

ENSEMBLE LEARNING

David Kauchak
CS158 – Fall 2016

Admin

Assignment grading

Assignment 9

Midterm 2

Final project

- ▣ No formal class Tuesday: figure out project ideas
- ▣ 11/23 (Wed) submit project proposal

Quick exercise

Write down on the paper (don't write your name):

- 1) Something you're happy about right now
- 2) Something you're worried about right now

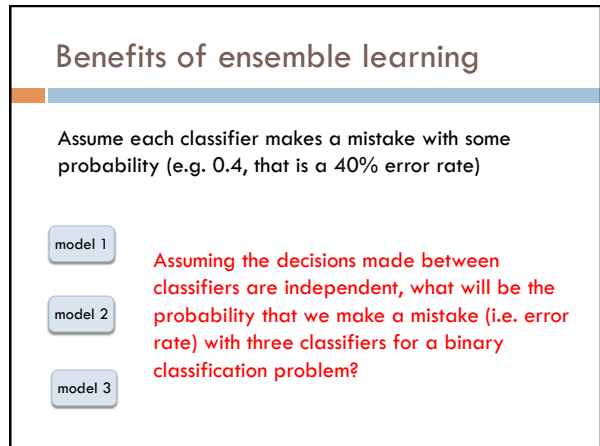
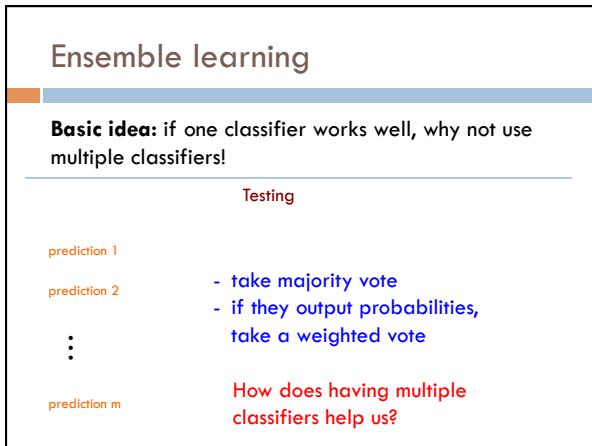
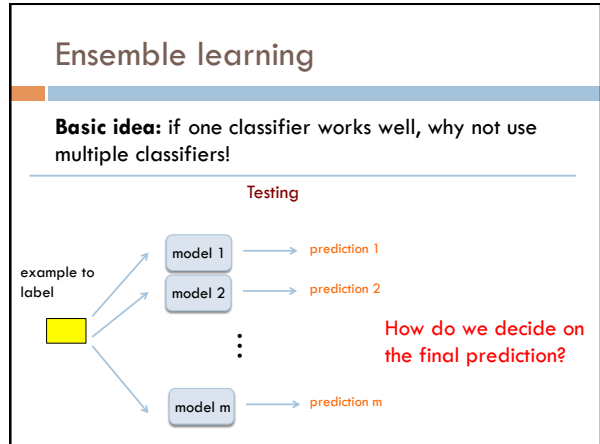
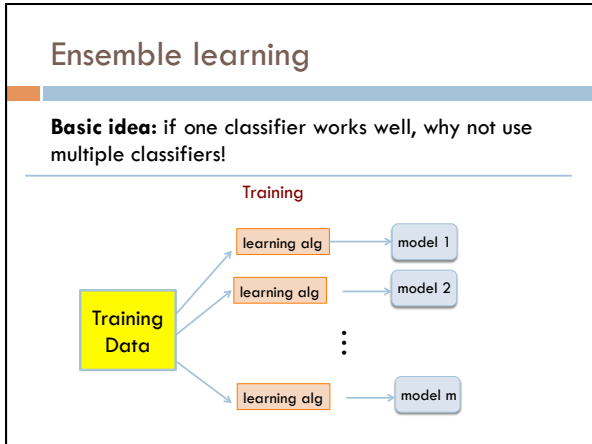
Fold the piece of paper

I'll collect them, redistribute them and we'll read them out loud

If you don't want to participate, just leave the paper blank

Ensemble learning

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



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

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!

