

Boosting

Can we make dumb learners smart?

Aarti Singh

Machine Learning 10-315

Mar 2, 2022

Slides Courtesy: Carlos Guestrin, Freund & Schapire



MACHINE LEARNING DEPARTMENT



Why boost weak learners?

Goal: Classify movie review sentiment

“I'm a fan of TV movies in general and this was one of the good ones”

“Long, boring. Never have I been so glad to see ending credits roll”

“I don't know why I like this movie, but I never get tired.”

- Easy to find “rules of thumb” that are better than random chance.

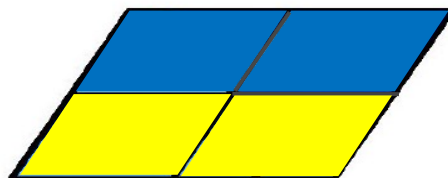
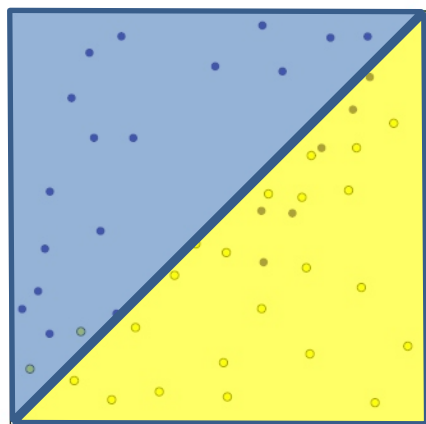
E.g. If ‘good’ occurs in utterance, then predict ‘positive’

- Hard to find single highly accurate prediction rule.

e.g. “This movie is terrible but it has some **good** effects”

Fighting the bias-variance tradeoff

- **Simple (a.k.a. weak) learners** e.g., naïve Bayes, logistic regression, decision stumps (or shallow decision trees)



weak learners
- high bias ✓
- low var

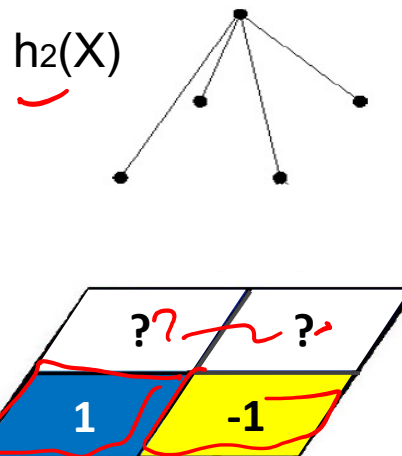
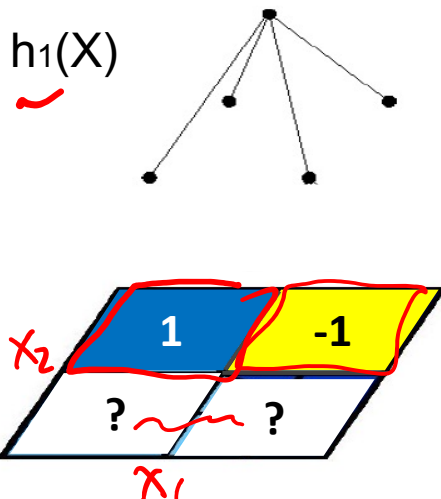
Are good 😊 - don't usually overfit

Are bad ☹️ - can't solve hard learning problems

- **Can we make weak learners good???**

Voting (Ensemble Methods)

- Instead of learning a single (weak) classifier, learn **many weak classifiers** that are **good at different parts of the input space**
- **Output class:** (Weighted) vote of each classifier
 - Classifiers that are most “sure” will vote with more conviction ✓ $\leftarrow \alpha_t$
 - Classifiers will be most “sure” about a particular part of the space ✓ $\leftarrow h_t$
 - On average, do better than single classifier!



$$H: X \rightarrow Y \text{ } (-1, 1)$$

$$H(X) = h_1(X) + h_2(X)$$

$$H(X) = \text{sign}\left(\sum_t \alpha_t h_t(X)\right)$$

final classifier
weights

Voting (Ensemble Methods)

- Instead of learning a single (weak) classifier, learn **many weak classifiers** that are **good at different parts of the input space**
- **Output class:** (Weighted) vote of each classifier
 - Classifiers that are most “sure” will vote with more conviction
 - Classifiers will be most “sure” about a particular part of the space
 - On average, do better than single classifier!
- **But how do you ???**
 - force classifiers h_t to learn about different parts of the input space?
 - weigh the votes of different classifiers? α_t

Boosting [Schapire'89]

- **Idea:** given a weak learner, run it multiple times on (reweighted) training data, then let learned classifiers vote
- On each iteration t :
 - weight $D_t(i)$ for each training example i , based on how incorrectly it was classified
 - Learn a weak hypothesis – h_t
 - A weight for this hypothesis – α_t
- Final classifier:
$$H(X) = \text{sign}(\sum \alpha_t h_t(X))$$
- **Practically useful**
- **Theoretically interesting**

*hypothesis = learner
= predictor*

Learning from weighted data

- Consider a weighted dataset
 - $D(i)$ – weight of i th training example (\mathbf{x}^i, y^i)
 - Interpretations:
 - i th training example counts as $D(i)$ examples
 - If I were to “resample” data, I would get more samples of “heavier” data points
- Now, in all calculations, whenever used, i th training example counts as $D(i)$ “examples”
 - e.g., in MLE redefine $Count(Y=y)$ to be weighted count

Unweighted data

$$Count(Y=y) = \sum_{i=1}^m \mathbf{1}(Y^i=y)$$

Weights $D(i)$

$$Count(Y=y) = \sum_{i=1}^m D(i) \mathbf{1}(Y^i=y)$$

AdaBoost [Freund & Schapire'95]

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

Initialize $D_1(i) = 1/m$. **Initially equal weights**

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

- Train weak learner using distribution D_t .

Naïve bayes, decision stump

- Get weak classifier $h_t : X \rightarrow \mathbb{R}$.

- Choose $\alpha_t \in \mathbb{R}$. *weight for classifier h_t*

- Update:

*weights
on data points*

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \left\{ \begin{array}{ll} e^{-\alpha_t} & \text{if } y_i = h_t(x_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(x_i) \end{array} \right\}$$

$$= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

**Increase weight
if wrong on pt i
 $y_i h_t(x_i) = -1 < 0$**

where Z_t is a normalization factor

↓
 $e^{-\alpha_t}$ y_i $h_t(x_i)$
 +1 +1
 -1 -1

AdaBoost [Freund & Schapire'95]

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

Initialize $D_1(i) = 1/m$. **Initially equal weights**

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

- Train weak learner using distribution D_t . **Naïve bayes, decision stump**

- Get weak classifier $h_t : X \rightarrow \mathbb{R}$.

- Choose $\alpha_t \in \mathbb{R}$. **Magic (+ve)**

- Update:

$$\underline{D_{t+1}(i)} = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

**Increase weight
if wrong on pt i
 $y_i h_t(x_i) = -1 < 0$**

where Z_t is a normalization factor

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

**Weights for all
pts must sum to 1
 $\sum_t \underline{D_{t+1}(i)} = 1$**

AdaBoost [Freund & Schapire'95]

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

Initialize $D_1(i) = 1/m$. **Initially equal weights**

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

- Train weak learner using distribution D_t . **Naïve bayes, decision stump**

- Get weak classifier $h_t : X \rightarrow \mathbb{R}$.

