# Decision Trees, cont.

# Boosting

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

October 1st, 2007

1

---

# A Decision Stump

mpg values:  bad   good

root

22   18

pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0   0 | 4   17 | 1   0 | 8   0 | 9   1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

17 > 4

2

mpg values:  bad  good

root

22  18

pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0  0 | 4  17 | 1  0 | 8  0 | 9  1 |
| Predict bad | pchance = 0.135 | Predict bad | Predict bad | pchance = 0.085 |

| maker = america | maker = asia | maker = europe | horsepower = low | horsepower = medium | horsepower = high |
|---|---|---|---|---|---|
| 0  10 | 2  5 | 2  2 | 0  0 | 0  1 | 9  0 |
| Predict good | pchance = 0.317 | pchance = 0.717 | Predict bad | Predict good | Predict bad |

| horsepower = low | horsepower = medium | horsepower = high | acceleration = low | acceleration = medium | acceleration = high |
|---|---|---|---|---|---|
| 0  4 | 2  1 | 0  0 | 1  0 | 0  1 | 1  1 |
| Predict good | pchance = 0.894 | Predict bad | Predict bad | Predict good | pchance = 0.717 |

| acceleration = low | acceleration = medium | acceleration = high | modelyear = 70to74 | modelyear = 75to78 | modelyear = 79to83 |
|---|---|---|---|---|---|
| 1  0 | 1  1 | 0  0 | 0  1 | 1  0 | 0  0 |
| Predict bad | (unexpandable) | Predict bad | Predict good | Predict bad | Predict bad |

Predict bad

*no more distinctions*

3

©Carlos Guestrin 2005-2007

# Basic Decision Tree Building Summarized

BuildTree(*DataSet,Output*)

- If all output values are the same in *DataSet*, return a leaf node that says "predict this unique output"
- If all input values are the same, return a leaf node that says "predict the majority output"
- Else find attribute *X* with highest Info Gain
- Suppose *X* has $n_X$ distinct values (i.e. X has arity $n_X$).
  - Create and return a non-leaf node with $n_X$ children.
  - The *i*'th child should be built by calling
        BuildTree(*DS_i,Output*)
     Where *DS_i* built consists of all those records in DataSet for which X = *i*th distinct value of X.

©Carlos Guestrin 2005-2007

4

mpg values: bad good

root
22 18
pchance = 0.001

| | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 1 | 40 | 2.50 |
| Test Set | 74 | 352 | 21.02 |

horsepower = high
0

Predict good   pchance = 0.311   pchance = 0.717   Predict bad   Predict good   Predict bad

| horsepower = low | horsepower = medium | horsepower = high | acceleration = low | acceleration = medium | acceleration = high |
|---|---|---|---|---|---|
| 0  4 | 2  1 | 0  0 | 1  0 | 0  1 | 1  1 |
| Predict good | pchance = 0.894 | Predict bad | Predict bad | Predict good | pchance = 0.717 |

| acceleration = low | acceleration = medium | acceleration = high | modelyear = 70to74 | modelyear = 75to78 | modelyear = 79to83 |
|---|---|---|---|---|---|
| 1  0 | 1  1 | 0  0 | 0  1 | 1  0 | 0  0 |
| Predict bad | (unexpandable) | Predict bad | Predict good | Predict bad | Predict bad |
| | Predict bad | | | | |

5

---

mpg values: bad good

root
22 18
pchance = 0.001

| | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 1 | 40 | 2.50 |
| Test Set | 74 | 352 | 21.02 |

horsepower = high
0

Predict good   pchance = 0.311   pchance = 0.717   Predict bad   Predict good   Predict bad

| horsepower = low | horsepower = medium | horsepower = high | acceleration = low | acceleration = medium | acceleration = high |
|---|---|---|---|---|---|
| 0 | | 0  0 | 1  0 | 0  1 | 1  1 |

The test set error is much worse than the training set error…

…why?

