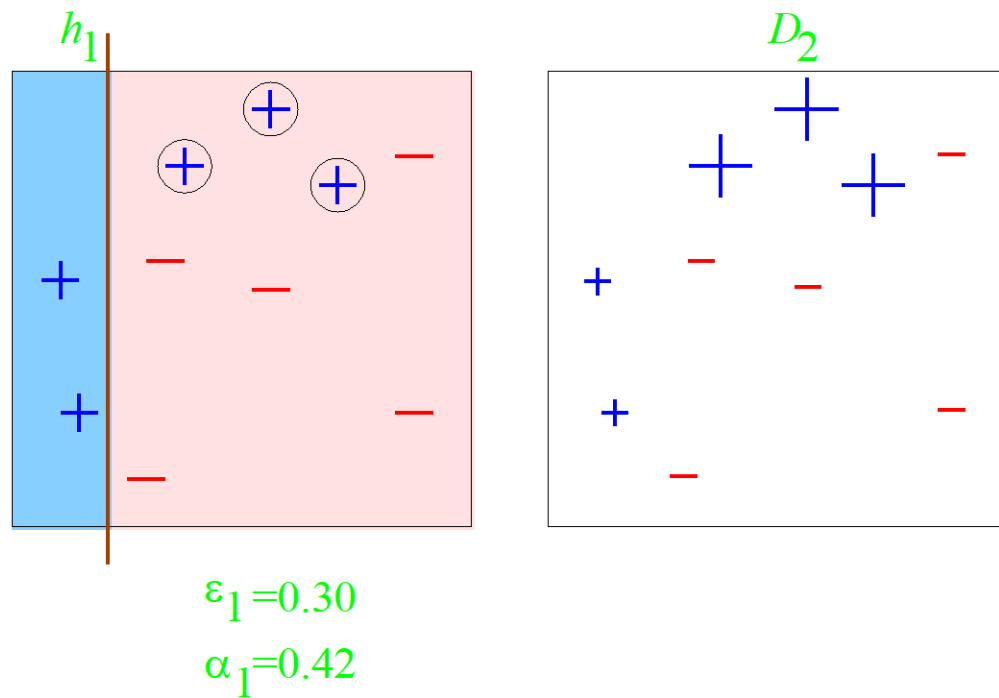
- an example of a boosting method
- simultaneously learns:
 - the classifiers themselves
 - (majority vote) weight for each classifiers

