

Module 5: Classification

Support Vector Machines

≈ perceptron + regularization

STAT/BIOSTAT 527, University of Washington

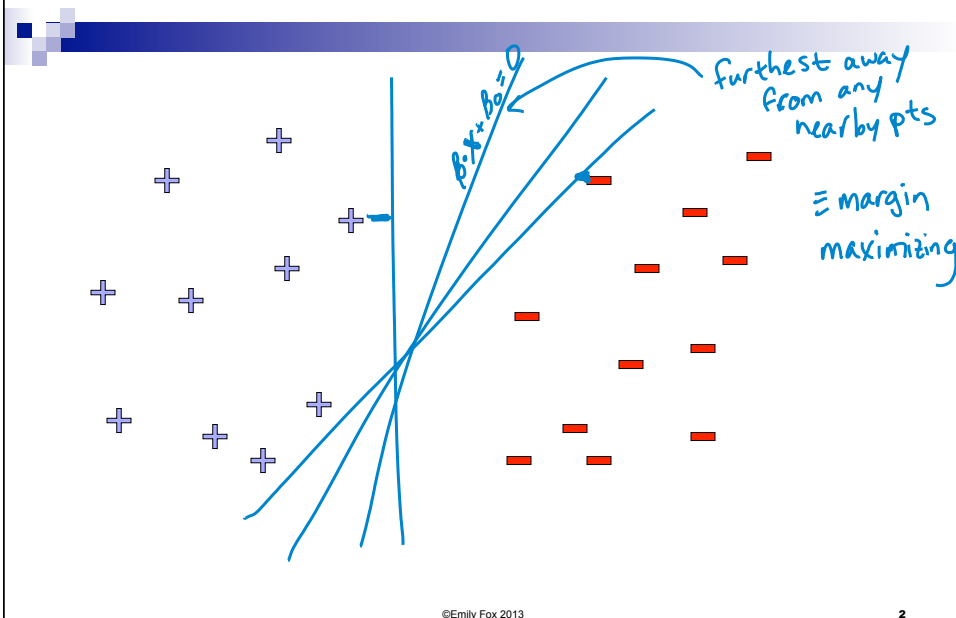
Emily Fox

June 6th, 2013

©Emily Fox 2013

1

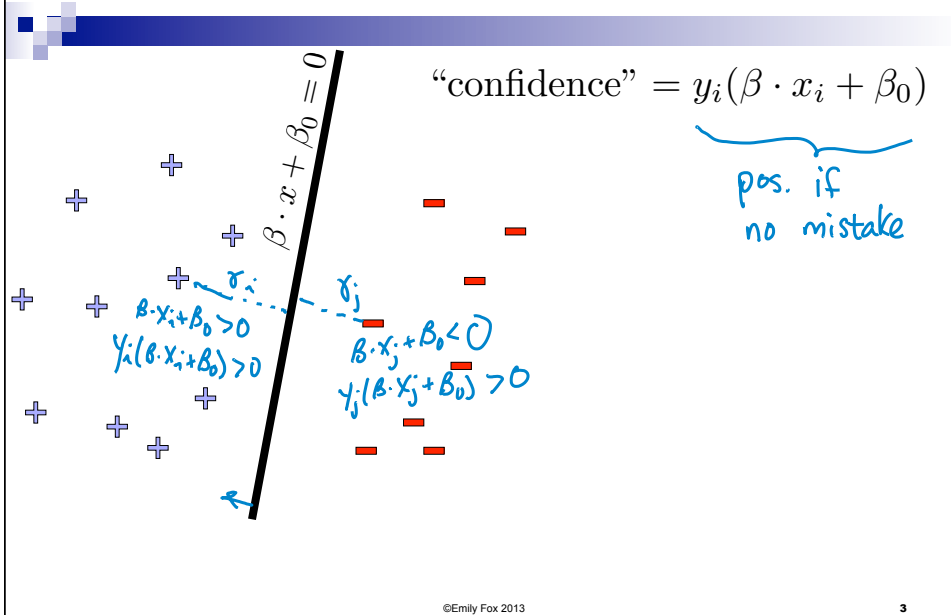
Linear classifiers – Which line is better?



