

ENSEMBLE LEARNING

David Kauchak
CS158 – Spring 2022

1

Admin

- Assignment 9
- Midterm 2
- Final project
 - Monday (4/18) submit project proposal

2

Ensemble learning

Basic idea: if one classifier works well, why not use multiple classifiers!

4

Ensemble learning

Basic idea: if one classifier works well, why not use multiple classifiers!

Training

```

graph LR
    TD[Training Data] --> LA1[learning alg]
    TD --> LA2[learning alg]
    TD --> LA3[learning alg]
    LA1 --> M1[model 1]
    LA2 --> M2[model 2]
    LA3 --> Mm[model m]
  
```

5

Ensemble learning

Basic idea: if one classifier works well, why not use multiple classifiers!

Testing

How do we decide on the final prediction?

6

Ensemble learning

Basic idea: if one classifier works well, why not use multiple classifiers!

Testing

- take majority vote
- if they output probabilities, take a weighted vote

How does having multiple classifiers help?

7

Benefits of ensemble learning

Assume each classifier makes a mistake with some probability (e.g. 0.4, that is a 40% error rate)

Assuming the decisions made between classifiers are independent, what will be the probability that we make a mistake (i.e. error rate) with three classifiers for a binary classification problem?

8

Benefits of ensemble learning

Assume each classifier makes a mistake with some probability (e.g. 0.4, that is a 40% error rate)

model 1	model 2	model 3	prob
C	C	C	.6*.6*.6=0.216
C	C	I	.6*.6*.4=0.144
C	I	C	.6*.4*.6=0.144
C	I	I	.6*.4*.4=0.096
I	C	C	.4*.6*.6=0.144
I	C	I	.4*.6*.4=0.096
I	I	C	.4*.4*.6=0.096
I	I	I	.4*.4*.4=0.064

9

Benefits of ensemble learning

Assume each classifier makes a mistake with some probability (e.g. 0.4, that is a 40% error rate)

model 1	model 2	model 3	prob
C	C	C	.6*.6*.6=0.216
C	C	I	.6*.6*.4=0.144
C	I	C	.6*.4*.6=0.144
C	I	I	.6*.4*.4=0.096
I	C	C	.4*.6*.6=0.144
I	C	I	.4*.6*.4=0.096
I	I	C	.4*.4*.6=0.096
I	I	I	.4*.4*.4=0.064

0.096+
0.096+
0.096+
0.064 =
35% error!

10

Benefits of ensemble learning

3 classifiers in general, for r = probability of mistake for individual classifier:

$$p(error) = 3r^2(1-r) + r^3$$

binomial distribution

r	p(error)
0.4	0.35
0.3	0.22
0.2	0.10
0.1	0.028
0.05	0.0073

11

Benefits of ensemble learning

5 classifiers in general, for r = probability of mistake for individual classifier:

$$p(error) = 10r^3(1-r)^2 + 5r^4(1-r) + r^5$$

r	p(error) 3 classifiers	p(error) 5 classifiers
0.4	0.35	0.32
0.3	0.22	0.16
0.2	0.10	0.06
0.1	0.028	0.0086
0.05	0.0073	0.0012

12

Benefits of ensemble learning

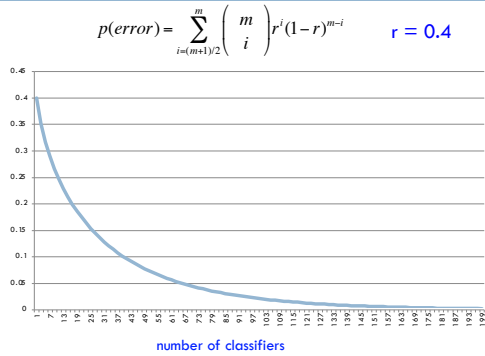
m classifiers in general, for r = probability of mistake for individual classifier:

$$p(error) = \sum_{i=(m+1)/2}^m \binom{m}{i} r^i (1-r)^{m-i}$$

(cumulative probability distribution for the binomial distribution)

13

Given enough classifiers...



14

What is the catch?

Assume each classifier makes a mistake with some probability (e.g. 0.4, that is a 40% error rate)

model 1

model 2

model 3

Assuming the decisions made between classifiers are independent, what will be the probability that we make a mistake (i.e. error rate) with three classifiers for a binary classification problem?

15

What is the catch?

Assume each classifier makes a mistake with some probability (e.g. 0.4, that is a 40% error rate)

model 1

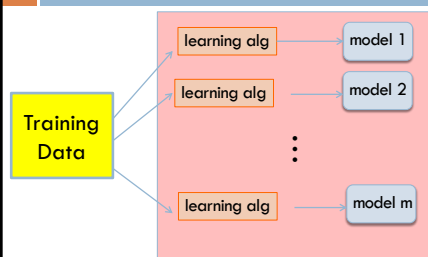
model 2

model 3

Assuming the decisions made between classifiers are independent, what will be the probability that we make a mistake (i.e. error rate) with three classifiers for a binary classification problem?

16

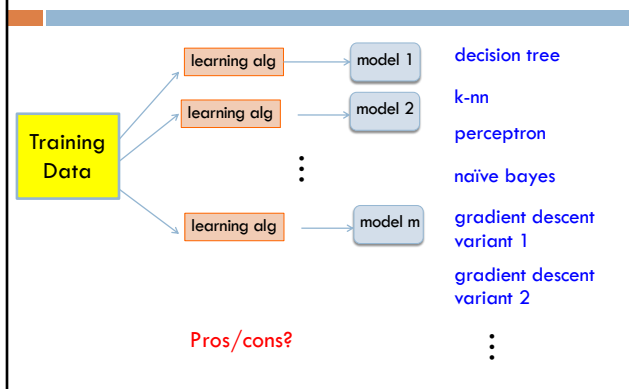
Obtaining independent classifiers



Where do we get m independent classifiers?