AdaBoost: Toy Example

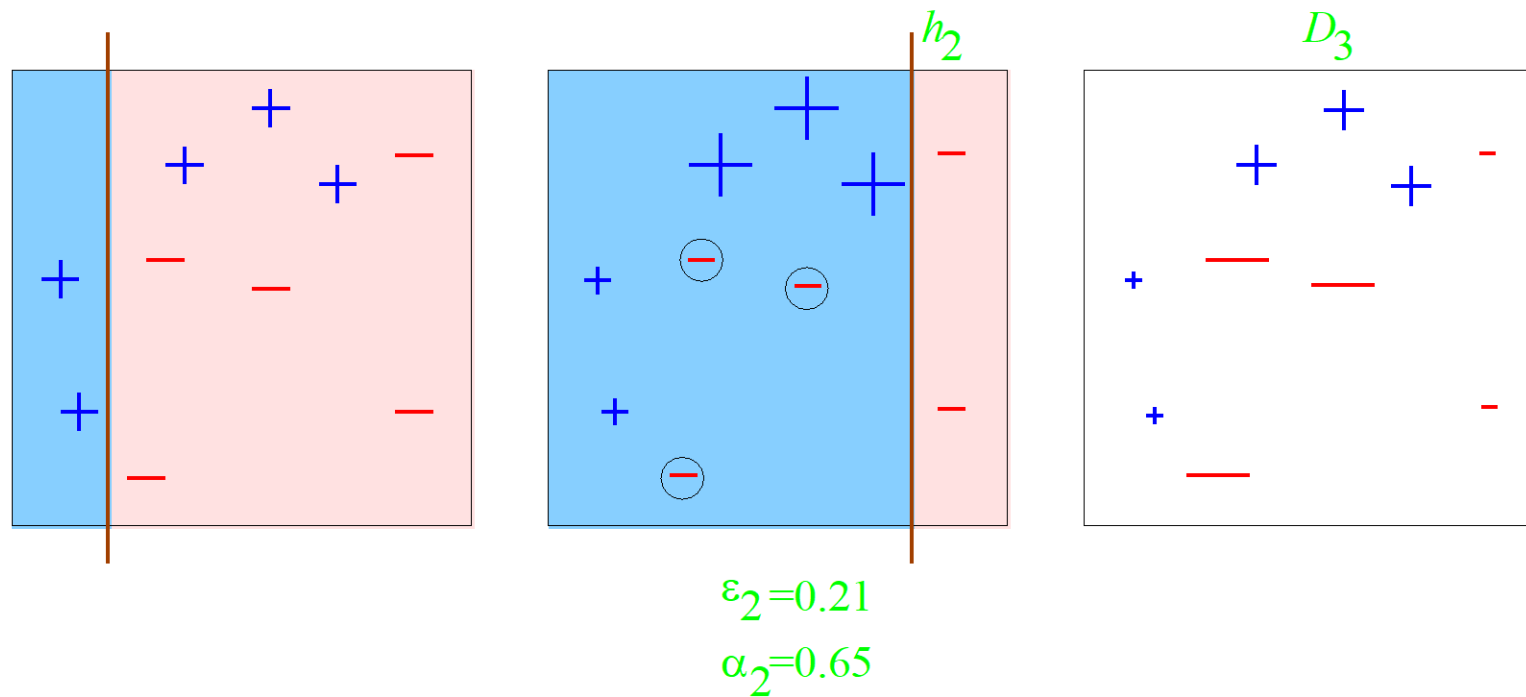


weak classifiers = vertical or horizontal half-planes

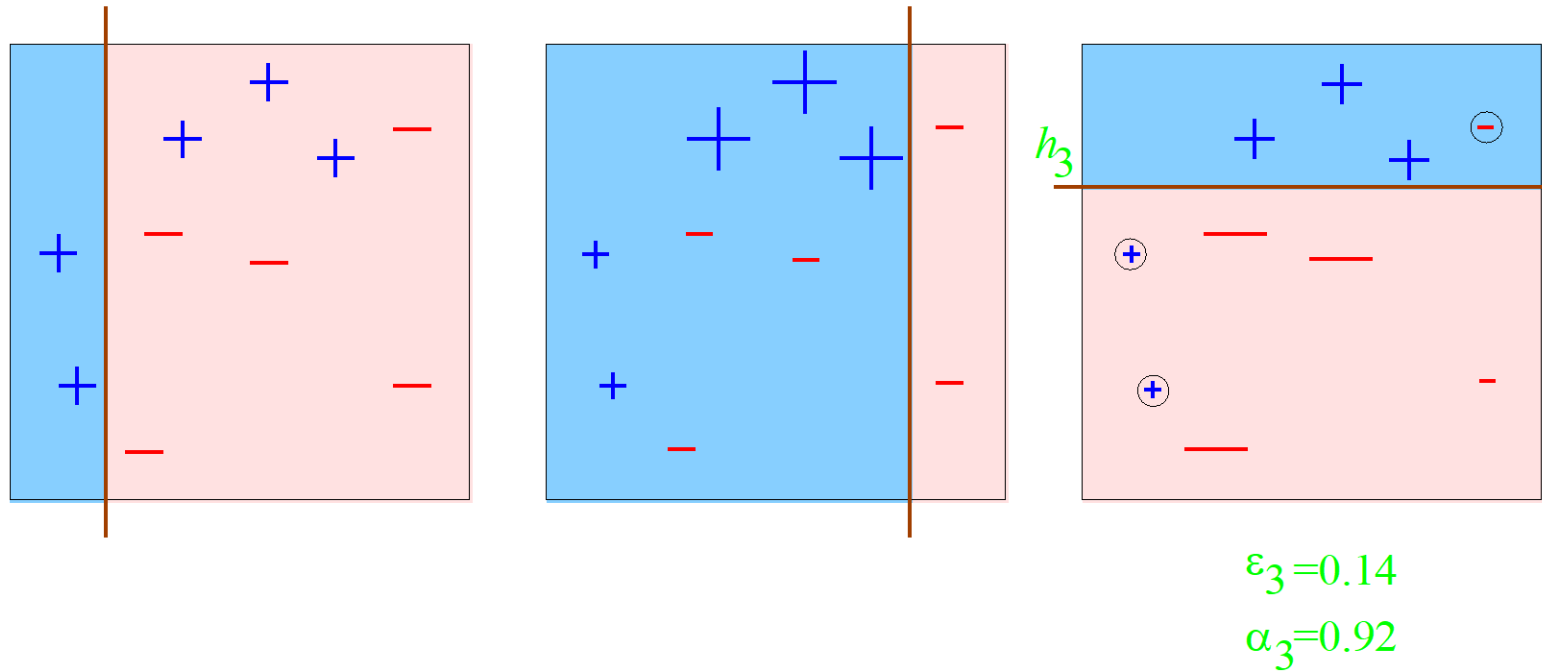
AdaBoost: Toy Example



AdaBoost: Toy Example

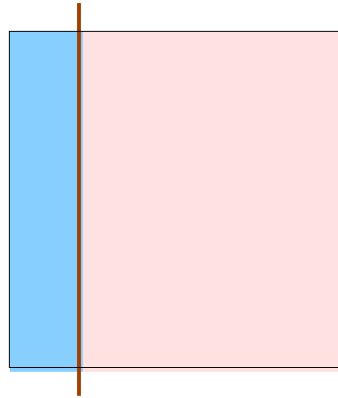


AdaBoost: Toy Example

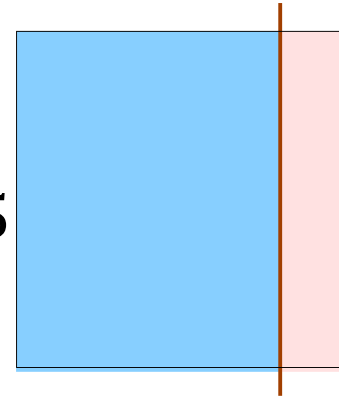


Piazza Poll 2-4