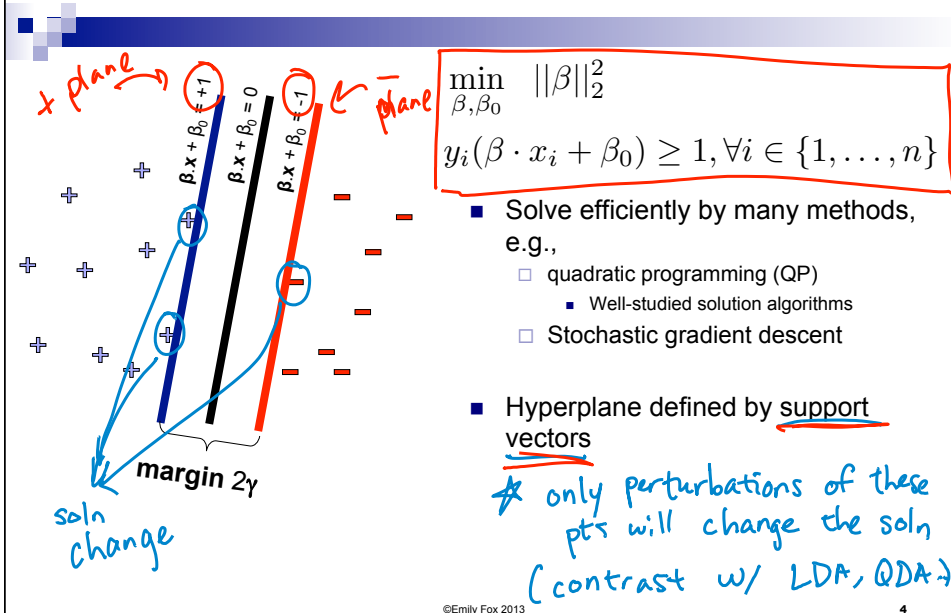
©Emily Fox 2013

2

Pick the one with the largest margin!



Support vector machines (SVMs)



What if the data are still not linearly separable?

Now support vectors are
- on margin
or
- misclassified pts

$$\min_{\beta, \beta_0} \|\beta\|_2^2$$

$$y_i(\beta \cdot x_i + \beta_0) \geq 1, \forall i \in \{1, \dots, n\}$$

- If data are not linearly separable, some points don't satisfy margin constraint:
 - $\exists i$ s.t. $y_i(\beta \cdot x_i + \beta_0) < 1$
 - or $1 - y_i(\beta \cdot x_i + \beta_0) > 0$ like hinge loss
- How bad is the violation?

$$\text{constraint violation} = \begin{cases} 0 & \text{if } 1 - y_i(\beta \cdot x_i + \beta_0) \leq 0 \\ 1 - y_i(\beta \cdot x_i + \beta_0) & \text{otherwise} \end{cases}$$
- Tradeoff margin violation with $\|\beta\|$:

$$\min_{\beta, \beta_0} \|\beta\|_2^2 + C \sum_{i=1}^n (1 - y_i(\beta \cdot x_i + \beta_0))_+$$

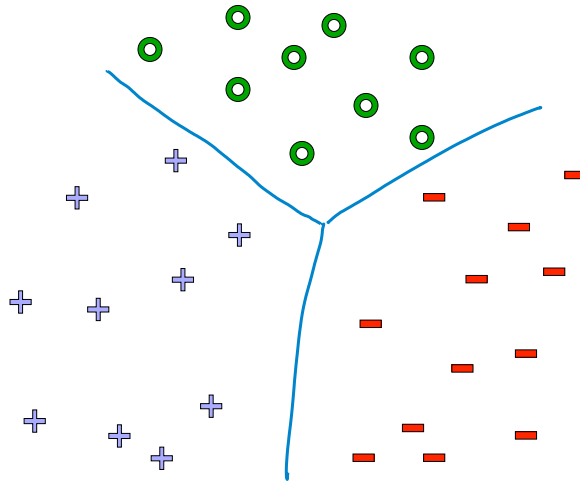
SVM objective

©Emily Fox 2013 5

Demo!

<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

What about multiple classes?