17

Idea 1: different learning methods



18

Idea 1: different learning methods

Pros:

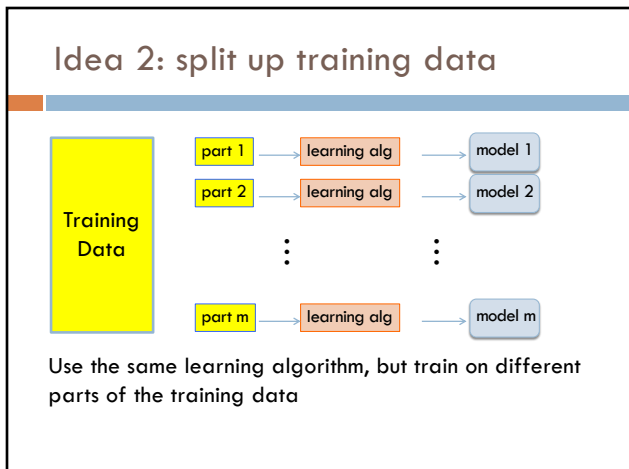
- ▣ Lots of existing classifiers already
- ▣ Can work well for some problems

Cons/concerns:

- ▣ Often, classifiers are not independent, that is, **they make the same mistakes!**
 - e.g. many of these classifiers are linear models
 - voting won't help us if they're making the same mistakes

19

Idea 2: split up training data



20

Idea 2: split up training data

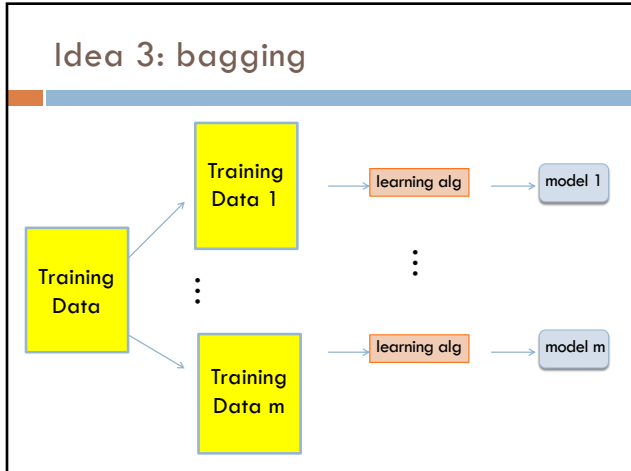
Pros:

- ▣ Learning from different data, so can't overfit to same examples
- ▣ Easy to implement
- ▣ fast

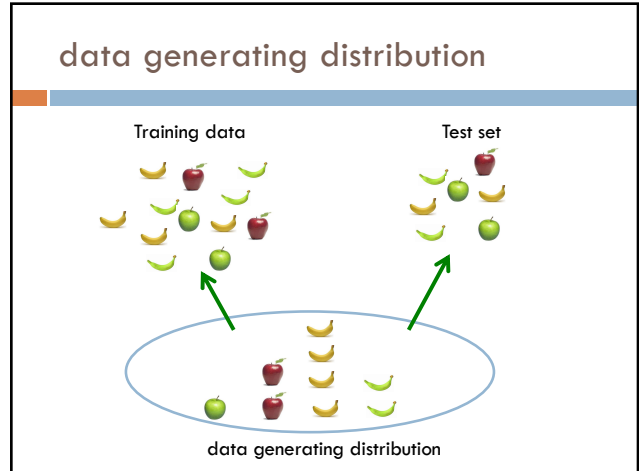
Cons/concerns:

- ▣ Each classifier is only training on a small amount of data
- ▣ Not clear why this would do any better than training on full data and using good regularization