- Choose $\alpha_t \in \mathbb{R}$. **Magic (+ve)**

- Update:

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

**Increase weight
if wrong on pt i
 $y_i h_t(x_i) = -1 < 0$**


where Z_t is a normalization factor

Output the final classifier:

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


What α_t to choose for hypothesis h_t ?

Weight Update Rule:
$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$


$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$
 [Freund & Schapire'95]

Weighted training error

$$\epsilon_t = P_{i \sim D_t(i)} [h_t(x^i) \neq y^i] = \sum_{i=1}^m D_t(i) \underbrace{\delta(h_t(x_i) \neq y_i)}_{\substack{\text{Does } h_t \text{ get } i^{\text{th}} \text{ point wrong}}}$$

Handwritten red notes:
- Above the delta function: $1_{h_t(x_i) \neq y_i}$
- Under the sum term $D_t(i)$: 

$\epsilon_t = 0$ if h_t perfectly classifies all weighted data pts

$\epsilon_t = 1$ if h_t perfectly wrong \Rightarrow $-h_t$ perfectly right

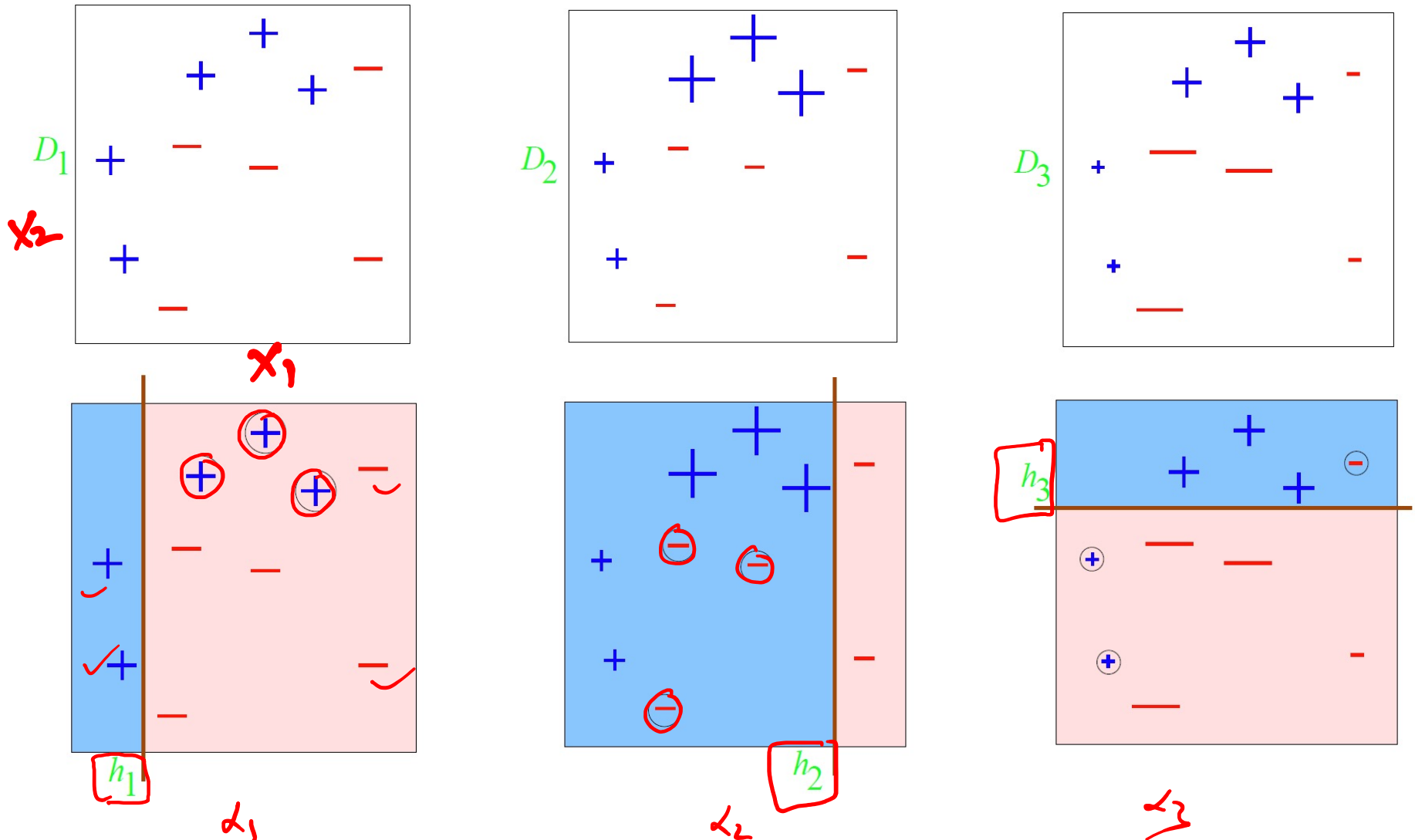
$\epsilon_t = 0.5$

$\alpha_t = \infty$

$\alpha_t = -\infty$

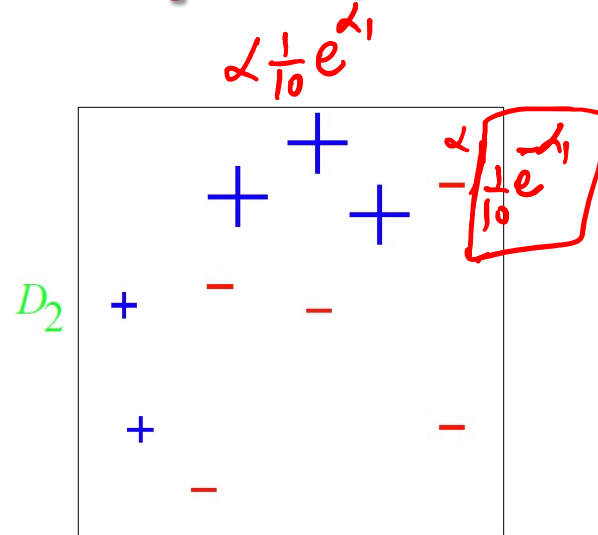
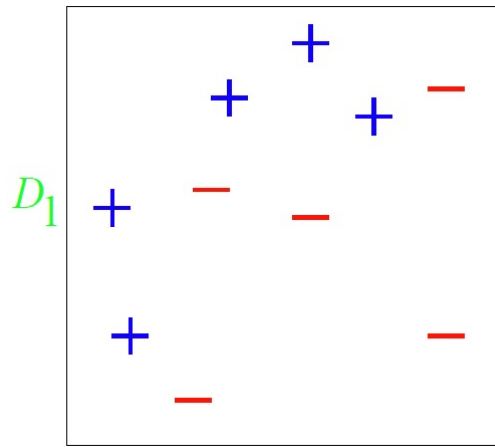
$\alpha_t = 0$

Boosting Example (Decision Stumps)