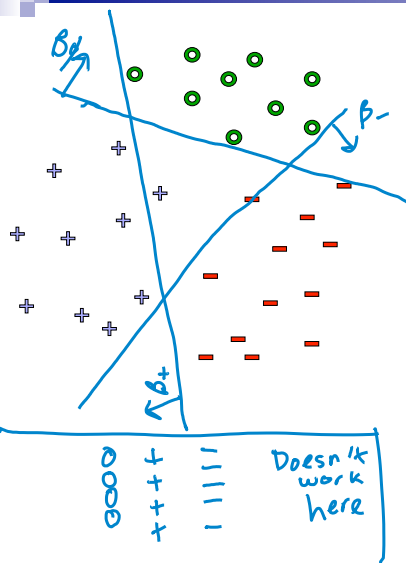


want boundaries like this

©Emily Fox 2013

7

One against All



Learn 3 classifiers:

$+$ vs. $\{0, -\}$	β_+
$-$ vs. $\{0, +\}$	β_-
0 vs. $\{+, -\}$	β_0

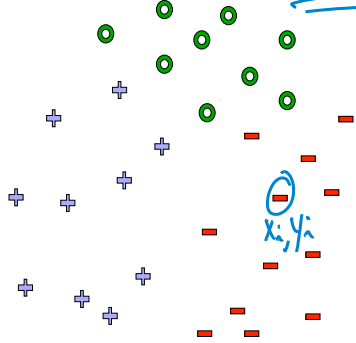
Output for x^* :
 $\beta^{(1)}, \dots, \beta^{(k)}$ for k classes
 $y^* = \arg \max_k \beta^{(k)} \cdot x^* + \beta_0$
smallest margin

©Emily Fox 2013

8

Learn 1 classifier: Multiclass SVM

Simultaneously learn 3 sets of weights



$\beta^{(k)}, \beta_i^{(k)}$ for each class
 jointly optimize to achieve

$$\beta^{(y_i)} \cdot x_i + \beta_0^{(y_i)} \geq 1 + \beta^{(y')} \cdot x_i + \beta_0^{(y')} \quad \forall y' \neq y_i, \forall i$$

correct class wins

exceeds plane w/ margin 1 just like before

$$\beta^{(y_i)} \cdot x_i + \beta_0^{(y_i)} \geq \beta^{(y')} \cdot x_i + \beta_0^{(y')} + 1, \forall y' \neq y_i, \forall i$$

©Emily Fox 2013

9

Learn 1 classifier: Multiclass SVM

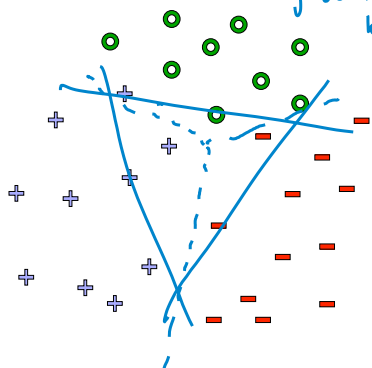
$$\beta^{(y_i)} \cdot x_i + \beta_0^{(y_i)} \geq \beta^{(y')} \cdot x_i + \beta_0^{(y')} + 1, \forall y' \neq y_i, \forall i$$

$$\min_{\beta, \beta_0} \sum_y \beta^{(y)} \cdot \beta^{(y)} + C \sum_i \xi_i$$

$$\xi_i = \sum_{y \neq y_i} (\beta^{(y)} \cdot x_i + \beta_0^{(y)} + 1 - (\beta^{(y_i)} \cdot x_i + \beta_0^{(y_i)}))_+$$

just like "1" before

same soln technique as for 2 class setting



©Emily Fox 2013

10

turned sideways

What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Non-linearly separable case
 - Hinge loss
 - A.K.A. adding slack variables
- SVMs = Perceptron + L2 regularization
- Can optimize SVMs with SGD
 - Many other approaches possible
- Handling multiple classes

Reading

- Hastie, Tibshirani, Friedman – 12.1-12.3

Module 5: Classification

Boosting

STAT/BIOSTAT 527, University of Washington

Emily Fox

June 6th, 2013

©Emily Fox 2013

13

Fighting the Bias-Variance Tradeoff

- **Simple (a.k.a. weak) learners are good**
 - e.g., naïve Bayes, logistic regression, shallow decision trees ... *decision stumps* | —
 - 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...** *how?*

©Emily Fox 2013

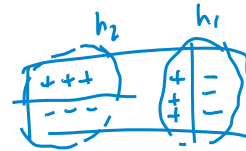
14

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!