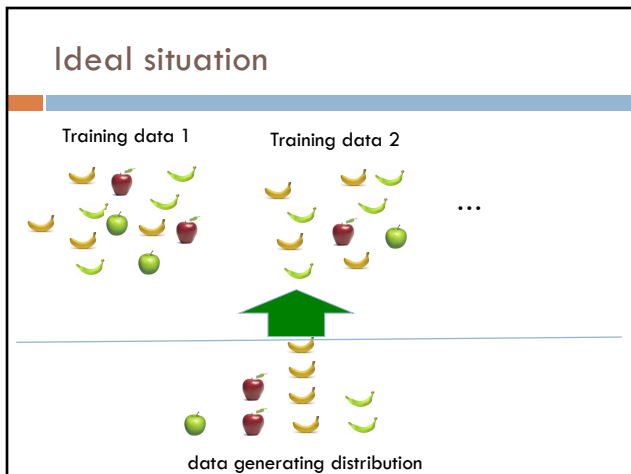
21



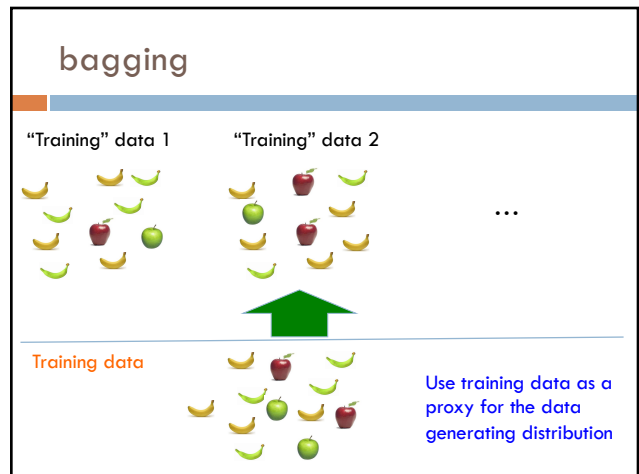
22



23



24



25

sampling with replacements

"Training" data 1

Training data



26

sampling with replacements

"Training" data 1

pick a random example from the
real training data

Training data



27

sampling with replacements

"Training" data 1



add it to the new "training" data

Training data



28

sampling with replacements

"Training" data 1



put it back (i.e. leave it) in the
original training data


Training data



29

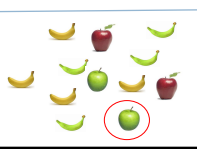
sampling with replacements

"Training" data 1



pick another random example


Training data



30

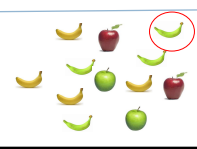
sampling with replacements

"Training" data 1



pick another random example


Training data



31

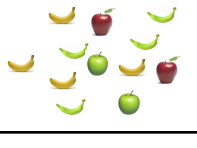
sampling with replacements

"Training" data 1



keep going until you've created
a new "training" data set

Training data



32

bagging

create m "new" training data sets by sampling with replacement from the original training data set (called m "bootstrap" samples)

train a classifier on each of these data sets

to classify, take the majority vote from the m classifiers

33

bagging concerns

Training Data

Training Data 1

⋮

Training Data m

Won't these all be basically the same?

34

bagging concerns

For a data set of size n , what is the probability that a given example will **NOT** be select in a "new" training set sampled from the original?

Training data

35

bagging concerns

What is the probability it isn't chosen the first time?

$$1 - 1/n$$

Training data

36

bagging concerns

What is the probability it isn't chosen the **any** of the n times?

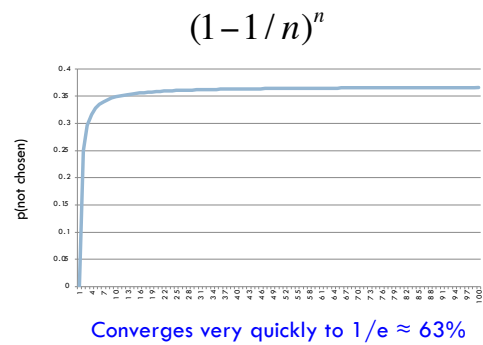
$$(1 - 1/n)^n$$

Each draw is independent and has the same probability

Training data

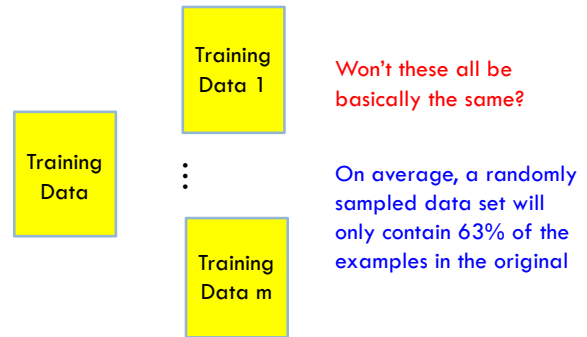
37

probability of overlap



38

bagging overlap



39

When does bagging work

Let's say 10% of our examples are noisy (i.e. don't provide good information)

For each of the "new" data set, what proportion of noisy examples will they have?

- They'll still have ~10% of the examples as noisy
- However, these examples will only represent about two-thirds of the original noisy examples

For some classifiers that have trouble with noisy classifiers, this can help

40

When does bagging work

Bagging tends to reduce the *variance* of the classifier

By voting, the classifiers are more robust to noisy examples

Bagging is most useful for classifiers that are:

- Unstable: small changes in the training set produce very different models
- Prone to overfitting

Often has similar effect to regularization


41

Idea 4: boosting

training data			"training" data 2			"training" data 3		
Data	Label	Weight	Data	Label	Weight	Data	Label	Weight
■	0	0.2	■	0	0.1	■	0	0.05
■	0	0.2	■	0	0.1	■	0	0.2
■	1	0.2	■	1	0.4	■	1	0.2
■	1	0.2	■	1	0.1	■	1	0.05
■	0	0.2	■	0	0.3	■	0	0.5

42

"Strong" learner




Given

- a reasonable amount of training data
- a target error rate ϵ
- a failure probability p

A **strong learning algorithm** will produce a classifier with error rate $< \epsilon$ with probability $1-p$

43

"Weak" learner



Given

- a reasonable amount of training data
- a failure probability p

A **weak learning algorithm** will produce a classifier with error rate < 0.5 with probability $1-p$

Weak learners are much easier to create!

44

weak learners for boosting

Data	Label	Weight
■	0	0.2
■	0	0.2
■	1	0.2
■	1	0.2
■	0	0.2

weak learning algorithm → weak classifier

Which of our algorithms can handle weights?

Need a weak learning algorithm that can handle **weighted** examples

45

boosting: basic algorithm

Training:
start with equal example weights

for some number of iterations:

- learn a weak classifier and save
- change the example weights

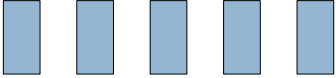
Classify:

- get prediction from all learned weak classifiers
- weighted vote based on how well the weak classifier did when it was trained (i.e. in relation to training error)


46

boosting basics

Start with equal weighted examples

Weights: 

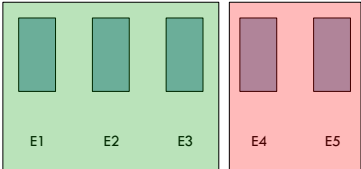
Examples: E1 E2 E3 E4 E5

Learn a weak classifier: 


47

Boosting

classified correct classified incorrect

Weights: 

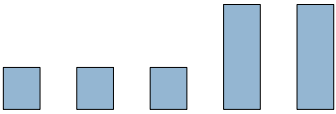
Examples: E1 E2 E3 E4 E5

 We want to reweight the examples and then learn another weak classifier

How should we change the example weights?