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

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

Benefits of ensemble learning

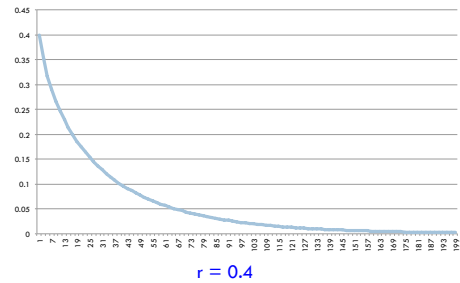
m classifiers in general, for r = probability of mistake for individual classifier:

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

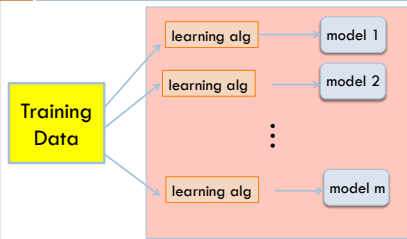
(cumulative probability distribution for the binomial distribution)

Given enough classifiers...

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

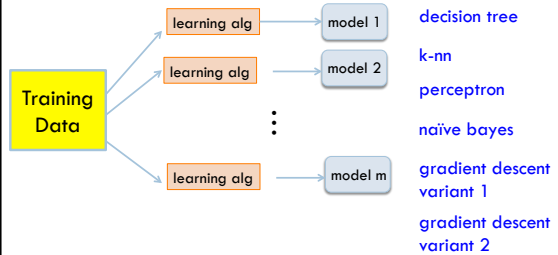


Obtaining independent classifiers



Where do we get m independent classifiers?

Idea 1: different learning methods



Pros/cons?

⋮

Idea 1: different learning methods

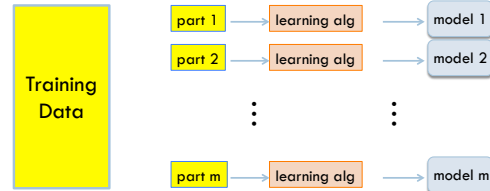
Pros:

- ▣ 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

Idea 2: split up training data



Use the same learning algorithm, but train on different parts of the training data

Idea 2: split up training data

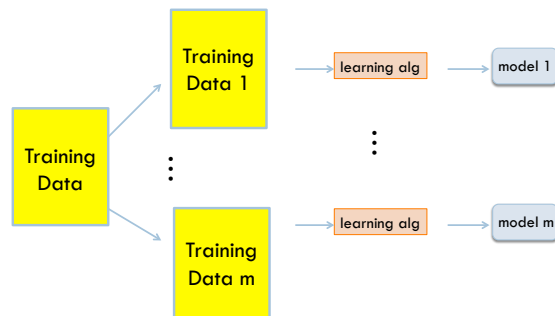
Pros:

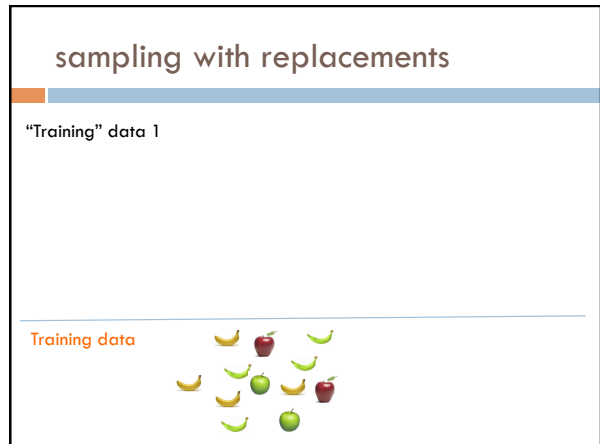
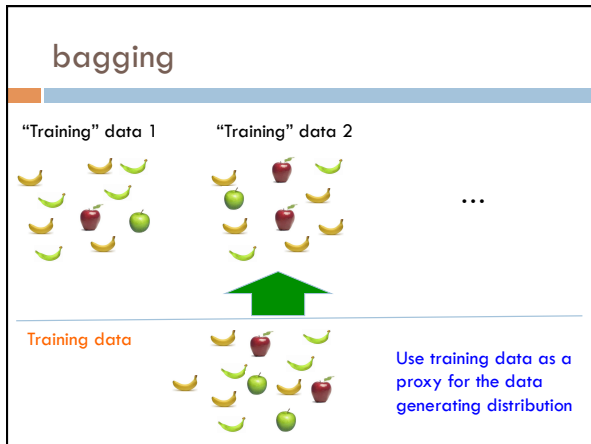
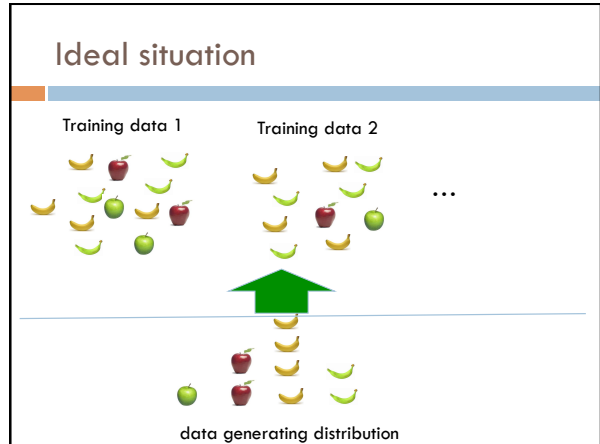
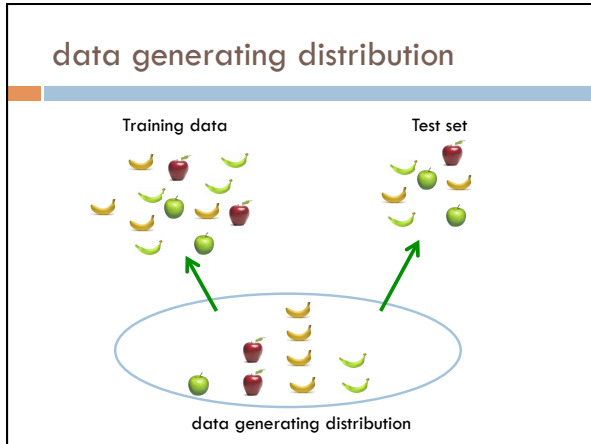
- ▣ Learning from different data, so can't overfit to same examples
- ▣ Easy to implement
- ▣ fast

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

Idea 3: bagging






sampling with replacements

"Training" data 1

pick a random example from the real training data

Training data

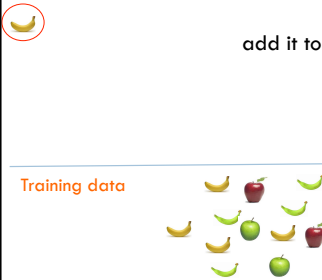


sampling with replacements

"Training" data 1

add it to the new "training" data

Training data




sampling with replacements

"Training" data 1

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

Training data




sampling with replacements

"Training" data 1


pick another random example

Training data



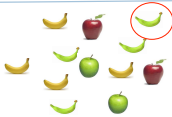
sampling with replacements

"Training" data 1




pick another random example

Training data




sampling with replacements

"Training" data 1



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

Training data



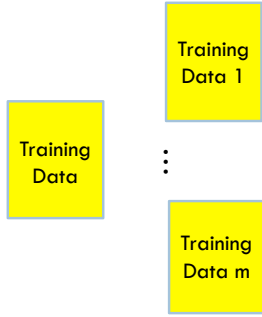
bagging

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

bagging concerns



Won't these all be basically the same?