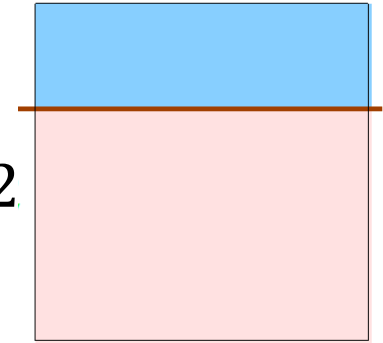
$$H_{final} = \text{sign}(0.42$$



+0.65



+0.92



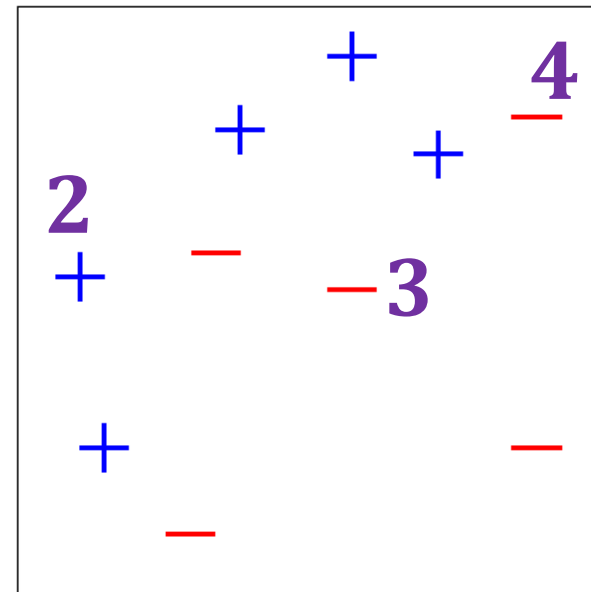
What is the final classification of points 2, 3, 4?

A.

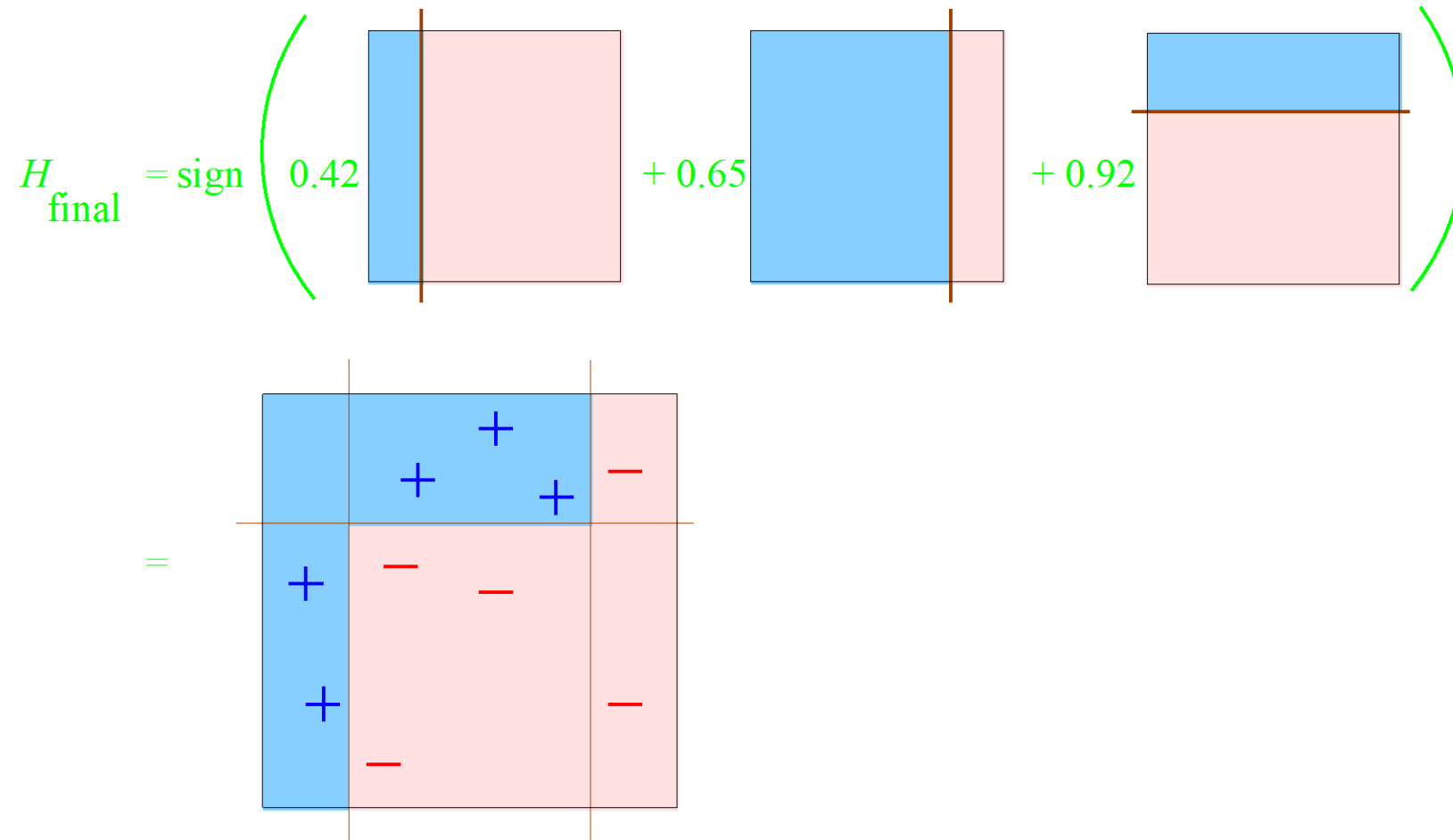
=

B.

C.



AdaBoost: Toy Example



AdaBoost

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

Initialize $D_1(i) = 1/m$.