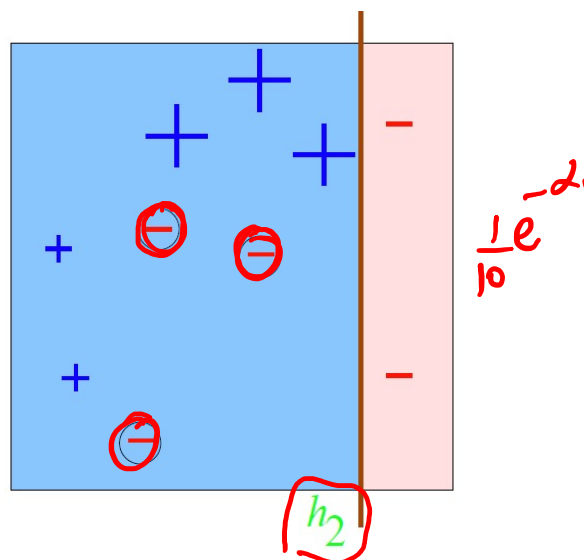
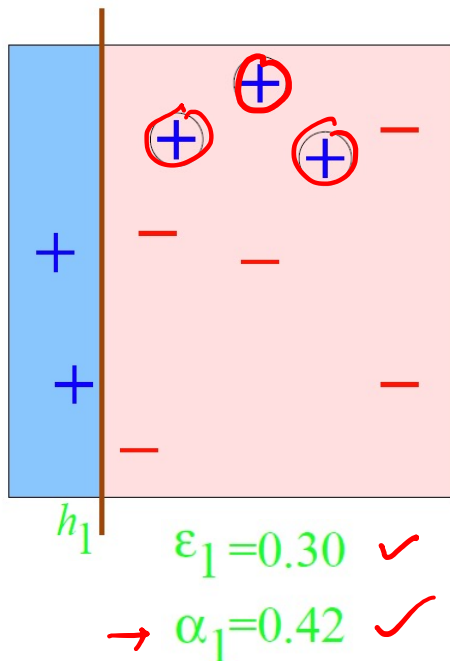
Boosting Example (Decision Stumps)

$$\mathcal{L} = \frac{1}{2} \ln \frac{e}{1-e} \quad \checkmark$$



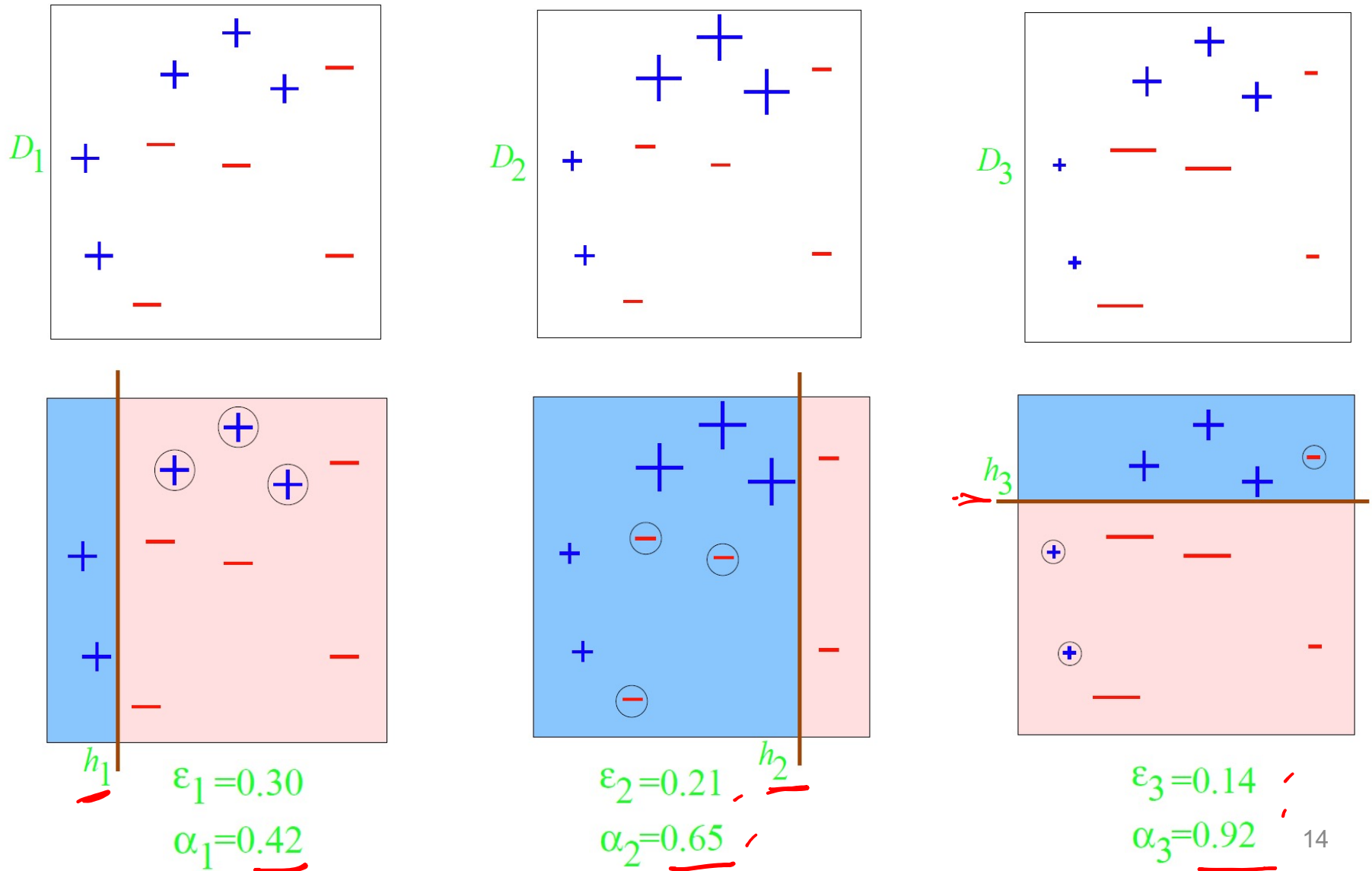
$$\epsilon_2 = \frac{3 \times \frac{1}{10} e^{-0.42}}{7 \times \frac{1}{10} e^{-0.42} + 3 \times \frac{1}{10} e^{0.42}}$$

