

Boosting

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 24, 2013

©Carlos Guestrin 2005-2013

1

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 too badly
- **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...**

©Carlos Guestrin 2005-2013

2

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!

$$h_t: X \rightarrow Y: \{-1, +1\}$$

α_t is strength of vote

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

weak classifier
weight of vote

- But how do you ???**

- force classifiers to learn about different parts of the input space?
- weigh the votes of different classifiers? α_t

©Carlos Guestrin 2005-2013

3

Boosting [Schapire, 1989]

decision stumps, DT, LR, Naive Bayes
... your choice

$$h_t: X \rightarrow \{-1, +1\}$$

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

$$h_t(x) \rightarrow \{-1, +1\}, Y \in \{-1, +1\}$$

- On each iteration t :

- weight each training example by how incorrectly it was classified so far
- Learn a hypothesis h_t ← focus on 'difficult' parts of the space ⇒ increase weight
- A strength for this hypothesis α_t

$$y_i h_t(x_i) > 0 \Rightarrow \text{correct classification}$$

$$y_i h_t(x_i) < 0 \Rightarrow \text{incorrect classification}$$

- Final classifier:

Output must be +1 or -1

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

- Practically useful**
- Theoretically interesting**

©Carlos Guestrin 2005-2013

4

Learning from weighted data

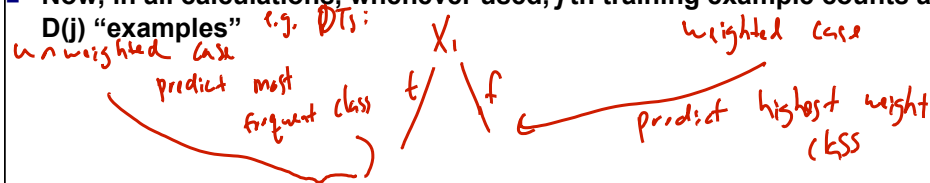
Sometimes not all data points are equal

- Some data points are more equal than others

Consider a weighted dataset

- $D(j)$ – weight of j th training example (x^j, y^j)
- Interpretations:
 - j th training example counts as $D(j)$ examples
 - If I were to “resample” data, I would get more samples of “heavier” data points

Now, in all calculations, whenever used, j th training example counts as $D(j)$ “examples” e.g. DTs:



©Carlos Guestrin 2005-2013

5

AdaBoost

- Initialize weights to uniform dist: $D_1(j) = 1/N$

- For $t = 1 \dots T$

- Train weak learner h_t on distribution D_t over the data

- Choose weight $\alpha_t > 0$ (usually)

- Update weights: *magic!* for endpoint j

$$D_{t+1}(j) = \frac{D_t(j) \exp(-\alpha_t y^j h_t(x^j))}{Z_t}$$

- Where Z_t is normalizer:

$$Z_t = \sum_{j=1}^N D_t(j) \exp(-\alpha_t y^j h_t(x^j))$$

- Output final classifier:

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

if $y^j h_t(x^j) > 0$
 \Rightarrow correct class
 \Rightarrow weight decreases

if $y^j h_t(x^j) < 0$
 \Rightarrow incorrect class
 \Rightarrow weight increase

©Carlos Guestrin 2005-2013

6

Picking Weight of Weak Learner

- Weigh h_t higher if it did well on training data (weighted by D_t):

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

e.g. $\epsilon_t = 1 \Rightarrow$ perfectly wrong \Rightarrow

$\alpha_t = -\infty \Rightarrow h_t$ is perfectly wrong
 $\Rightarrow -h_t$ is perfectly right

- Where ϵ_t is the weighted training error:

$$\epsilon_t = \sum_{j=1}^n D_t(j) \mathbb{1}[h_t(x^j) \neq y^j]$$

weight

if $\epsilon_t = 0 \Rightarrow h_t$ perfect on weighted data \Rightarrow perfect in general $\Rightarrow \alpha_t = \infty$

if $\epsilon_t = \frac{1}{2} \Rightarrow h_t$ as good as random $\Rightarrow \alpha_t = 0$

indicator fa. $\mathbb{1}(x=y) = \begin{cases} 1 & \text{if } x=y \\ 0 & \text{otherwise} \end{cases}$

©Carlos Guestrin 2005-2013

7

Why choose α_t for hypothesis h_t this way?