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

- Train weak learner using distribution D_t .
- Get weak hypothesis $h_t : X \rightarrow \{-1, +1\}$ with error

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i].$$

- Choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$.
- Update:

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(x_i) \neq y_i \end{cases} \\ &= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \end{aligned}$$

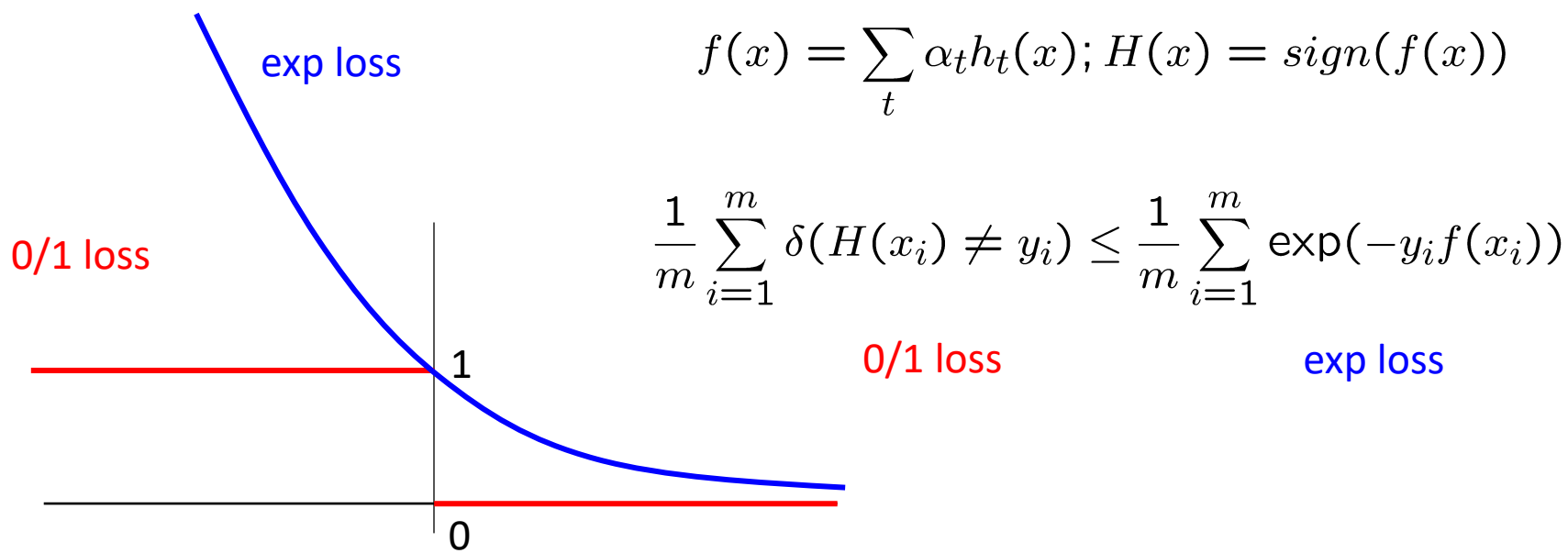
where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

Output the final hypothesis:

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

Analysis for Boosting

Choice of α_t and hypothesis h_t obtained by coordinate descent on exp loss
(convex upper bound on 0/1 loss)



Analysis for Boosting

Analysis reveals:

If each weak learner h_t is slightly better than random guessing ($\epsilon_t < 0.5$),

then training error of AdaBoost decays exponentially fast in number of rounds T .

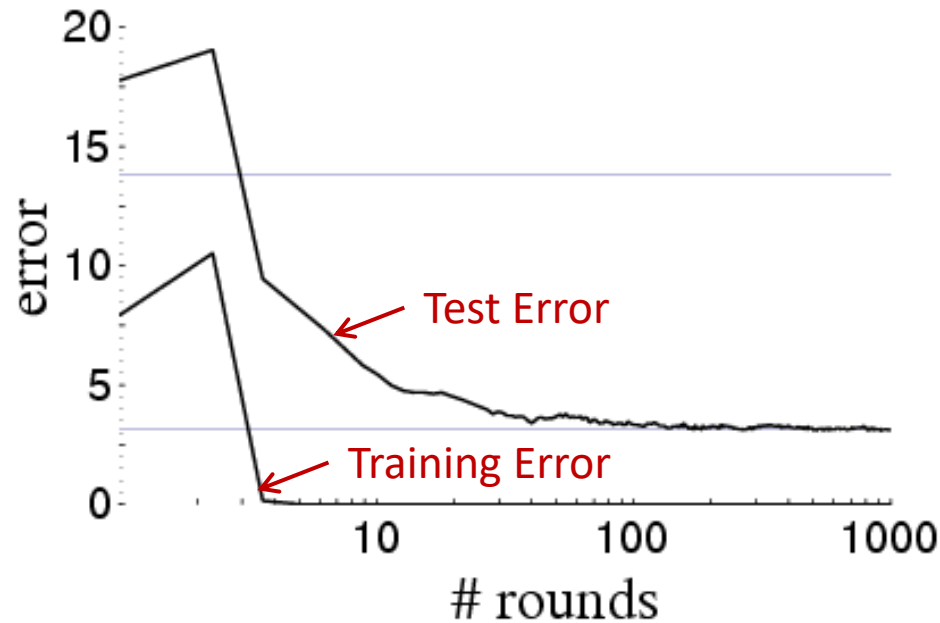
$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \exp \left(-2 \sum_{t=1}^T (1/2 - \epsilon_t)^2 \right)$$

Training Error

What about test error?