\mathcal{L}_2

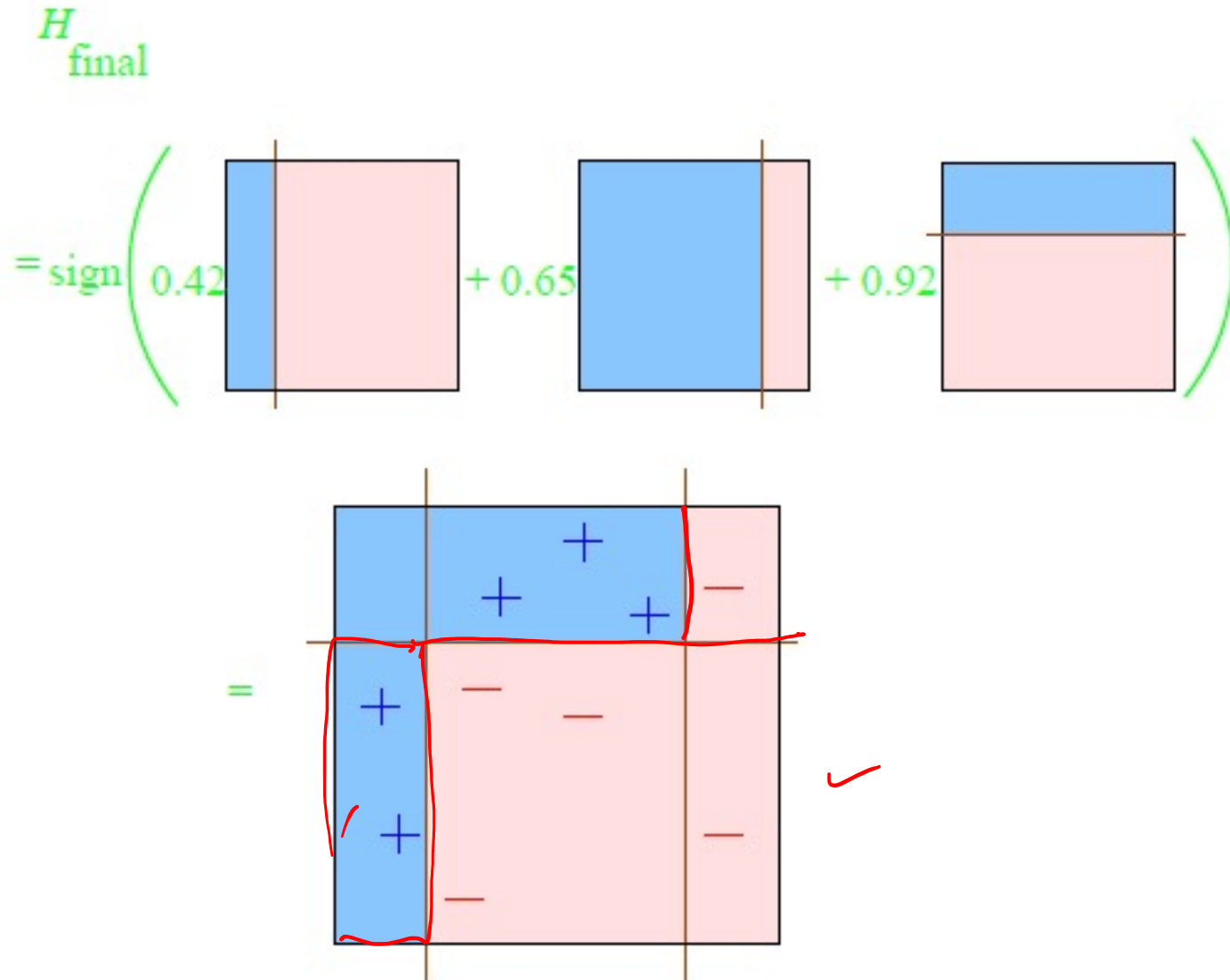


➤ What's the error on the weighted training data, ϵ_2 ?

Boosting Example (Decision Stumps)



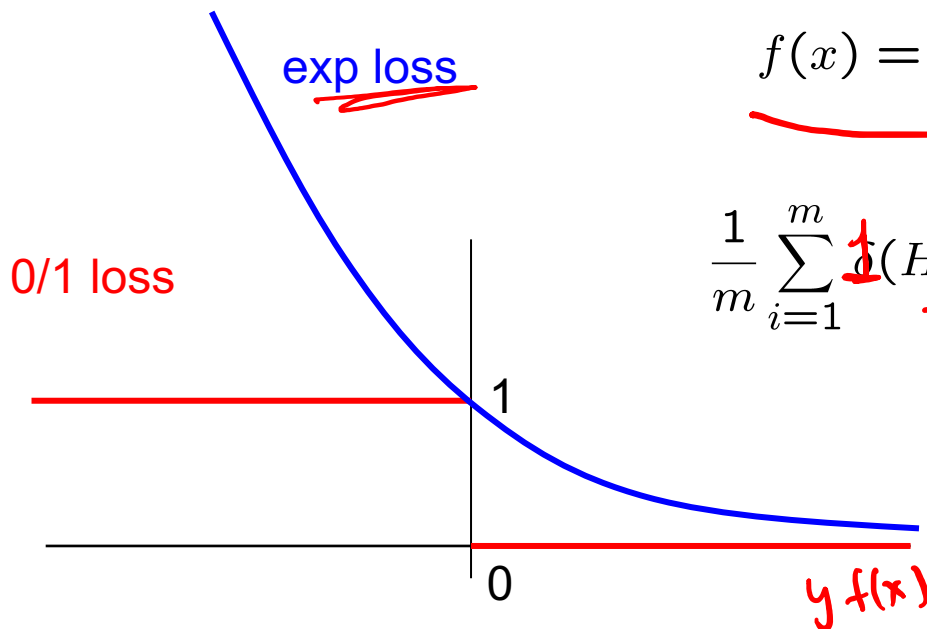
Boosting Example (Decision Stumps)



Analysis for Boosting

$$\mathcal{J}(\alpha_t, h_t)$$

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



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

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

0/1 loss

exp loss

Analysis for Boosting

Analysis reveals:

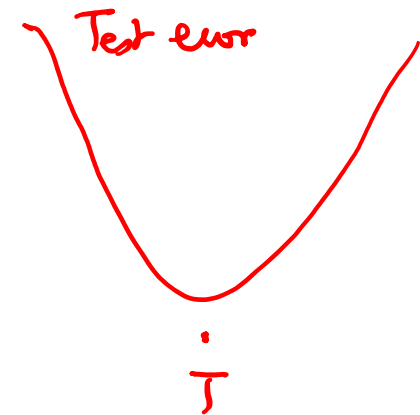
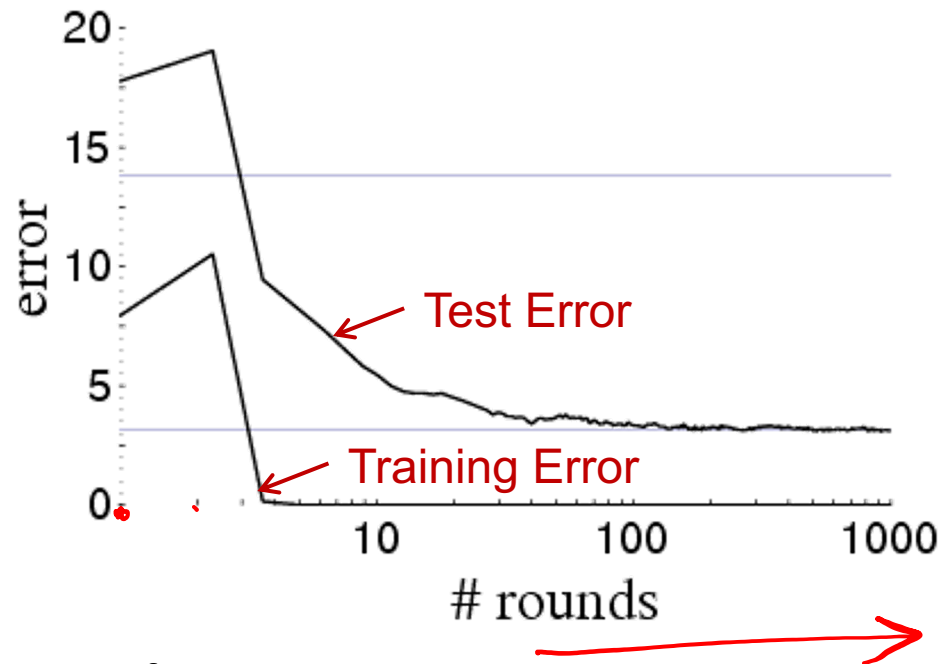
- If each weak learner h_t is slightly better than random guessing ($\epsilon_t < 0.5$), then training error of AdaBoost decays exponentially fast in number of rounds T .

$$\underbrace{\frac{1}{m} \sum_{i=1}^m \delta(H(x_i) \neq y_i)}_{\text{Training Error}} \leq \underbrace{\exp \left(-2 \sum_{t=1}^T (\underbrace{1/2}_{\text{red}} - \underbrace{\epsilon_t}_{\text{red}})^2 \right)}_{\text{red wavy line}}$$

What about test error?