48

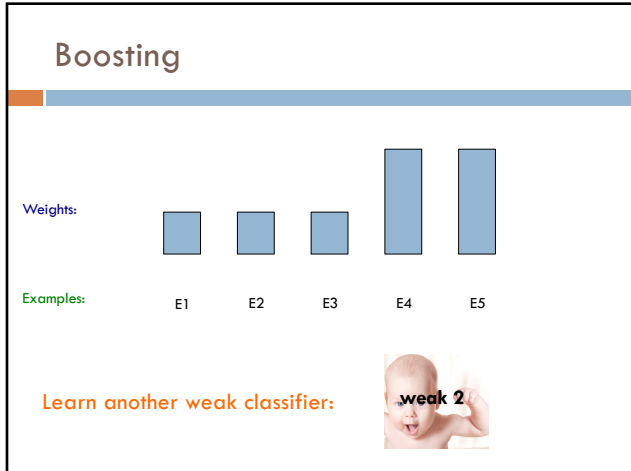
Boosting

Weights: 

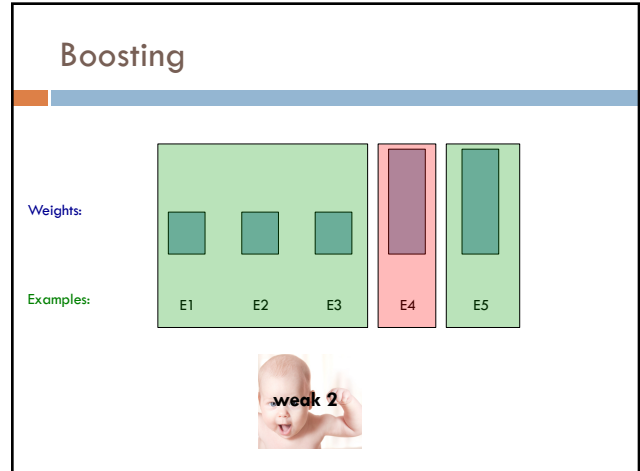
Examples: E1 E2 E3 E4 E5

- decrease the weight for those we're getting correct
- increase the weight for those we're getting incorrect

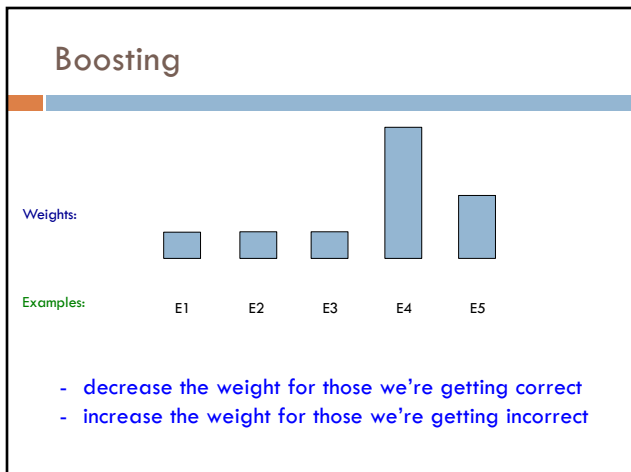
49



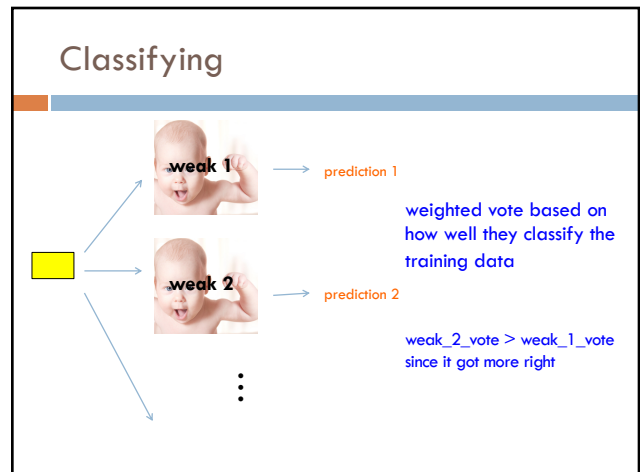
50



51



52



53

Notation

- x_i example i in the training data
- w_i weight for example i , we will enforce:
 $w_i \geq 0$
 $\sum_{i=1}^n w_i = 1$
- $\text{classifier}_k(x_i)$ +1/-1 prediction of classifier k example i

54

AdaBoost: train

for $k = 1$ to *iterations*:

- classifier_k = learn a weak classifier based on weights
- calculate weighted error for this classifier

$$\epsilon_k = \sum_{i=1}^n w_i * \mathbb{1}[\text{label}_i \neq \text{classifier}_k(x_i)]$$

- calculate "score" for this classifier:

$$\alpha_k = \frac{1}{2} \log\left(\frac{1 - \epsilon_k}{\epsilon_k}\right)$$

- change the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

55

AdaBoost: train

classifier_k = learn a weak classifier based on weights

weighted error for this classifier is:

$$\epsilon_k = \sum_{i=1}^n w_i * \mathbb{1}[\text{label}_i \neq \text{classifier}_k(x_i)]$$

What does this say?