= 0.717

= 79to83

Predict bad   (unexpandable)   Predict bad   Predict good   Predict bad   Predict bad

Predict bad

6

# Decision trees & Learning Bias

| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|---|---|---|---|---|---|---|---|
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | low | medium | low | medium | 70to74 | asia |
| bad | 4 | low | medium | low | low | 70to74 | asia |
| bad | 8 | high | high | high | low | 75to78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 8 | high | medium | high | high | 79to83 | america |
| bad | 8 | high | high | high | low | 75to78 | america |
| good | 4 | low | low | low | low | 79to83 | america |
| bad | 6 | medium | medium | medium | high | 75to78 | america |
| good | 4 | medium | low | low | low | 79to83 | america |
| good | 4 | low | low | medium | high | 79to83 | america |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 4 | low | medium | low | medium | 75to78 | europe |
| bad | 5 | medium | medium | medium | medium | 75to78 | europe |

*(handwritten annotations)*

Can fit training set perfectly

i.e., $\ell_{train} = 0$

very complex tree

no label noise!

if $X^1 \& X^2$

$\therefore X^1 = X^2 \Rightarrow Y^1 = Y^2$

agree on features

agree on labels

---

# Decision trees will overfit

- Standard decision trees are have no learning biased
  - Training set error is always zero!
    - (If there is no label noise)
  - Lots of variance
  - Will definitely overfit!!!
  - Must bias towards simpler trees
- Many strategies for picking simpler trees:
  - Fixed depth
  - Fixed number of leaves
  - Or something smarter…

mpg values: bad good

Consider this split

©Carlos Guestrin 2005-2007

9

# A chi-square test

mpg values: bad good

| maker | | bad | good | | |
|---|---|---|---|---|---|
| | america | 0 | 10 | | H( mpg \| maker = america ) = 0 |
| | asia | 2 | 5 | | H( mpg \| maker = asia ) = 0.863121 |
| | europe | 2 | 2 | | H( mpg \| maker = europe ) = 1 |

H(mpg) = 0.702467   H(mpg|maker) = 0.478183

IG(mpg|maker) = 0.224284

- Suppose that mpg was completely uncorrelated with maker.
- What is the chance we'd have seen data of at least this apparent level of association anyway?

©Carlos Guestrin 2005-2007

10

# A chi-square test

mpg values:  bad  good

| maker | america | 0 | 10 | | H( mpg \| maker = america ) = 0 |
|---|---|---|---|---|---|
| | asia | 2 | 5 | | H( mpg \| maker = asia ) = 0.863121 |
| | europe | 2 | 2 | | H( mpg \| maker = europe ) = 1 |

H(mpg) = 0.702467   H(mpg\|maker) = 0.478183

IG(mpg\|maker) = 0.224284

- Suppose that mpg was completely uncorrelated with maker.
- What is the chance we'd have seen data of at least this apparent level of association anyway?

By using a particular kind of chi-square test, the answer is 7.2%

(Such simple hypothesis tests are very easy to compute, unfortunately, not enough time to cover in the lecture,

but in your homework, you'll have fun! :))

---

# Using Chi-squared to avoid overfitting

- Build the full decision tree as before
- But when you can grow it no more, start to prune:
  - Beginning at the bottom of the tree, delete splits in which $p_{chance}$ > *MaxPchance*
  - Continue working you way up until there are no more prunable nodes

*MaxPchance*  is a magic parameter you must specify to the decision tree, indicating your willingness to risk fitting noise

# Pruning example

- With MaxPchance = 0.1, you will see the following MPG decision tree:

mpg values: bad good

root

22  18

pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0  0 | 4  17 | 1  0 | 8  0 | 9  1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

*decision stump*

Note the improved test set accuracy compared with the unpruned tree

|  | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 5 | 40 | 12.50 |
| Test Set | 56 | 352 | 15.91 |

---

# MaxPchance

- Technical note MaxPchance is a regularization parameter that helps us bias towards simpler models



Expected Test set Error

*error test*

Decreasing    MaxPchance    Increasing

High Bias    *magic value*    High Variance

We'll learn to choose the value of these magic parameters soon!