$h_t: X \rightarrow \{+1, -1\}$
 $\alpha_t = \text{strength of vote}$

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



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

Boosting [Schapire, 1989]

decision tree, log. reg., naive Bayes, ...

- 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 so far
 - Learn a hypothesis - h_t
 - A strength for this hypothesis - α_t

$y_i h_t(x_i) > 0 \Rightarrow \text{correct class}$
 $y_i h_t(x_i) < 0 \Rightarrow \text{incorrect}$

- Final classifier:

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

- Practically useful**
- Theoretically interesting**

Learning from Weighted Data

- Sometimes not all data points are equal

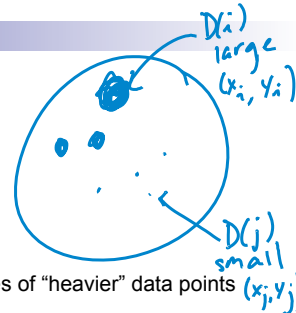
- Some data points are more equal than others

- Consider a weighted dataset

- $D(i)$ – weight of i th training example (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”

Example: decision trees

unweighted case:

predict most frequent class



weighted case:

predict highest weight class

©Emily Fox 2013

17

AdaBoost

- Initialize weights to uniform dist: $D_1(i) = 1/n$

- For $t = 1 \dots T$

- Train weak learner h_t on distribution D_t over the data

- Choose weight $\alpha_t > 0$ (usually) ... will get back to magic for now

- Update weights:

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

- Where Z_t is normalizer:

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

- Output final classifier:

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

If $y_i h_t(x_i) > 0$
 \Rightarrow correct class
 \Rightarrow weight decreases

If $y_i h_t(x_i) < 0$
 \Rightarrow incorrect
 \Rightarrow weight inc.

☆☆ focus more here

©Emily Fox 2013

18

Picking Weight of Weak Learner

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

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

- Where ϵ_t is the weighted training error:

$$\epsilon_t = \sum_{i=1}^n D_t(i) \mathbb{I}[h_t(x_i) \neq y_i]$$

Annotations: $D_t(i)$ is weight; $\mathbb{I}[h_t(x_i) \neq y_i]$ is made mistake; y_i is true class.

CASES

$\epsilon_t = 1 \Rightarrow$ perfectly wrong

$d_t = -\infty$
 h_t perfectly wrong
 $-h_t$ perfectly right

$\epsilon_t = 0 \Rightarrow$ perfect

$\alpha_t = \infty$

$\epsilon_t = \frac{1}{2} \Rightarrow$ random guess

$\alpha_t = 0$

©Emily Fox 2013

19

Why choose α_t for hypothesis h_t this way?

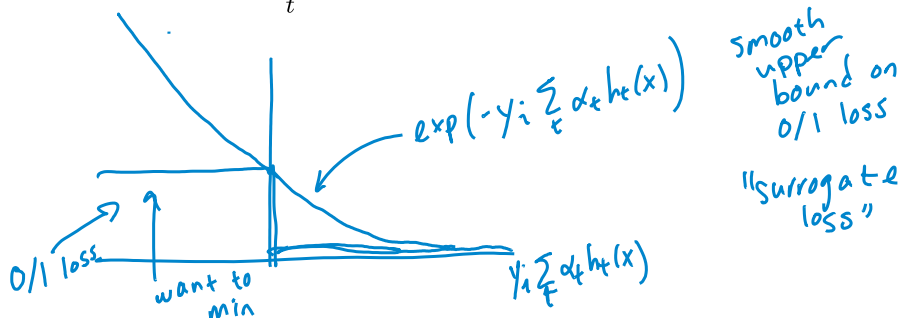
[Schapire, 1989]

Training error of final classifier is bounded by:

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

Annotation: $\mathbb{I}[H(x_i) \neq y_i]$ is 0/1 loss.

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



©Emily Fox 2013

20

Why choose α_t for hypothesis h_t this way?