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

bagging concerns

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

$$1 - 1/n$$

Training data

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

probability of overlap

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

Converges very quickly to $1/e \approx 63\%$

bagging overlap

Training Data

Training Data 1

...

Training Data m

Won't these all be basically the same?

On average, a randomly sampled data set will only contain 63% of the examples in the original

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 a third of the original noisy examples

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

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

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

A green arrow points from the 'training data' column to the 'training data 2' column.

“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 $1-p$

“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!

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

boosting: basic algorithm

Training:
start with equal example weights

for some number of iterations:

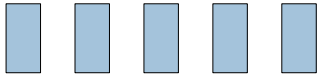
- learn a weak classifier and save
- change the example weights

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)

boosting basics

Start with equal weighted examples

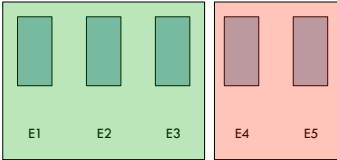
Weights: 

Examples: E1 E2 E3 E4 E5


Learn a weak classifier:  weak 1

Boosting

classified correct classified incorrect

Weights: 

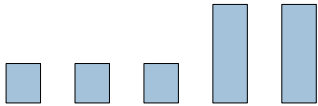
Examples: E1 E2 E3 E4 E5

 weak 1

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

How should we change the example weights?

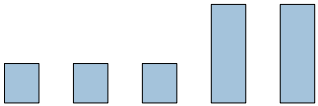
Boosting

Weights: 


Examples: E1 E2 E3 E4 E5

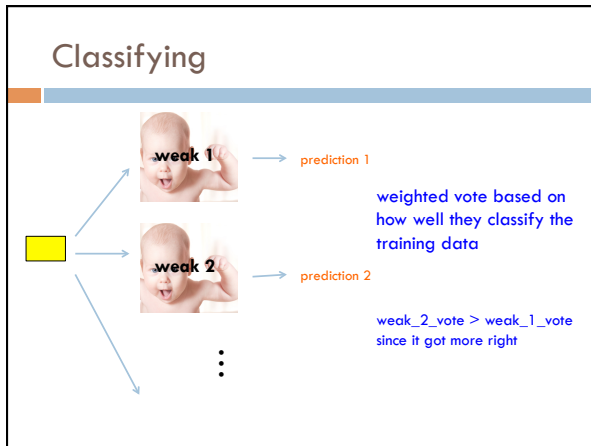
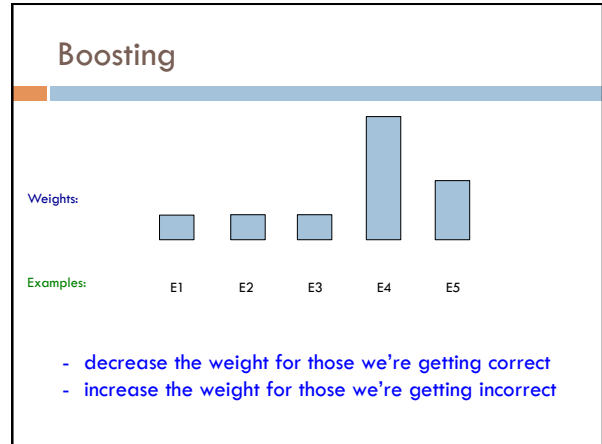
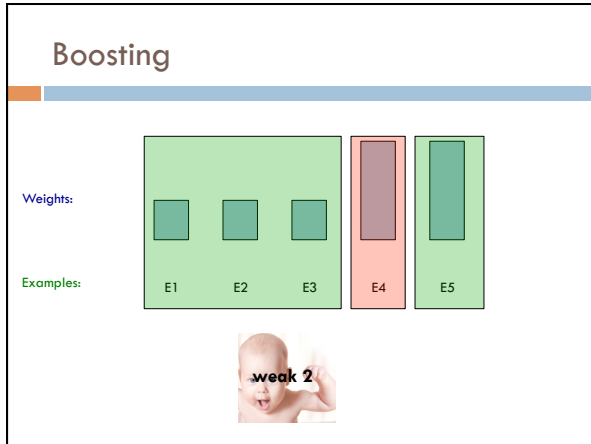
- decrease the weight for those we're getting correct
- increase the weight for those we're getting incorrect