Boosting results – Digit recognition

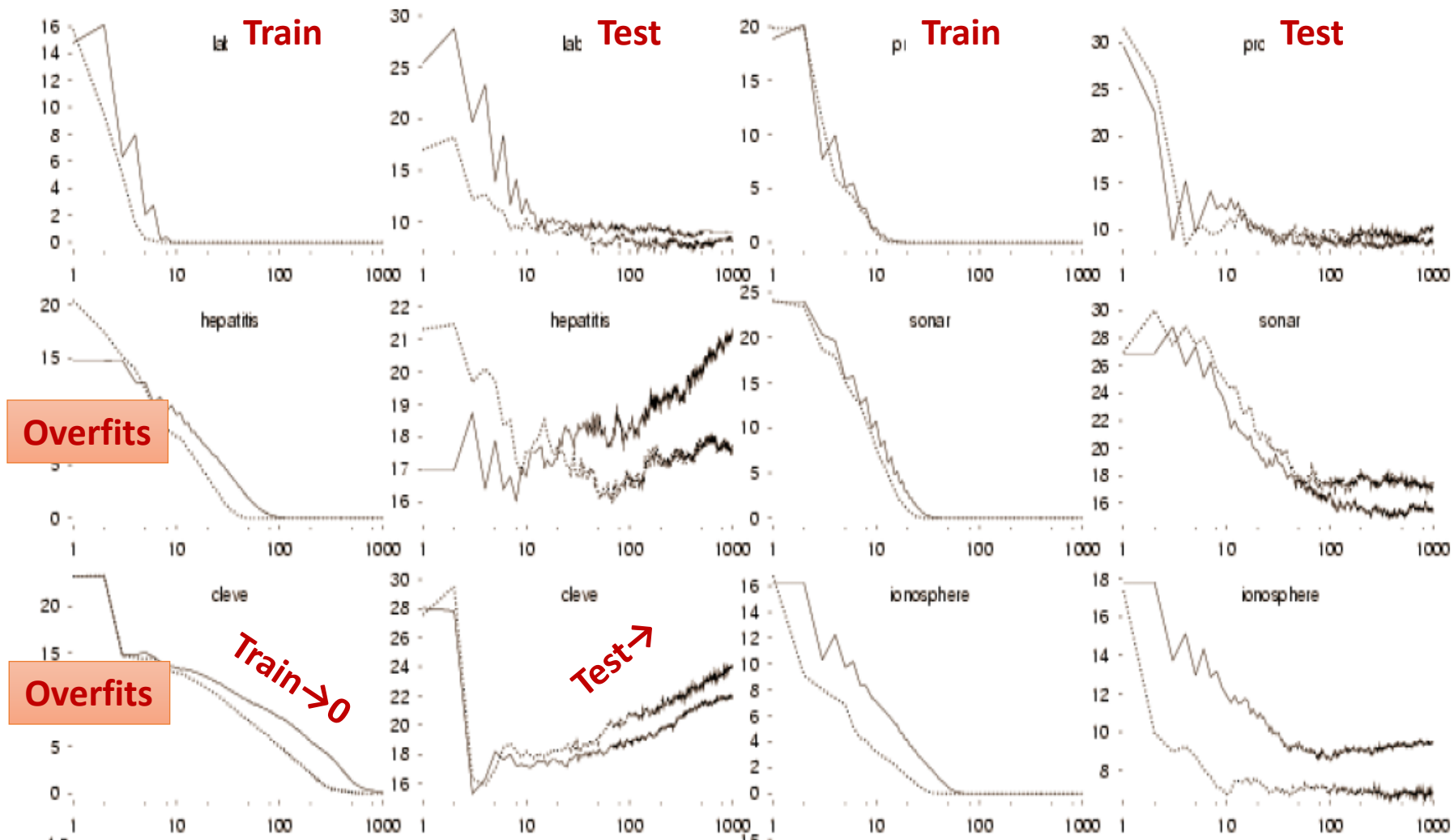
[Schapire, 1989]



Boosting often,

- Robust to overfitting
- Test set error decreases even after training error is zero
- If margin between classes is large, subsequent weak learners agree and hence more rounds does not necessarily imply that final classifier is getting more complex.

AdaBoost and AdaBoost.MH on Train (left) and Test (right) data from Irvine repository. [Schapire and Singer, ML 1999]



Boosting can overfit if margin between classes is too small (high label noise) or weak learners are too complex.

AdaBoost

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

Initialize $D_1(i) = 1/m$.

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

- Train weak learner using distribution D_t .
- Get weak hypothesis $h_t : X \rightarrow \{-1, +1\}$ with error

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i].$$

- Choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$.
- Update:

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(x_i) \neq y_i \end{cases} \\ &= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \end{aligned}$$

where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

Output the final hypothesis:

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

Dumb classifiers made Smart

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \prod_t Z_t = \prod_t \sqrt{1 - (1 - 2\epsilon_t)^2}$$
$$\leq \exp \left(-2 \sum_{t=1}^T \underbrace{(1/2 - \epsilon_t)^2}_{\text{grows as } \epsilon_t \text{ moves away from } 1/2} \right)$$

If each classifier is (at least slightly) better than random $\epsilon_t < 0.5$

AdaBoost will achieve zero training error exponentially fast (in number of rounds T) !!

What about test error?