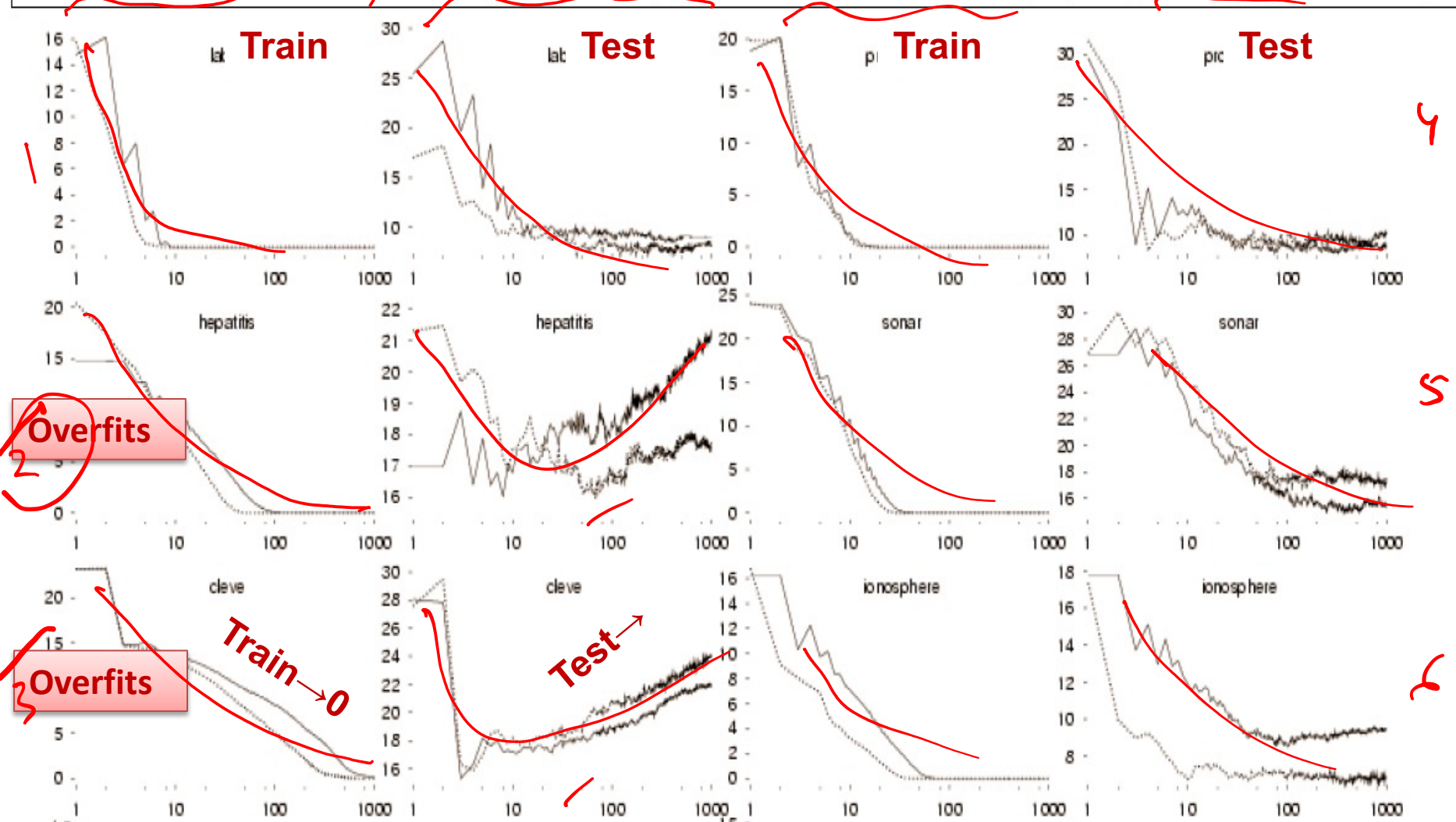
Boosting results – Digit recognition

[Schapire, 1989]



- Boosting often,
 - Robust to overfitting
 - Test set error decreases even after training error is zero
- If classes are well-separated, 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 classes not well separated (high label noise) or weak learners are too complex.

Handwritten red notes: $- - - + + +$ and $+ +$

Boosting and Logistic Regression

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

Logistic regression assumes:

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

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

And tries to maximize data likelihood: *conditioned*

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

Equivalent to minimizing log loss

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

Boosting and Logistic Regression

Logistic regression equivalent to minimizing log loss

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

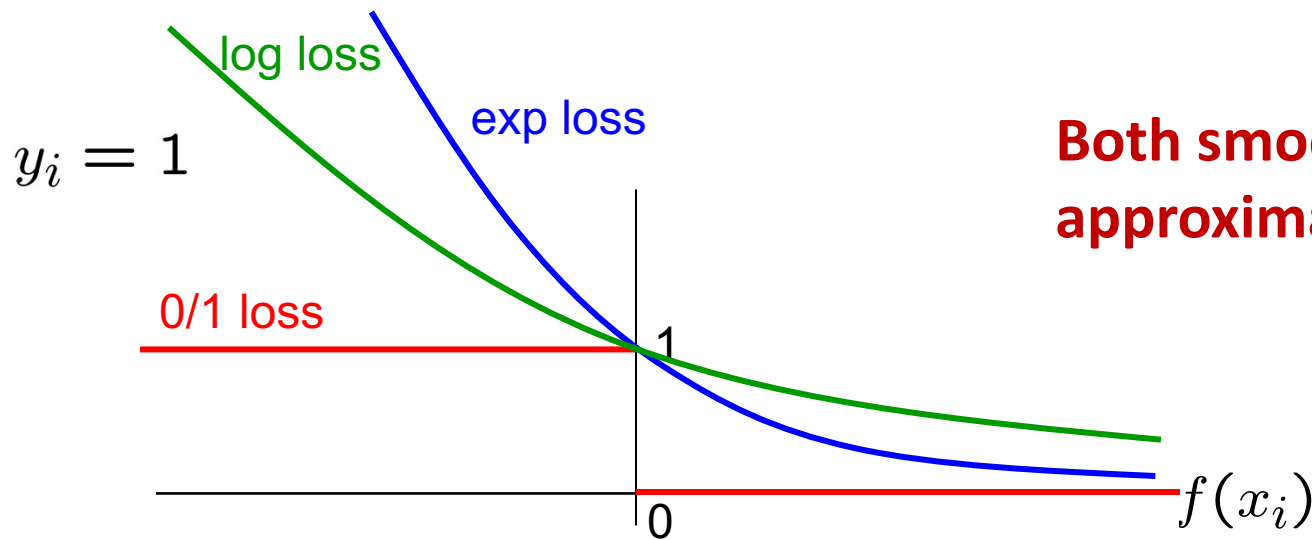
$$f(x) = w_0 + \sum_j w_j \underline{x_j}$$

Boosting minimizes similar loss function!!