Boosting

Weights: 

Examples: E1 E2 E3 E4 E5

Learn another weak classifier:  weak 2



Notation

x_i example i in the training data

w_i weight for example i , we will enforce:
 $w_i \geq 0$
 $\sum_{i=1}^n w_i = 1$

classifier $_k(x_i)$ +1/-1 prediction of classifier k example i

AdaBoost: train

for $k = 1$ to *iterations*:

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

$$\epsilon_k = \sum_{i=1}^n w_i * I[label_i \neq classifier_k(x_i)]$$

- calculate "score" for this classifier:

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

- change the example weights

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

AdaBoost: train

classifier_k = learn a weak classifier based on weights

weighted error for this classifier is:

$$\epsilon_k = \sum_{i=1}^n w_i * I[label_i \neq classifier_k(x_i)]$$

What does this say?

AdaBoost: train

classifier_k = learn a weak classifier based on weights

weighted error for this classifier is:

$$\epsilon_k = \sum_{i=1}^n w_i * I[label_i \neq \underbrace{classifier_k(x_i)}_{\text{prediction}}]$$

did we get the example wrong

weighted sum of the errors/mistakes

What is the range of possible values?

AdaBoost: train

classifier_k = learn a weak classifier based on weights

weighted error for this classifier is:

$$\epsilon_k = \sum_{i=1}^n w_i * I[label_i \neq \underbrace{classifier_k(x_i)}_{\text{prediction}}]$$

did we get the example wrong

weighted sum of the errors/mistakes

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

AdaBoost: train

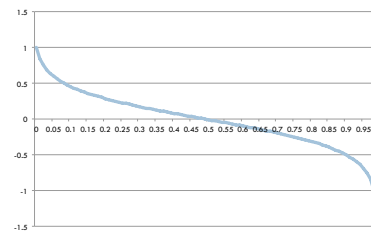
classifier_k = learn a weak classifier based on weights

“score” or weight for this classifier is:

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

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

AdaBoost: train



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

- ranges from $+\infty$ to $-\infty$
- for most reasonable values: ranges from ~ 1 to -1
- errors of 50% = 0

AdaBoost: classify

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

What does this do?

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 $> 50\%$

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 actually vote the opposite!

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?

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):

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

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?

AdaBoost: train

update the example weights

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

correct → positive
incorrect → negative

correct ?
incorrect

AdaBoost: train

update the example weights

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

correct → positive
incorrect → negative

correct → small value
incorrect → large value

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

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?

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?

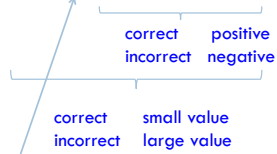
If the classifier was good (<50% error) α is positive:
 trust classifier output and move as normal

If the classifier was back (>50% error) α is negative
 classifier is so bad, consider opposite prediction of classifier

AdaBoost: train

update the example weights

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



If the classifier was good (<50% error) α is positive
 If the classifier was bad (>50% error) α is negative

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?

AdaBoost justification

update the example weights

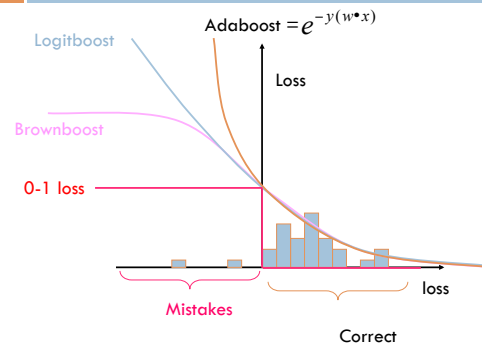
$$w_i = \frac{1}{Z} w_i \exp(-\alpha_k * label_i * classifier_k(x_i))$$