56

AdaBoost: train

classifier_k = learn a weak classifier based on weights

weighted error for this classifier is:

$$\epsilon_k = \sum_{i=1}^n w_i * \mathbb{1}[\text{label}_i \neq \underbrace{\text{classifier}_k(x_i)}_{\text{prediction}}]$$

What is the range
of possible values?

prediction

did we get the example wrong

weighted sum of the errors/mistakes

57

AdaBoost: train

classifier_k = learn a weak classifier based on weights

weighted error for this classifier is:

$$\epsilon_k = \sum_{i=1}^n w_i * 1[\text{label}_i \neq \text{prediction}]$$

Between 0 (if we get all examples right) and 1 (if we get them all wrong)

prediction

did we get the example wrong

weighted sum of the errors/mistakes

58

AdaBoost: train

classifier_k = learn a weak classifier based on weights

“score” or weight for this classifier is:

$$\alpha_k = \frac{1}{2} \log\left(\frac{1-\epsilon_i}{\epsilon_i}\right)$$

What does this look like (specifically for errors between 0 and 1)?

59

AdaBoost: train

$$\alpha_k = \frac{1}{2} \log\left(\frac{1-\epsilon_i}{\epsilon_i}\right)$$

- ranges from +∞ to -∞
- for most reasonable values: ranges from ~1 to -1
- errors of 50% = 0
- error > 50% = positive error < 50% = negative

60

AdaBoost: classify

$$\text{classify}(x) = \text{sign}\left(\sum_{k=1}^{\text{iterations}} \alpha_k * \text{classifier}_k(x)\right)$$

What does this do?

61

AdaBoost: classify

$$\text{classify}(x) = \text{sign} \left(\sum_{k=1}^{\text{iterations}} \alpha_k * \text{classifier}_k(x) \right)$$

The weighted vote of the learned classifiers weighted by α (remember α generally varies from ~ 1 to -1 training error)

What happens if a classifier has error $> 50\%$

62

AdaBoost: classify

$$\text{classify}(x) = \text{sign} \left(\sum_{k=1}^{\text{iterations}} \alpha_k * \text{classifier}_k(x) \right)$$

The weighted vote of the learned classifiers weighted by α (remember α generally varies from ~ 1 to -1 training error)

We vote the opposite!

63

AdaBoost: train, updating the weights

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

Remember, we want to enforce:

$$w_i \geq 0$$

$$\sum_{i=1}^n w_i = 1$$

Z is called the **normalizing constant**. It is used to make sure that the weights sum to 1

What should it be?

64

AdaBoost: train

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

Remember, we want to enforce:

$$w_i \geq 0$$

$$\sum_{i=1}^n w_i = 1$$

normalizing constant (i.e. the sum of the "new" w_i):

$$Z = \sum_{i=1}^n w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

65

AdaBoost: train

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

What does this do?

66

AdaBoost: train

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

correct → positive
incorrect → negative

correct ?
incorrect

67

AdaBoost: train

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

correct → positive
incorrect → negative

correct → small value
incorrect → large value

Note: only change weights based on current classifier (not all previous classifiers)

68

AdaBoost: train

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * \text{label}_i * \text{classifier}_k(x_i))$$

What does the α do?

69

AdaBoost: train

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * label_i * classifier_k(x_i))$$

What does the α do?

- If the classifier was good (<50% error) α is positive: trust classifier output and move as normal
- If the classifier was bad (>50% error) α is negative classifier is so bad, consider opposite prediction of classifier

70

AdaBoost justification

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * label_i * classifier_k(x_i))$$

Does this look like anything we've seen before?

72

AdaBoost justification

update the example weights

$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * label_i * classifier_k(x_i))$$

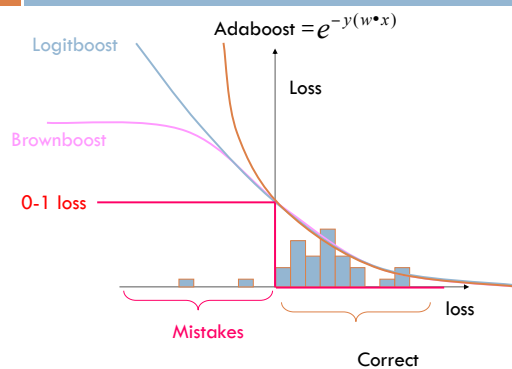
Exponential loss!

$$l(y, y') = \exp(-yy')$$

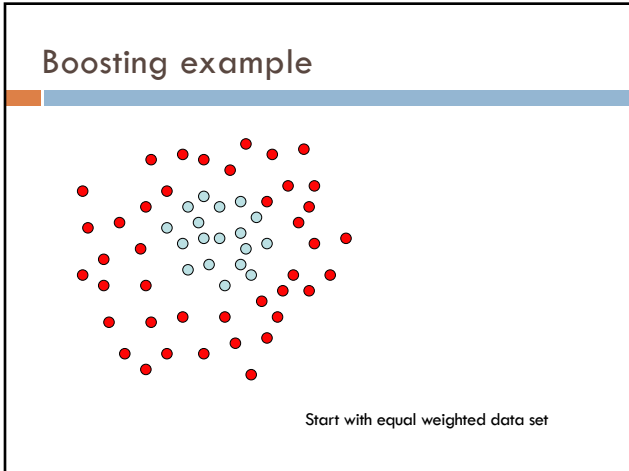
AdaBoost turns out to be another approach for minimizing the exponential loss!