[Schapire, 1989]

Training error of final classifier is bounded by:

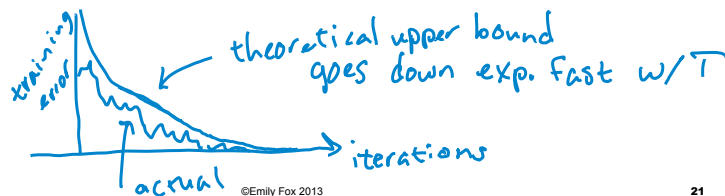
$$Z_t = \sum_{i=1}^n D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}[H(x_i) \neq y_i] \leq \frac{1}{n} \sum_{i=1}^n \exp(-y_i f(x_i)) = \prod_{t=1}^T Z_t$$

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

Pf: uses "telescopic sums"

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



©Emily Fox 2013

21

Why choose α_t for hypothesis h_t this way?

[Schapire, 1989]

Training error of final classifier is bounded by:

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}[H(x_i) \neq y_i] \leq \frac{1}{n} \sum_{i=1}^n \exp(-y_i f(x_i)) = \prod_{t=1}^T Z_t$$

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 .

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

©Emily Fox 2013

22

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 .

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

fake der. + set = 0

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

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

Strong, weak classifiers

counter examples exist, but in general it's fine

+	-
-	+
+	-
-	+
+	-
-	+
⋮	⋮

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

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

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}[H(x_i) \neq y_i] \leq \prod_{t=1}^T Z_t \leq \exp \left(-2 \sum_{t=1}^T (1/2 - \epsilon_t)^2 \right)$$

upper bound for $\epsilon_t < 1/2$

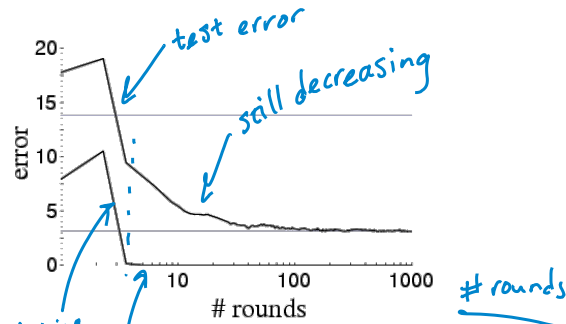
as ϵ_t moves away from $1/2$, gets bigger

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

Boosting results – Digit recognition

[Schapire, 1989]

0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9
 0 1 2 3 4 5 6 7 8 9



can't use training error to decide when to stop

Boosting often

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

- guess
 - validation data
 - CV

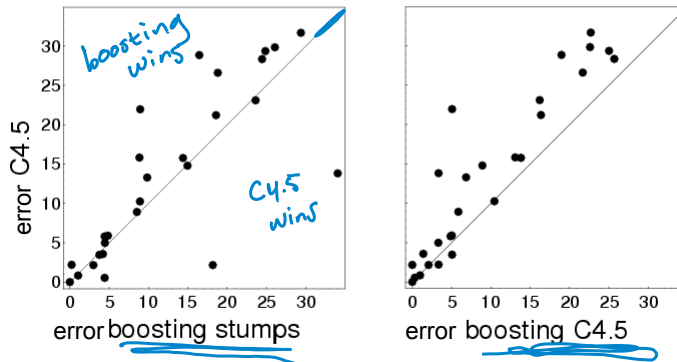
©Emily Fox 2013

25

Boosting: Experimental Results

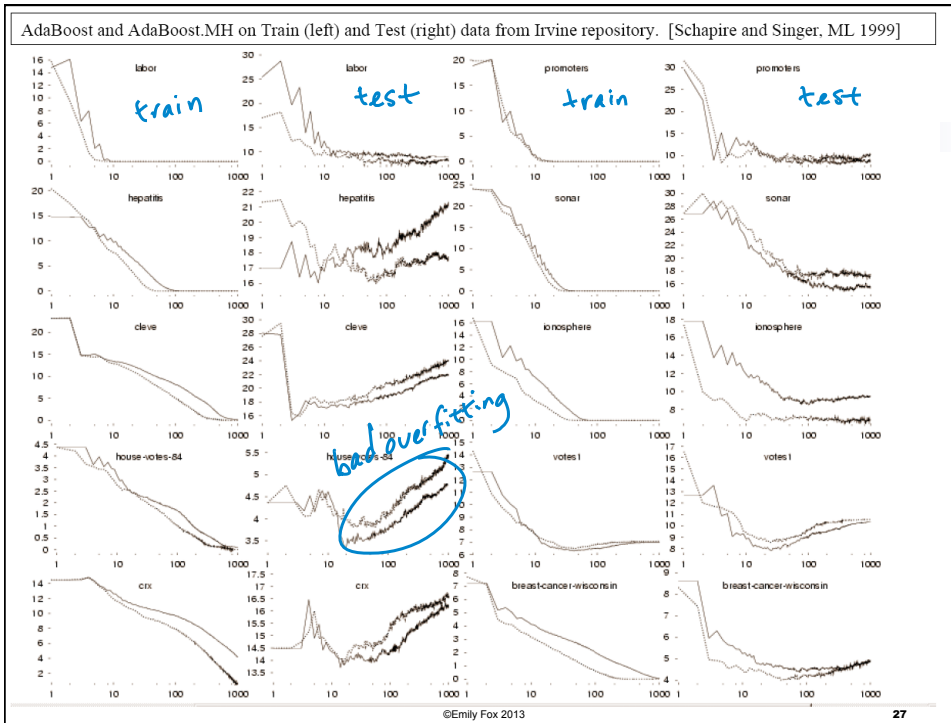
[Freund & Schapire, 1996]