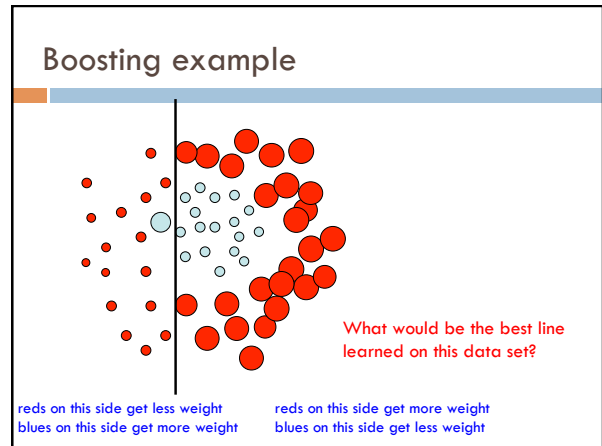
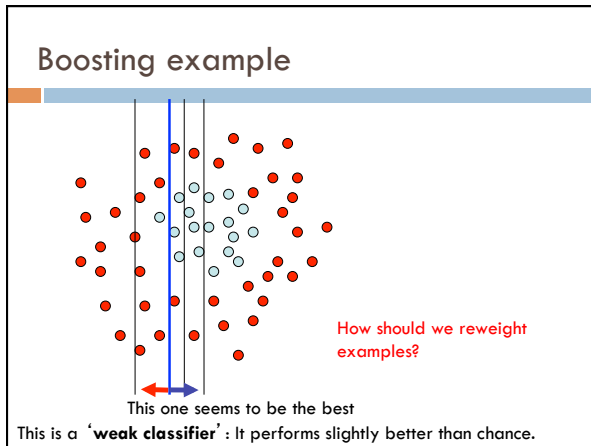
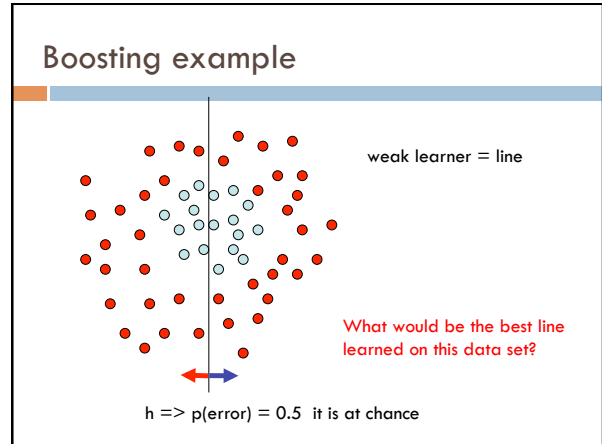
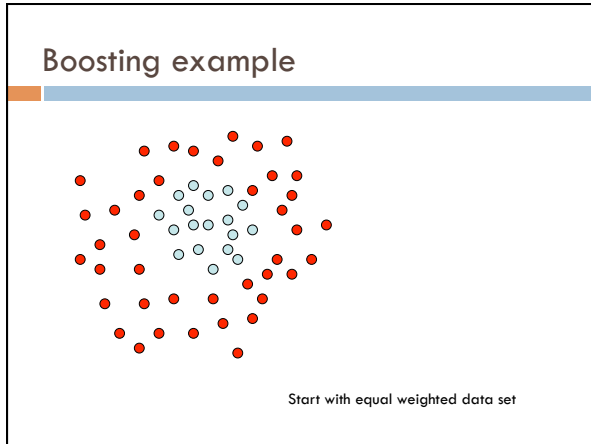
Exponential loss!