Generalization Error Bounds

[Freund & Schapire'95]

$$error_{true}(H) \leq error_{train}(H) + \tilde{O} \left(\sqrt{\frac{Td}{m}} \right) \quad \text{With high probability}$$

Generalization Error Bounds

[Freund & Schapire'95]

$$error_{true}(H) \leq error_{train}(H) + \tilde{O} \left(\sqrt{\frac{Td}{m}} \right) \quad \text{With high probability}$$

Boosting can overfit if T is large

Boosting often,

Contradicts experimental results

- Robust to overfitting
- Test set error decreases even after training error is zero

Need better analysis tools – margin based bounds

Margin Based Bounds

[Schapire, Freund, Bartlett, Lee'98]

$$error_{true}(H) \leq \hat{\Pr} [\text{margin}_f(x, y) \leq \theta] + \tilde{O} \left(\sqrt{\frac{d}{m\theta^2}} \right) \quad \text{With high probability}$$

Boosting increases the margin very aggressively since it concentrates on the hardest examples.

If margin is large, more weak learners agree and hence more rounds does not necessarily imply that final classifier is getting more complex.

Bound is independent of number of rounds T !

Boosting can still overfit if margin is too small (can perform arbitrarily close to random guessing) or weak learners are too complex

Boosting and Logistic Regression

Logistic regression assumes:

$$P(Y = 1|X) = \frac{1}{1 + \exp(f(x))}$$

$$f(x) = w_0 + \sum_j w_j x_j$$

And tries to maximize data likelihood:

$$P(\mathcal{D}|f) \stackrel{\text{iid}}{=} \prod_{i=1}^m \frac{1}{1 + \exp(-y_i f(x_i))}$$

Equivalent to minimizing log loss

$$-\log P(\mathcal{D}|f) = \sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

Boosting and Logistic Regression

Logistic regression equivalent to minimizing log loss

$$\sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

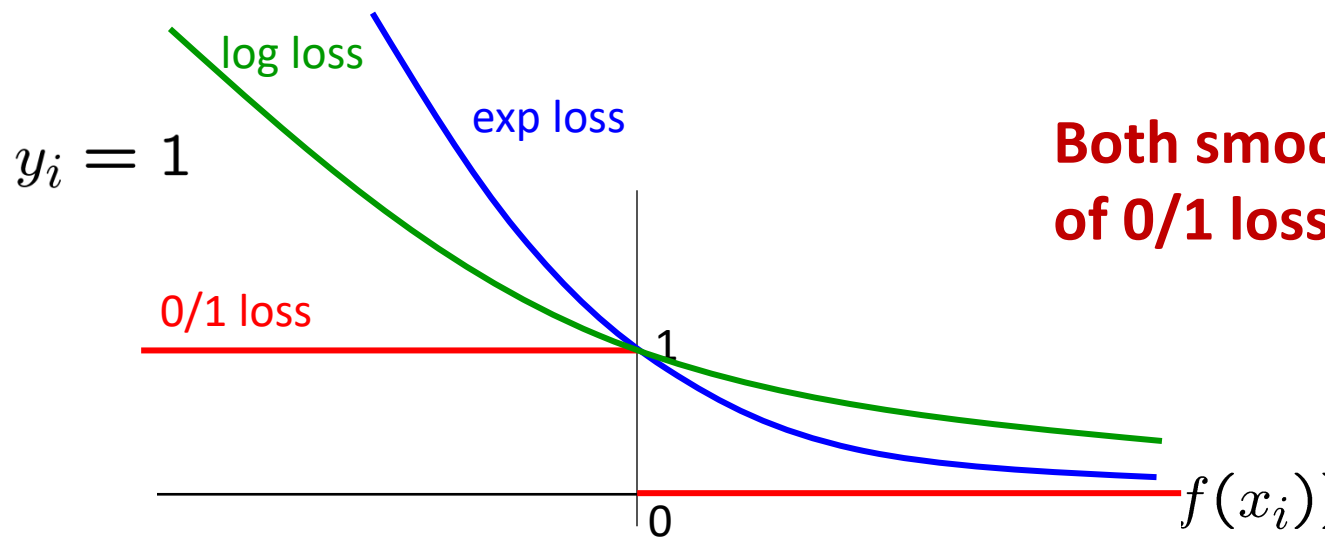
$$f(x) = w_0 + \sum_j w_j x_j$$

Boosting minimizes similar loss function!!

$$\frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i))$$

$$f(x) = \sum_t \alpha_t h_t(x)$$

Weighted average of weak learners



**Both smooth approximations
of 0/1 loss!**

Boosting and Logistic Regression

Logistic regression:

- Minimize log loss

$$\sum_{i=1}^m \ln(1 + \exp(-y_i f(x_i)))$$

- Define

$$f(x) = \sum_j w_j x_j$$

where x_j predefined features
(linear classifier)

- Jointly optimize over all weights $w_0, w_1, w_2 \dots$

Boosting:

- Minimize exp loss

$$\sum_{i=1}^m \exp(-y_i f(x_i))$$

- Define

$$f(x) = \sum_t \alpha_t h_t(x)$$

where $h_t(x)$ defined dynamically
to fit data
(not a linear classifier)

- Weights α_t learned per iteration
t incrementally

Hard & Soft Decision

Weighted average of weak learners $f(x) = \sum_t \alpha_t h_t(x)$

Hard Decision/Predicted label: $H(x) = \text{sign}(f(x))$

Soft Decision:
(based on analogy with
logistic regression)

$$P(Y = 1|X) = \frac{1}{1 + \exp(f(x))}$$

Bagging vs Boosting

Bagging

Resamples data points

Weight of each classifier
is the same

Only variance reduction

vs.