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

$$f(x) = \sum_t \alpha_t \underline{h_t(x)}$$

Weighted average of weak learners



Both smooth and convex 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

$$f(x) = \sum_t \alpha_t \underline{h_t(x)}$$


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)

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


Matlab example – decision tree

load ionosphere

% UCI dataset

% 34 features, 351 samples

% binary classification

rng(100)

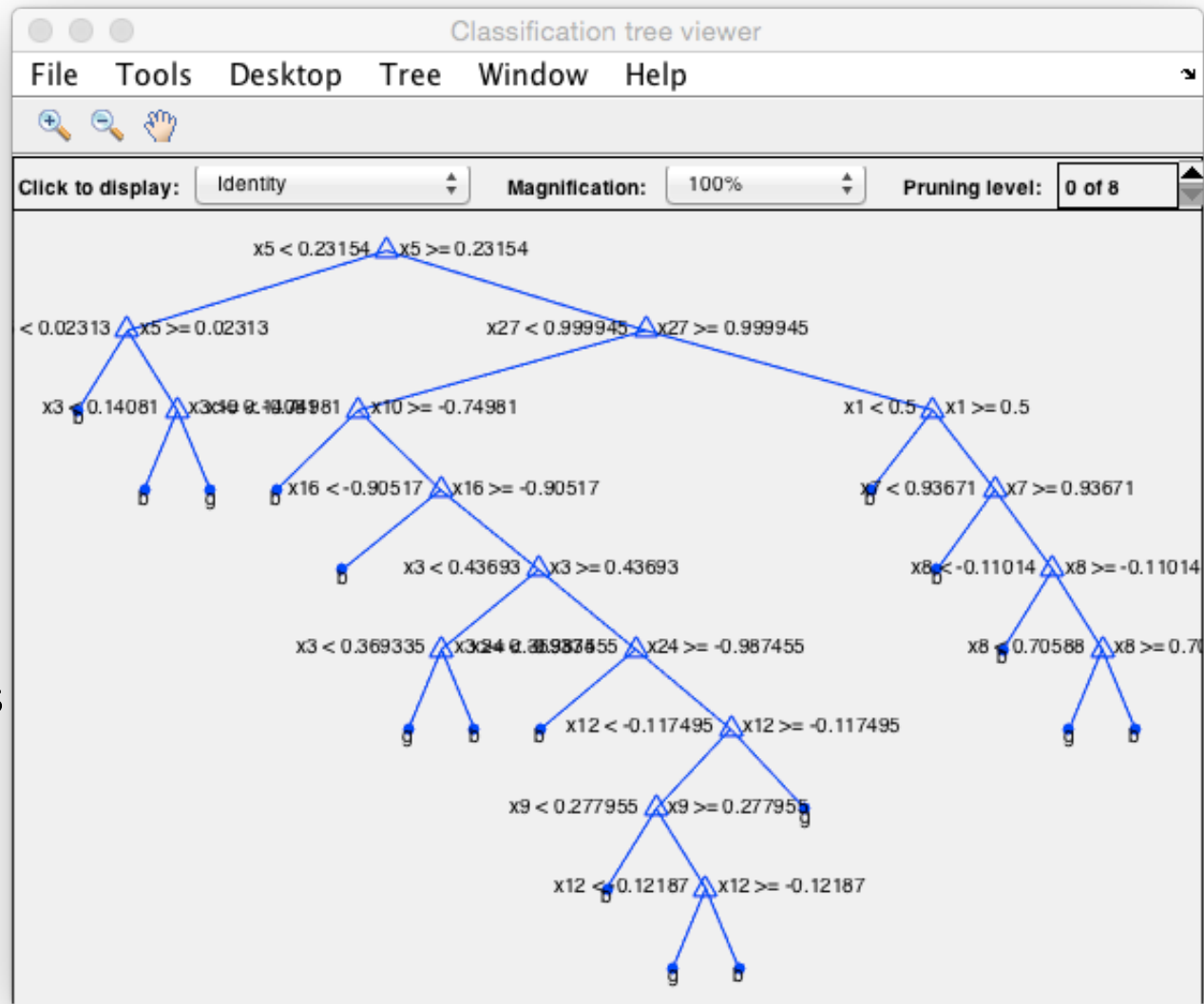
%Default MinLeafSize = 1

tc = fitctree(X,Y);

cvmodel = crossval(tc);

view(cvmodel.Trained{1},'Mode','graph')

kfoldLoss(cvmodel)



Validation error = 0.1254

Matlab example – decision tree

load ionosphere

% UCI dataset

% 34 features, 351 samples

% binary classification

rng(100)

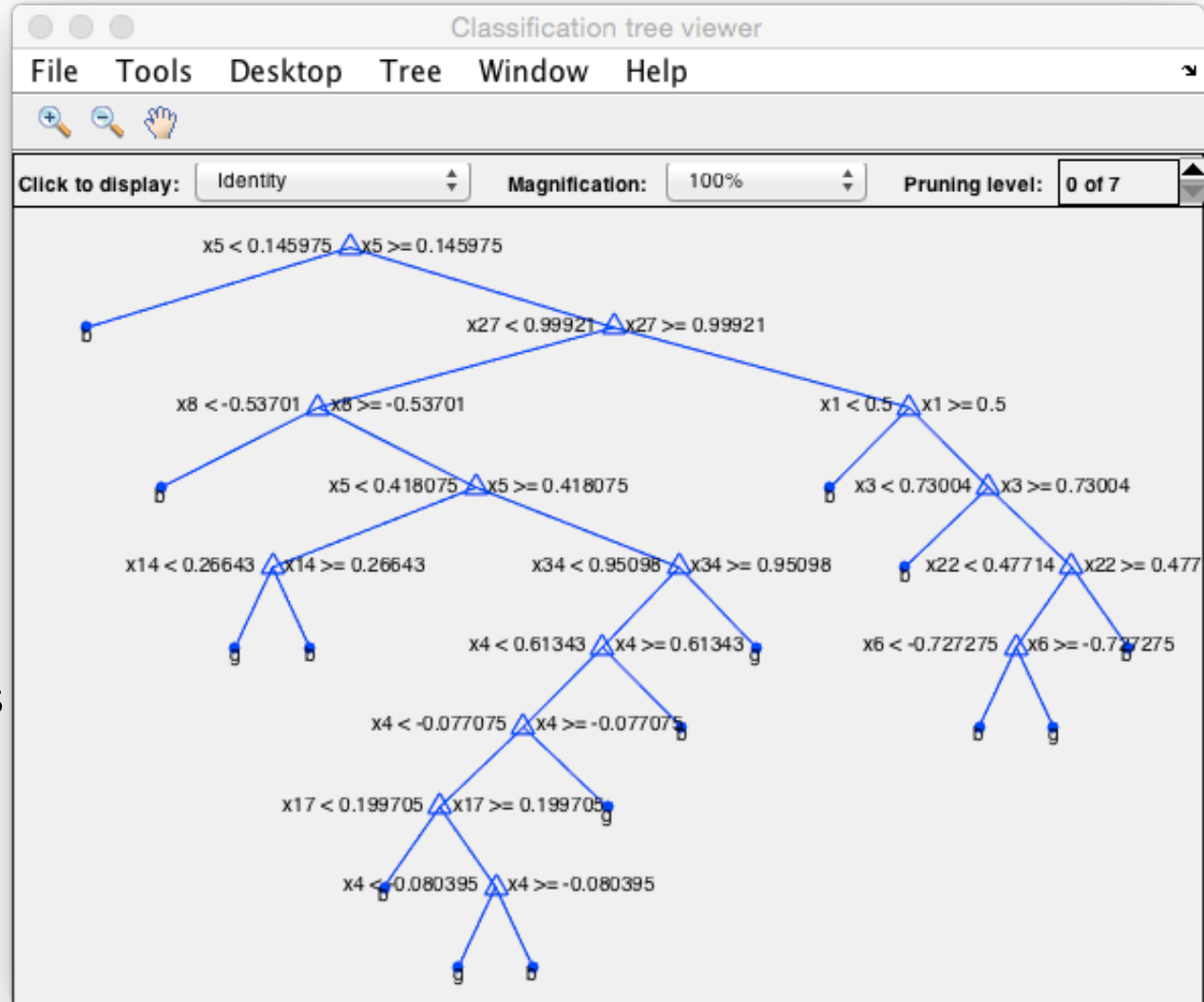
%Default MinLeafSize = 1

tc = fitctree(X,Y, 'MinLeafSize',2);