# Real-Valued inputs

- What should we do if some of the inputs are real-valued?

| mpg | cylinders | displacemen | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|-------------|------------|--------|--------------|-----------|---------|
| good | 4 | 97 | 75 | 2265 | 18.2 | 77 | asia |
| bad | 6 | 199 | 90 | 2648 | 15 | 70 | america |
| bad | 4 | 121 | 110 | 2600 | 12.8 | 77 | europe |
| bad | 8 | 350 | 175 | 4100 | 13 | 73 | america |
| bad | 6 | 198 | 95 | 3102 | 16.5 | 74 | america |
| bad | 4 | 108 | 94 | 2379 | 16.5 | 73 | asia |
| bad | 4 | 113 | 95 | 2228 | 14 | 71 | asia |
| bad | 8 | 302 | 139 | 3570 | 12.8 | 78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| good | 4 | 120 | 79 | 2625 | 18.6 | 82 | america |
| bad | 8 | 455 | 225 | 4425 | 10 | 70 | america |
| good | 4 | 107 | 86 | 2464 | 15.5 | 76 | europe |
| bad | 5 | 131 | 103 | 2830 | 15.9 | 78 | europe |
| | | | | | | | |

Infinite number of possible split values!!!

Finite dataset, only finite number of relevant splits!

Idea One: Branch on each possible real value

15

# "One branch for each numeric value" idea:



Hopeless: with such high branching factor will shatter the dataset and overfit

16

# Threshold splits

- Binary tree, split on attribute X
  - One branch: X < t
  - Other branch: X ≥ t

*(handwritten diagram: Year node splitting into <79 and ≥79; left branch to accel. splitting into <18 and ≥18; right branch to Year splitting into <82 and ≥82)*

---

# Choosing threshold split

- Binary tree, split on attribute X
  - One branch: X < t
  - Other branch: X ≥ t
- Search through possible values of *t*
  - Seems hard!!!
- But only finite number of *t*'s are important
  - Sort data according to X into $\{x_1, \ldots, x_m\}$
  - Consider split points of the form $x_i + (x_{i+1} - x_i)/2$

# A better idea: thresholded splits

- Suppose X is real valued

  *[handwritten: ← threshold]*

- Define *IG(Y|X:t)* as *H(Y) - H(Y|X:t)*
- Define *H(Y|X:t)* =

  $$H(Y|X < t)\, P(X < t) + H(Y|X >= t)\, P(X >= t)$$

  *[handwritten above: year < t year year]*

  - *IG(Y|X:t)* is the information gain for predicting Y if all you know is whether X is greater than or less than *t*

- Then define *IG\*(Y|X) = max$_t$ IG(Y|X:t)*
- For each real-valued attribute, use *IG\*(Y|X)* for assessing its suitability as a split

  *[handwritten: → then pick best $X_i = \arg\max_i IG^*(Y|X_i)$]*

- Note, may split on an attribute multiple times, with different thresholds

*[handwritten right margin: naive implementation $O(m^2)$ for m data points. Dynamic program $O(m)$]*

19

---

# Example with MPG

Information gains using the training set (40 records)

mpg values:  bad  good

| Input | Value | Distribution | Info Gain |
|-------|-------|--------------|-----------|
| cylinders | < 5 | | 0.48268 |
| | >= 5 | | |
| displacement | < 198 | | 0.428205 |
| | >= 198 | | |
| horsepower | < 94 | | 0.48268 |
| | >= 94 | | |
| weight | < 2789 | | 0.379471 |
| | >= 2789 | | |
| acceleration | < 18.2 | | 0.159982 |
| | >= 18.2 | | |
| modelyear | < 81 | | 0.319193 |
| | >= 81 | | |
| maker | america | | 0.0437265 |
| | asia | | |
| | europe | | |