Boosting

Reweights data points (modifies their
distribution)

Weight is dependent on
classifier's accuracy

Both bias and variance reduced –
learning rule becomes more complex
with iterations

Thanks to all of you!

Thanks to Our Course Staff!

Teaching Assistants



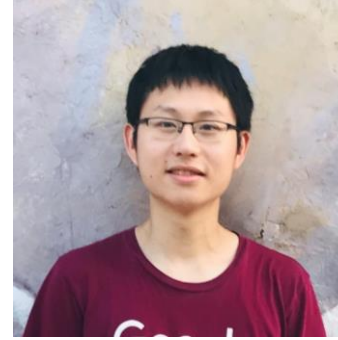
Alex
Singh



Annie
Hu



George
Brown



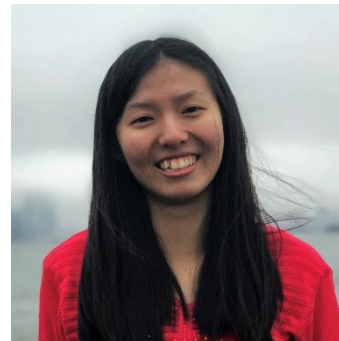
Haoran
Fei



Michell
Ma

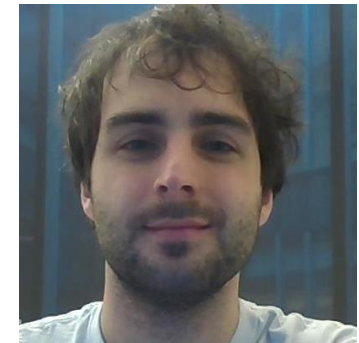


Nidhi
Jain



Vicky
Zeng

Education Associate



Daniel
Bird

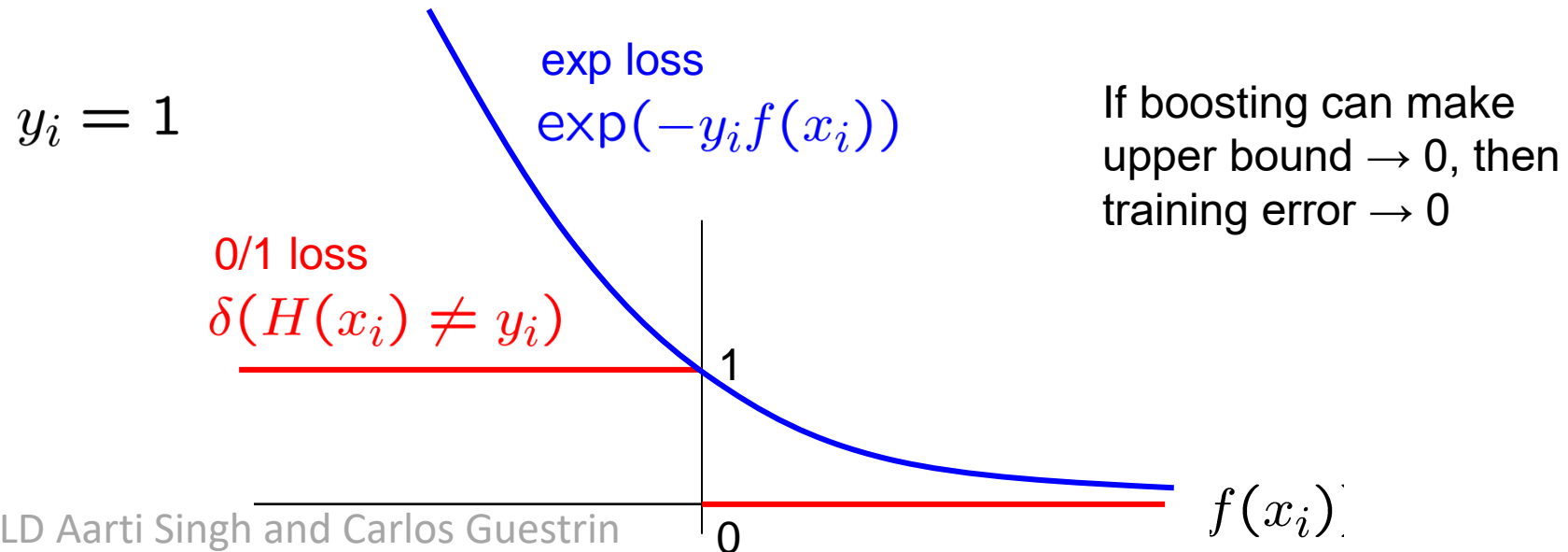
Additional Analysis

Analyzing training error

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \underbrace{\delta(H(x_i) \neq y_i)}_{\text{0/1 loss}} \leq \frac{1}{m} \sum_{i=1}^m \underbrace{\exp(-y_i f(x_i))}_{\text{exp loss}} \quad \text{Convex upper bound}$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$



Analyzing training error

Analysis reveals:

- What α_t to choose for hypothesis h_t ?
What h_t to choose?

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

ϵ_t - weighted training error

- If each weak learner h_t is slightly better than random guessing ($\epsilon_t < 0.5$), then training error of AdaBoost decays exponentially fast in number of rounds T .