cvmodel = crossval(tc);

view(cvmodel.Trained{1},'Mode','graph')

kfoldLoss(cvmodel)

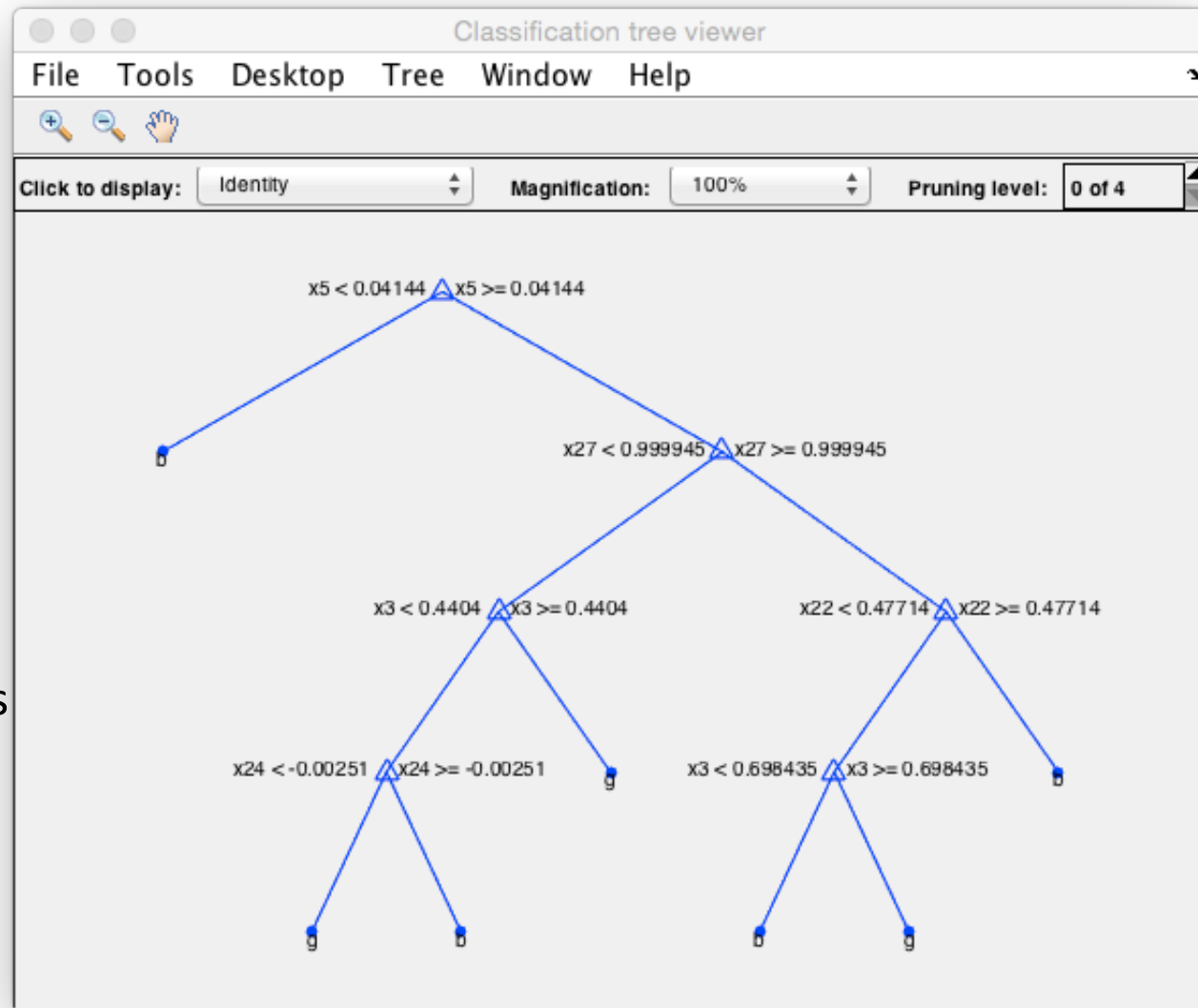


Validation error = 0.1168

Matlab example – decision tree

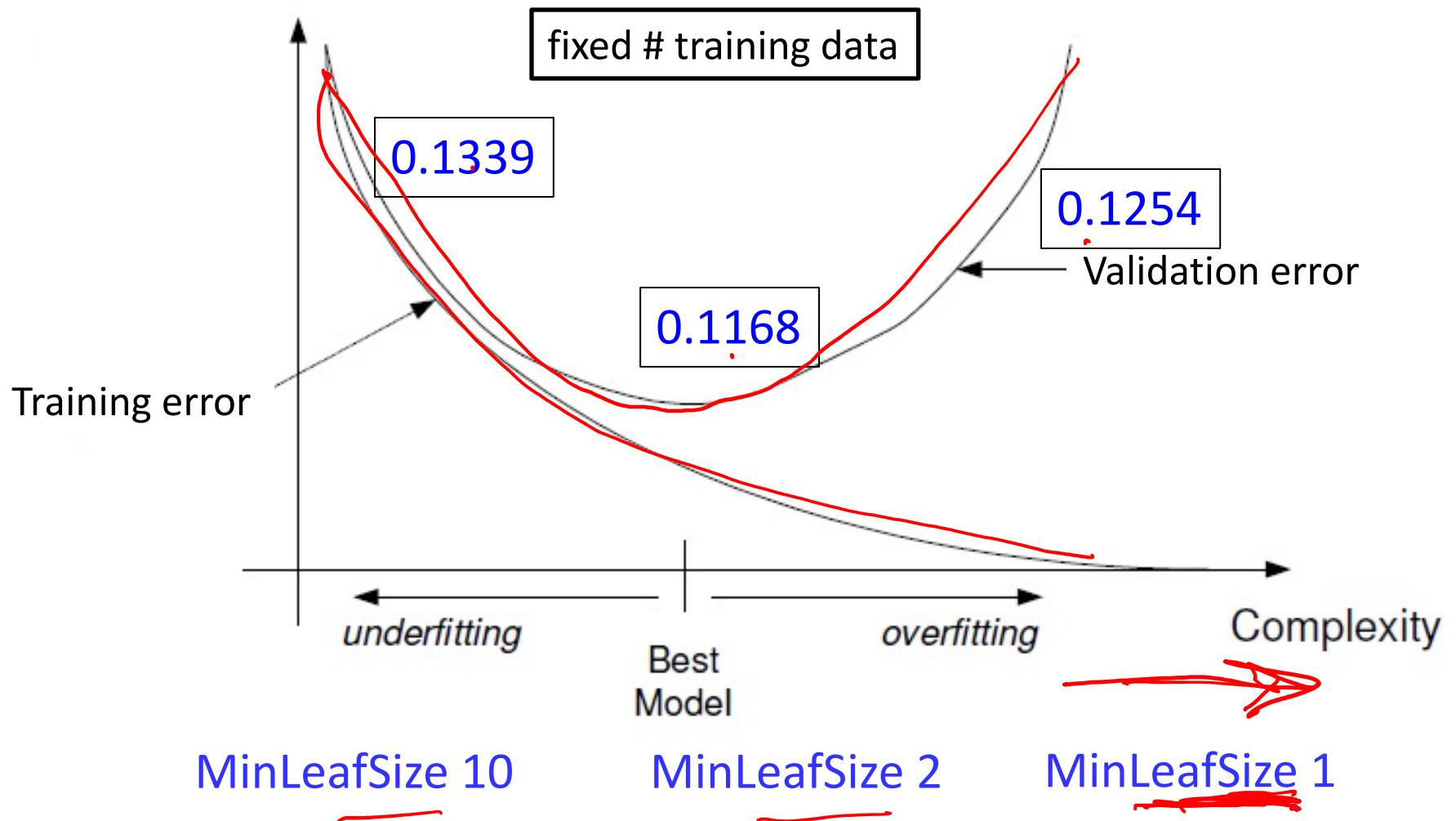
```
load ionosphere
% UCI dataset
% 34 features, 351 samples
% binary classification
rng(100)
```

```
%Default MinLeafSize = 1
tc = fitctree(X,Y, 'MinLeafSize',10);
cvmodel = crossval(tc);
view(cvmodel.Trained{1},'Mode','graph')
kfoldLoss(cvmodel)
```