*[handwritten: best values]*

*[handwritten: Cyl < 5 , Cyl >= 5]*

20

# Example tree using reals



mpg values: bad good

root
22 18
pchance = 0.000

cylinders < 5 | cylinders >= 5
4 17 | 18 1
pchance = 0.001 | pchance = 0.003

horsepower < 94 | horsepower >= 94 | acceleration < 19 | acceleration >= 19
1 17 | 3 0 | 18 0 | 0 1
pchance = 0.274 | Predict bad | Predict bad | Predict good

maker = america | maker = asia | maker = europe
0 10 | 0 5 | 1 2
Predict good | Predict good | pchance = 0.270

displacement < 116 | displacement >= 116
0 2 | 1 0
Predict good | Predict bad

---

# What you need to know about decision trees

- Decision trees are one of the most popular data mining tools
  - □ Easy to understand
  - □ Easy to implement          " interpretable "
  - □ Easy to use
  - □ Computationally cheap (to solve heuristically)
- Information gain to select attributes (ID3, C4.5,…)
- Presented for classification, can be used for regression and density estimation too
- Decision trees will overfit!!!
  - □ Zero bias classifier → Lots of variance
  - □ Must use tricks to find "simple trees", e.g.,
    - Fixed depth/Early stopping
    - Pruning
    - Hypothesis testing

# Acknowledgements

■ Some of the material in the decision trees presentation is courtesy of Andrew Moore, from his excellent collection of ML tutorials:

  ☐ http://www.cs.cmu.edu/~awm/tutorials

# Announcements

■ Homework 1 due Wednesday beginning of class
  ☐ started early, started early, started early, started early, started early, started early, started early, started early

■ Exam dates set:    *this room*
  ☐ Midterm: Thursday, Oct. 25th, 5-6:30pm, MM A14
  ☐ Final: Tuesday, Dec. 11, 05:30PM-08:30PM    *↳ Somewhere*

# Fighting the bias-variance tradeoff

- **Simple (a.k.a. weak) learners are good**
  - ☐ e.g., naïve Bayes, logistic regression, decision stumps (or shallow decision trees)
  - ☐ Low variance, don't usually overfit
- **Simple (a.k.a. weak) learners are bad**
  - ☐ High bias, can't solve hard learning problems

- Can we make weak learners always good???
  - ☐ **No!!!**
  - ☐ **But often yes…**

---

# 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!

$$H(x) : X \mapsto Y$$

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

weight

$$Y \in \{-1, +1\}$$

Simple learners:
$$h_t(x) : X \longrightarrow \{-1, +1\}$$
$$\longrightarrow [-1, +1]$$

- **But how do you ???**
  - ☐ force classifiers to learn about different parts of the input space?
  - ☐ weigh the votes of different classifiers?

# Boosting [Schapire, 1989]

- Idea: given a weak learner, run it multiple times on (reweighted) training data, then let learned classifiers vote

- On each iteration $t$:
  - weight each training example by how incorrectly it was classified
  - Learn a hypothesis – $h_t$
  - A strength for this hypothesis – $\alpha_t$

- Final classifier:

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

*first iteration*
*record it.*
*more important*

- **Practically useful**
- **Theoretically interesting**

27

---

# Learning from weighted data

- **Sometimes not all data points are equal**
  - Some data points are more equal than others
- **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., MLE for Naïve Bayes, redefine *Count(Y=y)* to be weighted count

*Prior*
$$\hat{P}(Y=y) = \frac{Count(Y=y)}{m}$$

*weighted:*   *indicator*
$$\hat{P}(Y=y) = \frac{\sum_{i=1}^{m} D(i)\, \mathbb{1}(y^i = y)}{\sum_{i=1} D(i)}$$

28

## Slide 1

*AdaBoost* (handwritten)

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

$\left( \frac{e^{-a} < 1 \text{ if } a > 0}{e^{a} > 1 \text{ if } a < 0} \right)$ (handwritten)

Initialize $D_1(i) = 1/m$. ⟵ uniform

For $t = 1, \ldots, T$: ⟵ iteration

- Train base learner using distribution $D_t$. (weak) ← as in previous slide
- Get base classifier $h_t : X \to \mathbb{R}$.
- Choose $\alpha_t \in \mathbb{R}$.
- Update:    *weight before* (handwritten)