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \exp \left(-2 \sum_{t=1}^T (1/2 - \epsilon_t)^2 \right)$$

Training Error

Analyzing training error

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_t Z_t$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$

Proof: Using Weight Update Rule

$$D_1(i) = 1/m.$$

$$D_2(i) = \frac{1}{m} \frac{e^{-\alpha_1 y_i h_1(x_i)}}{Z_1}$$

$$D_3(i) = \frac{1}{m} \frac{e^{-\alpha_1 y_i h_1(x_i)} e^{-\alpha_2 y_i h_2(x_i)}}{Z_1 Z_2}$$

$$D_{T+1}(i) = \frac{1}{m} \frac{\exp(-y_i f(x_i))}{\prod_t Z_t}$$

Wts of all pts add to 1

$$\sum_{i=1}^m D_{T+1}(i) = 1$$

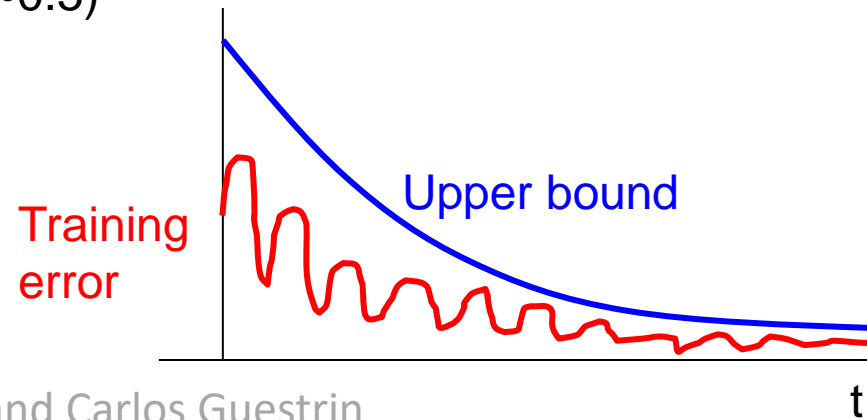
Analyzing training error

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_t Z_t$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$

If $Z_t < 1$, training error decreases exponentially (even though weak learners may not be good $\varepsilon_t \sim 0.5$)



What α_t to choose for hypothesis h_t ?

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i)) = \prod_t Z_t$$

Where $f(x) = \sum_t \alpha_t h_t(x)$; $H(x) = \text{sign}(f(x))$

If we minimize $\prod_t Z_t$, we minimize our training error

We can tighten this bound greedily, by choosing α_t and h_t on each iteration to minimize Z_t .

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

What α_t to choose for hypothesis h_t ?

We can minimize this bound by choosing α_t on each iteration to minimize Z_t .

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

For boolean target function, this is accomplished by [Freund & Schapire '97]:

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

Proof:

$$\begin{aligned} Z_t &= \sum_{i: y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t} + \sum_{i: y_i = h_t(x_i)} D_t(i) e^{-\alpha_t} \\ &= \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t} \end{aligned}$$

$$\frac{\partial Z_t}{\partial \alpha_t} = \epsilon_t e^{\alpha_t} - (1 - \epsilon_t) e^{-\alpha_t} = 0 \quad \Rightarrow \quad e^{2\alpha_t} = \frac{1 - \epsilon_t}{\epsilon_t}$$

Bounding the error with choice of α_t

We can minimize this bound by choosing α_t on each iteration to minimize Z_t .

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

For boolean target function, this is accomplished by [Freund & Schapire '97]:

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

Proof:

$$\begin{aligned} Z_t &= \sum_{i: y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t} + \sum_{i: y_i = h_t(x_i)} D_t(i) e^{-\alpha_t} \\ &= \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t} \\ &= 2\sqrt{\epsilon_t(1 - \epsilon_t)} = \sqrt{1 - (1 - 2\epsilon_t)^2} \end{aligned}$$

Dumb classifiers made Smart

Training error of final classifier is bounded by:

$$\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i) \leq \prod_t Z_t = \prod_t \sqrt{1 - (1 - 2\epsilon_t)^2}$$
$$\leq \exp \left(-2 \sum_{t=1}^T \underbrace{(1/2 - \epsilon_t)^2}_{\text{grows as } \epsilon_t \text{ moves away from } 1/2} \right)$$

If each classifier is (at least slightly) better than random $\epsilon_t < 0.5$

AdaBoost will achieve zero training error exponentially fast (in number of rounds T) !!

What about test error?

Generalization Error Bounds

[Freund & Schapire'95]

$$error_{true}(H) \leq error_{train}(H) + \tilde{O} \left(\sqrt{\frac{Td}{m}} \right) \quad \text{With high probability}$$

Boosting can overfit if T is large

Boosting often,

Contradicts experimental results

- Robust to overfitting
- Test set error decreases even after training error is zero

Need better analysis tools – margin based bounds

Margin Based Bounds

[Schapire, Freund, Bartlett, Lee'98]

$$error_{true}(H) \leq \hat{\Pr} [\text{margin}_f(x, y) \leq \theta] + \tilde{O} \left(\sqrt{\frac{d}{m\theta^2}} \right) \quad \text{With high probability}$$

Boosting increases the margin very aggressively since it concentrates on the hardest examples.

If margin is large, more weak learners agree and hence more rounds does not necessarily imply that final classifier is getting more complex.

Bound is independent of number of rounds T !

Boosting can still overfit if margin is too small (can perform arbitrarily close to random guessing) or weak learners are too complex