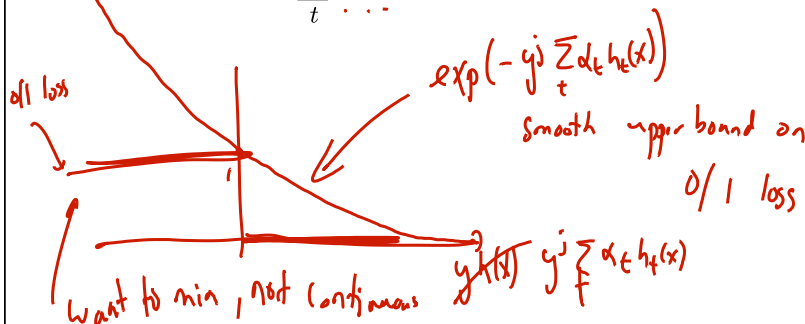
[Schapire, 1989]

Training error of final classifier is bounded by:

unweighted train error $\rightarrow \frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j))$

0/1 loss

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



©Carlos Guestrin 2005-2013

8

Why choose α_t for hypothesis h_t this way?

[Schapire, 1989]

Training error of final classifier is bounded by:

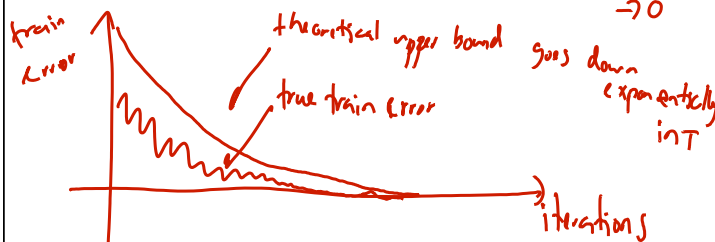
$$Z_t = \sum_{j=1}^N D_t(j) \exp(-\alpha_t y^j h_t(x^j))$$

training error

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j)) = \prod_{t=1}^T Z_t$$

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

if $Z_t < 1 \forall t \Rightarrow$ as $t \rightarrow \infty$, training error $\rightarrow 0$



Proof by
homework:
use "telescopic
sums"

©Carlos Guestrin 2005-2013

9

Why choose α_t for hypothesis h_t this way?

[Schapire, 1989]

Training error of final classifier is bounded by:

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j)) = \prod_{t=1}^T Z_t$$

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 .

$$Z_t = \sum_{j=1}^N D_t(j) \exp(-\alpha_t y^j h_t(x^j))$$

©Carlos Guestrin 2005-2013

10

Why choose α_t for hypothesis h_t this way?

[Schapire, 1989]

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

$$Z_t = \sum_{j=1}^N D_t(j) \exp(-\alpha_t y^j h_t(x^j))$$

*take derivative
set to 0*

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

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

You'll prove this in your homework! ☺

©Carlos Guestrin 2005-2013

11

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):

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x^j) \neq y^j] \leq \prod_{t=1}^T Z_t \leq \exp \left(-2 \sum_{t=1}^T (1/2 - \epsilon_t)^2 \right)$$

$\epsilon_t < \frac{1}{2}$ weak classifier must be strictly better than random

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

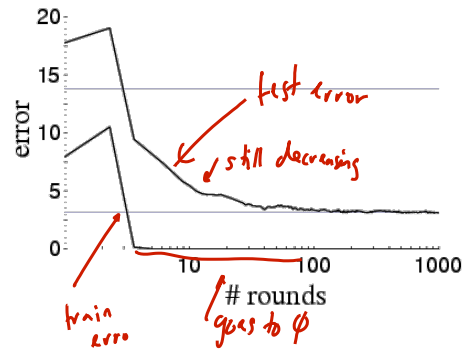
how well does he do on training data?

©Carlos Guestrin 2005-2013

12

Boosting results – Digit recognition

[Schapire, 1989]



Can't use
train error
to decide
when to stop

of rounds?
→ out of a hat
→ single
validation data
→ cross-validation

■ Boosting often

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

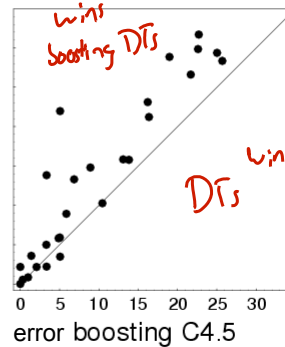
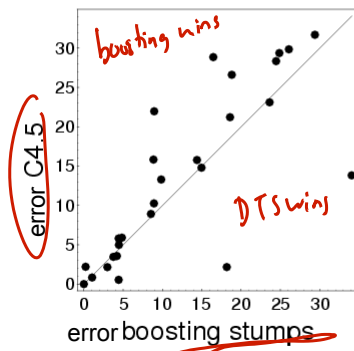
©Carlos Guestrin 2005-2013