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

*if normalized* $Z_t = 1$ (handwritten)

where $Z_t$ is a normalization factor    *normalizer* $\sum_i D_{t+1}(i) = 1$

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

Output the final classifier:

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

Handwritten right margin:
if $y_i = +1$
$\alpha_t > 0$
if $h_t(x_i)$ is correct:
$\Rightarrow h_t(x_i) > 0$
$\Rightarrow -\alpha_t y_i h_t(x_i) < 0$
$\Rightarrow D_{t+1}(i)$ reduced
if $h_t(x_i)$ is incorrect
$\Rightarrow -\alpha_t y_i h_t(x_i) > 0$
$\Rightarrow D_{t+1}(i)$ increase

Figure 1: The boosting algorithm AdaBoost.

## Slide 2

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

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

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

- Train base learner using distribution $D_t$.
- Get base classifier $h_t : X \to \mathbb{R}$.
- Choose $\alpha_t \in \mathbb{R}$. ⟵ $\boxed{\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)}$
- Update:

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

$\epsilon_t$ is weighted error of $h_t(x)$ at iteration $t$:

really trust

$$\boxed{\epsilon_t = P_{i \sim D_i}[\mathbf{x}^i \neq y^i]}$$

$$\boxed{\epsilon_t = \frac{1}{\sum_{i=1}^{m} D_t(i)} \sum_{i=1}^{m} D_t(i) \delta(h_t(x_i) \neq y_i)}$$

as $\epsilon_t \to 0$, $\alpha_t \to +\infty$

as $\epsilon_t \to 1$, $\alpha_t \to -\infty$ ← really trust opposite

if $\epsilon_t = 0.5$, $\alpha_t = 0$, random classifiers are bad $\Rightarrow$ zero weight

# What $\alpha_t$ to choose for hypothesis $h_t$?

[Schapire, 1989]

Training error of final classifier is bounded by:

*bound (count # mistakes in training)*

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

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

*if $y_i +1$*

$\delta(H(x_i) \neq y_i)$

$\forall_i \quad \delta(H(x_i) \neq y_i) \leq e^{-y_i f(x_i)}$

$y_i = +1$

$H(x_i)$

$y_i f(x_i)$

31

©Carlos Guestrin 2005-2007

---

# What $\alpha_t$ to choose for hypothesis $h_t$?

[Schapire, 1989]

*[magic of telescoping]*

*your homework answer*

Training error of final classifier is bounded by:

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

Where $f(x) = \sum_{t=1}^{T} \alpha_t h_t(x); H(x) = sign(f(x))$

*upper bound on train error*

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

$\prod_t Z_t$

*(some conditions apply...)*

*error train*

*iterations t*

32

©Carlos Guestrin 2005-2007

# What $\alpha_t$ to choose for hypothesis $h_t$?

Training error of final classifier is bounded by:

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

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

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

$Z_{t-1}$ doesn't depend on $\alpha_t, h_t$

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

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

33

---

# What $\alpha_t$ to choose for hypothesis $h_t$?

We can minimize this bound by choosing $\alpha_t$ 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]:

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

You'll prove this in your homework! ☺

34

# Strong, weak classifiers

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

- AdaBoost will achieve zero *training error* (exponentially fast):

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

*[handwritten: $e^{-2T\gamma^2}$ exponentially fast, error train]*

*[handwritten: $(1/2 - \epsilon_t)^2 \leftarrow$ how much better is $\epsilon_t$ than random]*

- Is it hard to achieve better than random training error?

*[handwritten: $T$]*

*[handwritten: $\left|\frac{1}{2} - \epsilon_t\right| \geq \gamma$]*

---

# Boosting results – Digit recognition

[Schapire, 1989]



*[handwritten annotations on graph: error test, error train, even when error train = 0, error test still decreasing... don't know when to stop]*

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