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



©Emily Fox 2013

26



Demo!

<http://cseweb.ucsd.edu/~yfreund/adaboost/>

Boosting and Logistic Regression

Logistic regression assumes:

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

$$f(x) = \beta_0 + \sum_j \beta_j h_j(x)$$

And tries to maximize data likelihood:

$$P(\mathcal{D}|H) = \prod_{i=1}^n \frac{1}{1 + \exp(-y_i f(x_i))} \Rightarrow \min_{\beta} -\log P(\mathcal{D}|H)$$

Equivalent to minimizing log loss

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

©Emily Fox 2013

29

Boosting and Logistic Regression

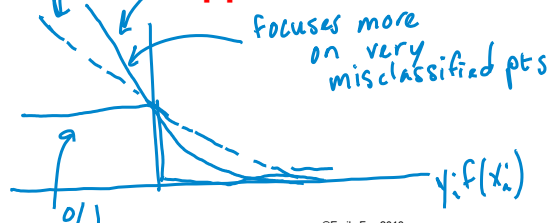
Logistic regression equivalent to minimizing log loss

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

Boosting minimizes similar loss function!!

$$\frac{1}{n} \sum_{i=1}^n \exp(-y_i f(x_i)) = \prod_{t=1}^T Z_t$$

Both smooth approximations of 0/1 loss!



©Emily Fox 2013

30

Logistic regression and Boosting

Logistic regression:

- Minimize loss fn

$$\sum_{i=1}^n \ln(1 + \exp(-y_i f(x_i))) \xleftrightarrow{\text{loss fn}} \sum_{i=1}^n \exp(-y_i f(x_i))$$

- Define

$$f(x) = \beta_0 + \sum_j \beta_j h(x_j)$$

where features x_j are predefined

- Weights β_j are learned in joint optimization

Boosting:

- Minimize loss fn

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

- 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 incrementally

©Emily Fox 2013

31

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

©Emily Fox 2013

32

Reading

- Hastie, Tibshirani, Friedman – 10.1-10.6

What you need to know from 527

- **Module 1: Preliminaries of Nonparametrics**
 - Loss functions and optimal predictions
 - Linear smoothers
 - Ridge regression, LASSO
 - Cross validation, etc.
- **Module 2: Splines and Kernel Methods**
 - Smoothing splines, penalized regression splines
 - Local polynomial regression, KDE
- **Module 3: Bayesian Nonparametrics**
 - Gaussian processes
 - Dirichlet process mixture of Gaussians

What you need to know from 527

- **Module 4: Nonparametrics with Multiple Predictors**

- Thin plate splines, tensor product splines
- GAMs
- Projection pursuit
- Multivariate kernels and KDE
- Regression trees

- **Module 5: Classification**

- Classification trees
- Logistic regression (also looked at nonparametrics for GLMs)
- LDA, QDA, KDE, naïve Bayes, mixture models
- Perceptron, SVM and with kernels for both
- Boosting

©Emily Fox 2013

35

THANK YOU!!!

- You have been a great, interactive class!
...especially for a 9am lecture =)
- We're looking forward to the poster session

- Thanks to Shirley, too!

©Emily Fox 2013

36