73

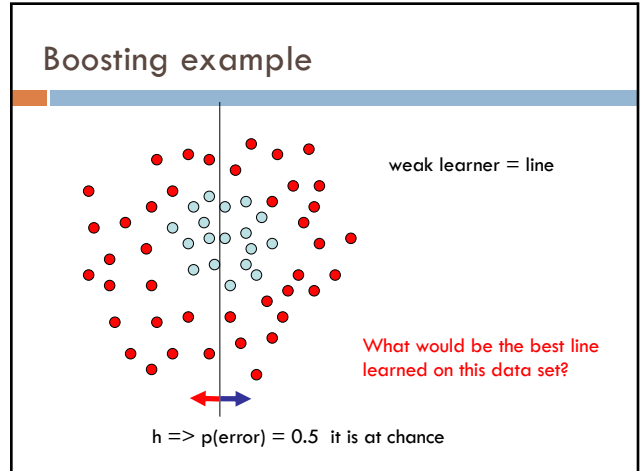
Other boosting variants



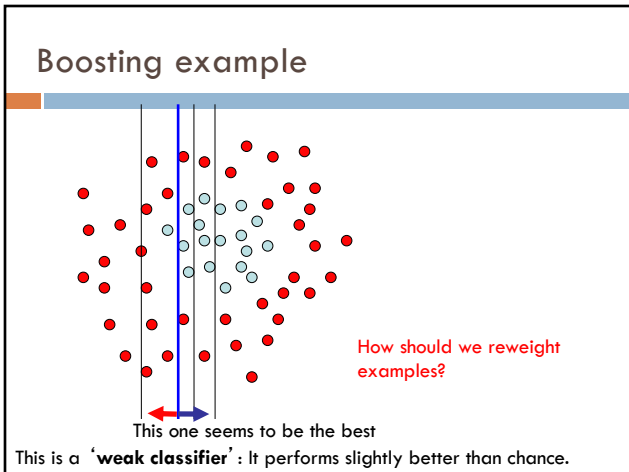
74



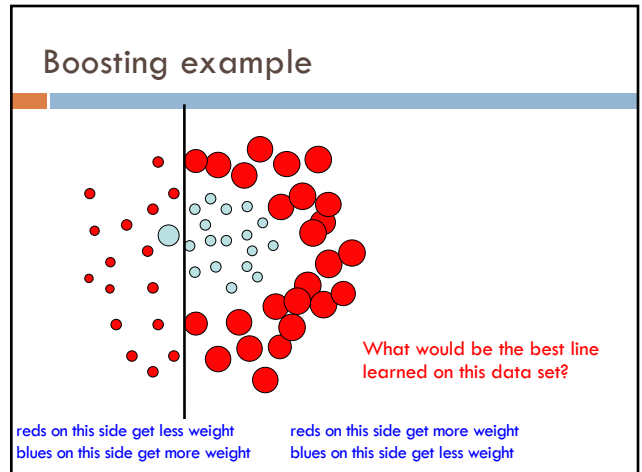
75



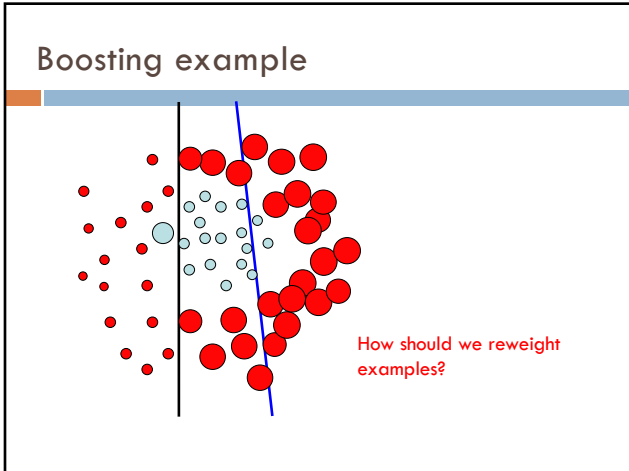
76



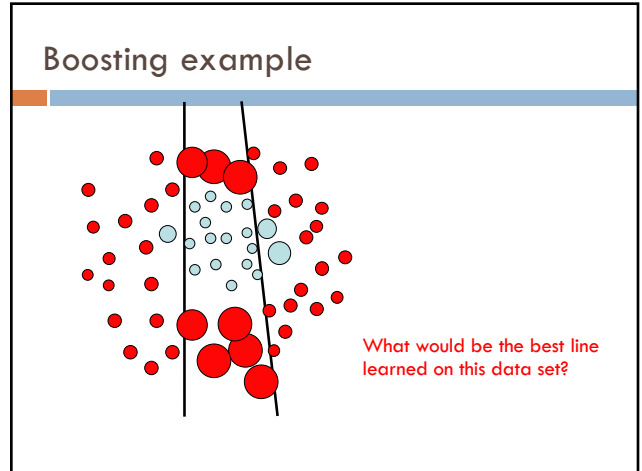
77



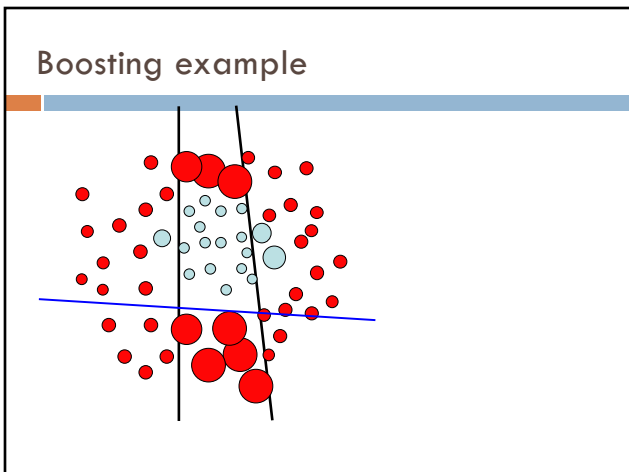
78



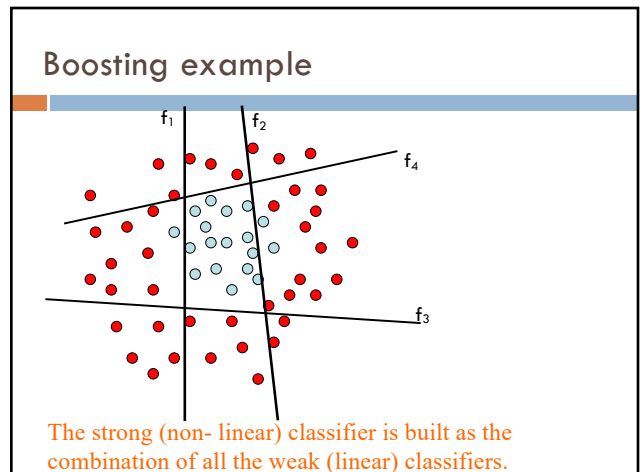
79



80



81



82

AdaBoost: train

for $k = 1$ to *iterations*:

- classifier_k = learn a weak classifier based on weights
- weighted error for this classifier is:
- "score" or weight for this classifier is:
- change the example weights

What can we use as a classifier?

83

AdaBoost: train

for $k = 1$ to *iterations*:

- classifier_k = learn a weak classifier based on weights
- weighted error for this classifier is:
- "score" or weight for this classifier is:
- change the example weights

- Anything that can train on weighted examples
 - For most applications, must be fast!
- Why?

84

AdaBoost: train

for $k = 1$ to *iterations*:

- classifier_k = learn a weak classifier based on weights
- weighted error for this classifier is:
- "score" or weight for this classifier is:
- change the example weights

- Anything that can train on weighted examples
- For most applications, must be fast!
 - Each iteration we have to train a new classifier

85

Boosted decision stumps

One of the most common classifiers to use is a decision tree:

- can use a shallow (2-3 level tree)
- even more common is a 1-level tree
 - called a **decision stump** 😊
 - asks a question about a single feature

What does the decision boundary look like for a decision stump?

86

Boosted decision stumps

One of the most common classifiers to use is a decision tree:

- can use a shallow (2-3 level tree)
- even more common is a 1-level tree
- called a **decision stump** 😊
- asks a question about a single feature

What does the decision boundary look like for boosted decision stumps?

87

Boosted decision stumps