Validation error = 0.1339

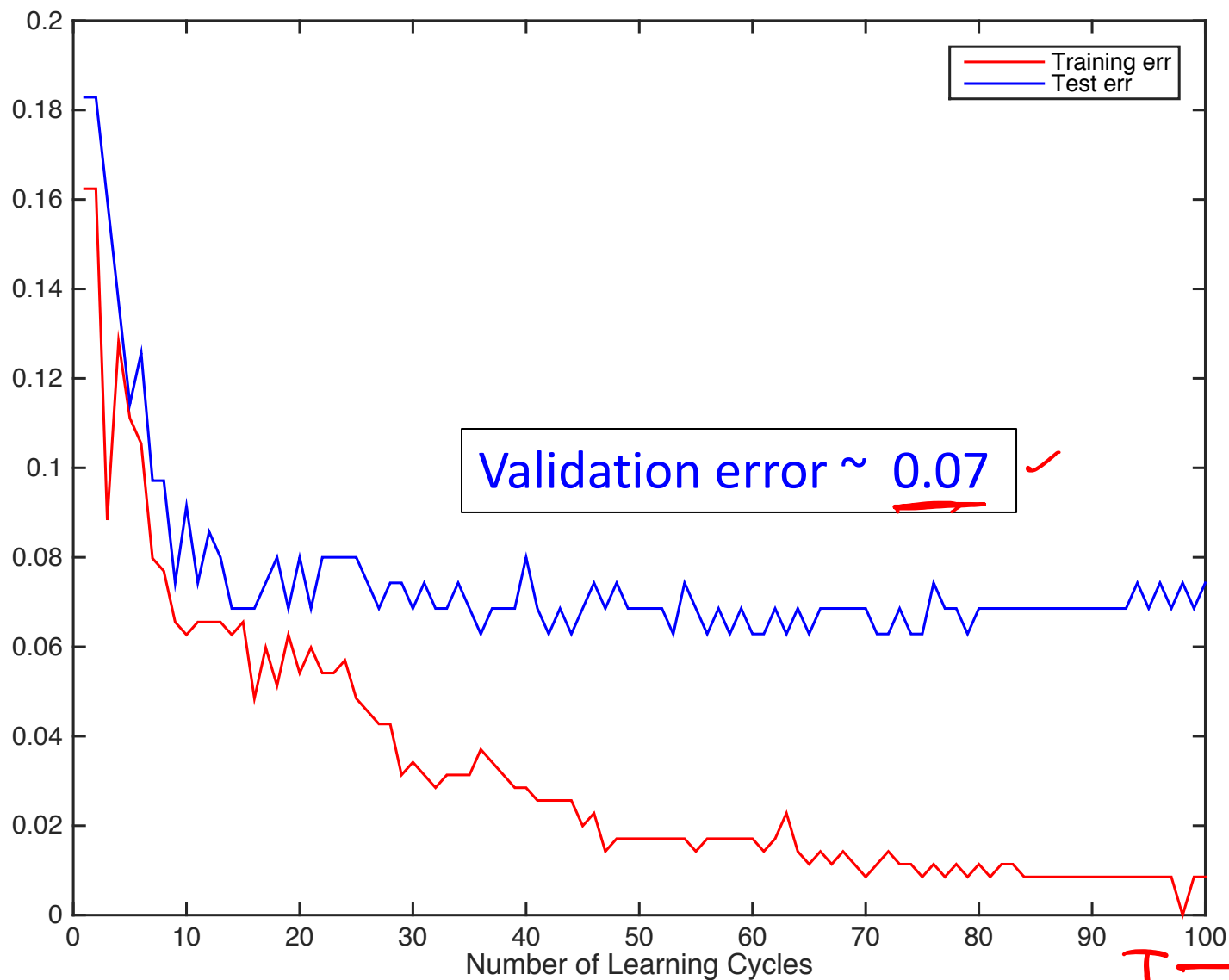
Matlab example – decision trees



Matlab example - boosting

- % UCI dataset
- % 34 features, 351 samples
- % binary classification
- load ionosphere;
- rng(2); % For reproducibility
- ClassTreeEns = fitensemble(X,Y,'AdaBoostM1',100,'Tree');
- rsLoss = resubLoss(ClassTreeEns,'Mode','Cumulative');
- plot(rsLoss,'r');
- hold on
- ClassTreeEns = fitensemble(X,Y,'AdaBoostM1',100,'Tree',...
- 'Holdout',0.5);
- genError = kfoldLoss(ClassTreeEns,'Mode','Cumulative');
- plot(genError,'b');
- xlabel('Number of Learning Cycles');
- legend('Training err', 'Test err')

Matlab example - boosting



T →

Bagging

[Breiman, 1996]

Related approach to combining classifiers:

1. Run independent weak learners on subsampled data (sample with replacement) from the training set
2. Average/vote over weak hypotheses

Bagging

vs.

Boosting

[Resamples data points

[Weight of each classifier is the same

[Only variance reduction ✓

[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

Boosting Summary

- Combine weak classifiers to obtain strong classifier
 - Weak classifier – slightly better than random on training data
 - Resulting very strong classifier – can eventually provide zero training error
- AdaBoost algorithm
- Boosting v. Logistic Regression
 - Similar loss functions
 - Single optimization (LR) v. Incrementally improving classification (B)
- Most popular application of Boosting:
 - Boosted decision stumps!
 - Very simple to implement, very effective classifier