13

Boosting: Experimental Results

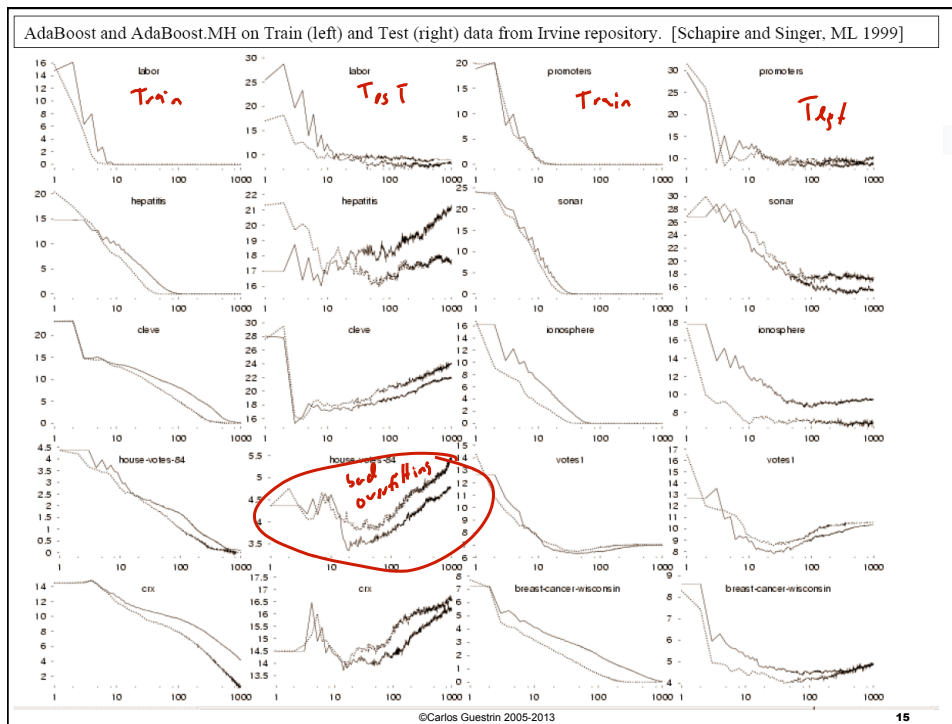
[Freund & Schapire, 1996]

Comparison of C4.5, Boosting C4.5, Boosting decision stumps (depth 1 trees), 27 benchmark datasets



©Carlos Guestrin 2005-2013

14



Boosting and Logistic Regression

Logistic regression assumes:

$$f(x) = w_0 + \sum_i w_i h_i(x)$$

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

And tries to maximize data likelihood:

$$\max_w P(\mathcal{D}|H) = \prod_{j=1}^N \frac{1}{1 + \exp(-y^j f(x^j))} \stackrel{\text{red}}{=} \min_w -\ln P(\mathcal{D}|H)$$

Equivalent to minimizing log loss

$$\sum_{j=1}^N \ln(1 + \exp(-y^j f(x^j)))$$

Boosting and Logistic Regression

Logistic regression equivalent to minimizing log loss

$$\sum_{j=1}^N \ln(1 + \exp(-y^j f(x^j)))$$

✓ LR loss function

Boosting minimizes similar loss function!!

$$\frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j)) = \prod_{t=1}^T Z_t$$

Both smooth approximations of 0/1 loss!



©Carlos Guestrin 2005-2013

17

Logistic regression and Boosting

Logistic regression:

- Minimize loss fn

$$\sum_{j=1}^N \ln(1 + \exp(-y^j f(x^j)))$$

- Define

$$f(x) = w_0 + \sum_i w_i h_i(x)$$

h_i(x) model

where features h_i are predefined

- Weights w_i are learned in joint optimization

Boosting:

- Minimize loss fn

$$\sum_{j=1}^N \exp(-y^j f(x^j))$$

- Define

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

where $h_t(x)$ defined dynamically to fit data

(not a linear classifier)

via weak classifiers

- Weights α_t learned incrementally

©Carlos Guestrin 2005-2013

18

What you need to know about Boosting

- Combine weak classifiers to obtain very 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