One of the most common classifiers to use is a decision tree:

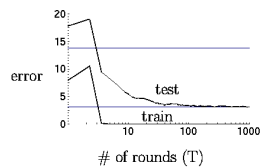
- can use a shallow (2-3 level tree)
- even more common is a 1-level tree
- called a **decision stump** 😊
- asks a question about a single feature
- **Linear classifier!**
- Each stump defines the weight for that dimension
 - If you learn multiple stumps for that dimension then it's the weighted average

88

Boosting in practice

Very successful on a wide range of problems

One of the keys is that boosting tends not to overfit, even for a large number of iterations



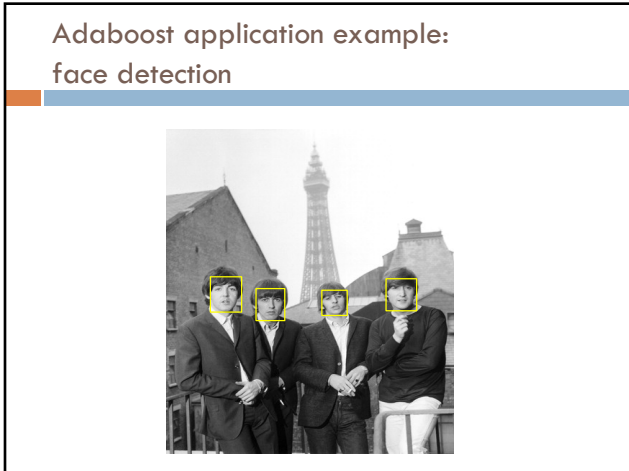
Using **<10,000** training examples can fit **>2,000,000** parameters!

89

Adaboost application example: face detection



90



91

Rapid Object Detection using a Boosted Cascade of Simple Features

Paul Viola
viola@merl.com
Mitsubishi Electric Research Labs
201 Broadway, 8th FL
Cambridge, MA 02139

Michael Jones
mjones@crl.dec.com
Compaq CRL
One Cambridge Center
Cambridge, MA 02142

[Rapid object detection using a boosted cascade of simple features](#)
P. Viola, M. Jones - ... Vision and Pattern Recognition, 2001. CVPR ... 2001 - ieeexplore.ieee.org ... overlap. Each partition yields a single final detection. The ... set. Experiments on a Real-World Test Set We tested our system on the MIT-CMU frontal face test set [1]. This set consists of 130 images with 507 labeled frontal faces. A ...
Cited by 8422 Related articles All 129 versions Cite Save More

92

[PDF] [Rapid object detection using a boosted cascade of simple features](#)
P. Viola, M. Jones - CVPR (1), 2001 - researchgate.net

This paper describes a machine learning approach for visual **object detection** which is capable of processing images extremely **rapidly** and achieving **high detection** rates. This work is distinguished by three key contributions. The first is the introduction of a new image ...
☆ [Cited by 19976](#) Related articles All 101 versions

To give you some context of importance:

The anatomy of a large-scale hypertextual web search engine
S. Brin, L. Page - Computer networks and ISDN systems, 1998 - Elsevier

In this paper, we present Google, a prototype of a **large-scale** search engine which makes heavy use of the structure present in hypertext. Google is designed to crawl and index the Web efficiently and produce much more satisfying search results than existing systems. The ...
☆ [Cited by 18931](#) Related articles All 220 versions Web of Science: 4067

or:

Modeling word burstiness using the Dirichlet distribution
B.E. Madsen, D. Kauchak, C. Elkan - Proceedings of the 22nd international ..., 2005 - dl.acm.org

Multinomial distributions are often used to model text documents. However, they do not capture well the phenomenon that words in a document tend to appear in bursts: if a word appears once, it is more likely to appear again. In this paper, we propose the Dirichlet compound multinomial model (DCM) as an alternative to the multinomial. The DCM model has one additional degree of freedom, which allows it to capture burstiness. We show experimentally that the DCM is substantially better than the multinomial at modeling text ...
☆ [Cited by 289](#) Related articles All 27 versions

93

“weak” learners

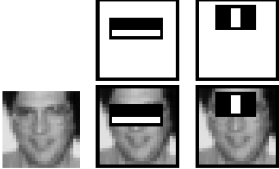
4 Types of “Rectangle filters”
(Similar to Haar wavelets
Papageorgiou, et al.)

Based on 24x24 grid:
160,000 features to choose from

$$g(x) = \text{sum(WhiteArea)} - \text{sum(BlackArea)}$$

94

"weak" learners

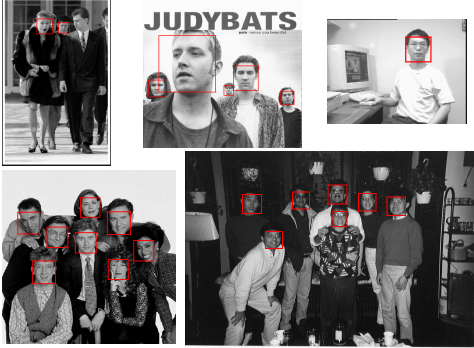


$$F(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x) + \dots$$

$$f_i(x) = \begin{cases} 1 & \text{if } g_i(x) > \theta_i \\ -1 & \text{otherwise} \end{cases}$$


95

Example output

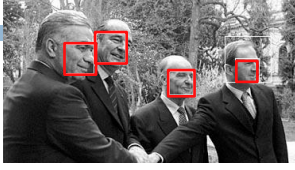


96

Solving other "Face" Tasks

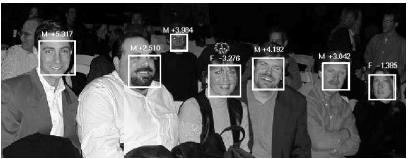


Facial Feature Localization



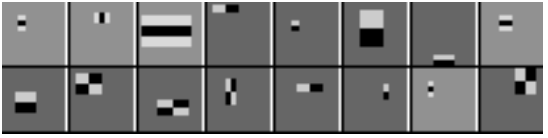

Profile Detection

Demographic Analysis



97

"weak" classifiers learned

98

Bagging vs Boosting

Journal of Artificial Intelligence Research 11 (1999) 169-198 Submitted 1/99; published 8/99

Popular Ensemble Methods: An Empirical Study

David Opitz
 Department of Computer Science
 University of Montana
 Missoula, MT 59812 USA

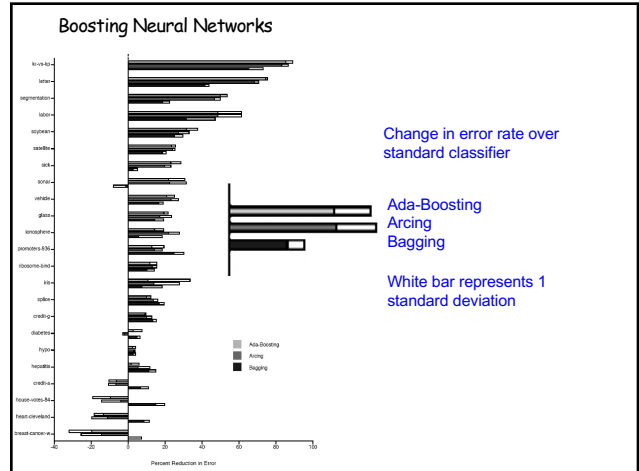
Richard Maclin
 Computer Science Department
 University of Minnesota
 Duluth, MN 55812 USA

OPITZ@CS.UMT.EDU

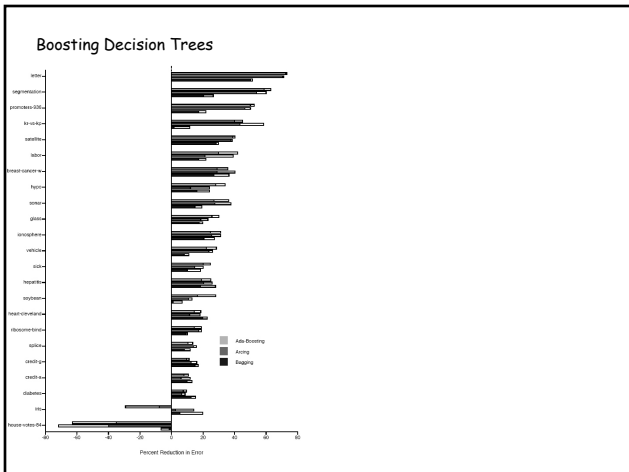
RMACLIN@D.UMN.EDU

<http://arxiv.org/pdf/1106.0257.pdf>

99



100



101