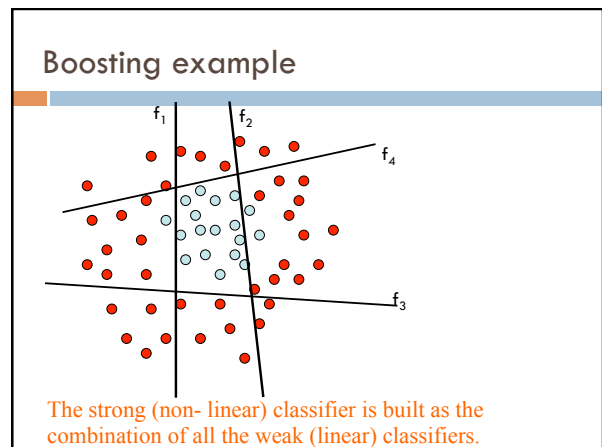
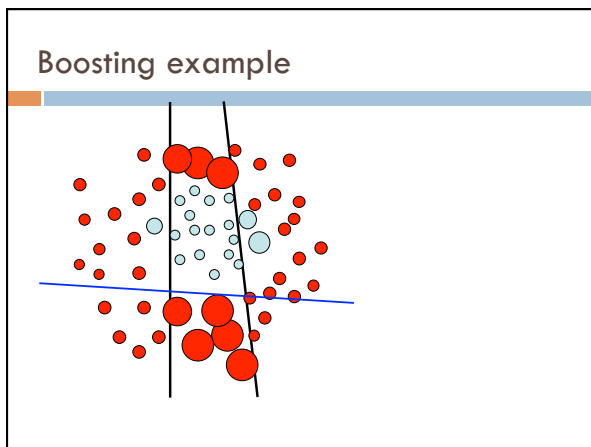
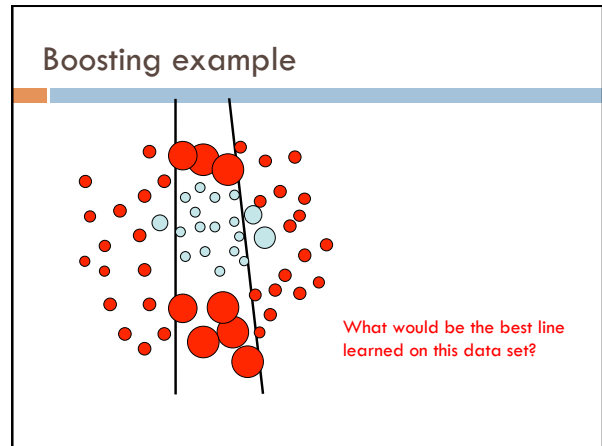
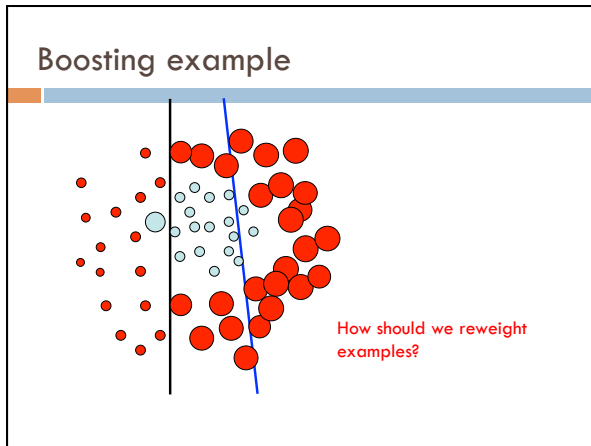
$$l(y, y') = \exp(-y y')$$

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

Other boosting variants







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?

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?

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

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?

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?

Boosted decision stumps

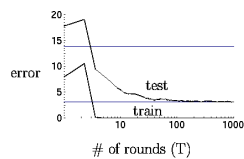
One of the most common classifiers to use is a decision tree:

- 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

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!

Adaboost application example: face detection



Adaboost application example: face detection

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 [11]. This set consists of 130 images with 507 labeled frontal **faces**. A ...
[Cited by 8422](#) [Related articles](#) [All 129 versions](#) [Cite](#) [Save](#) [More](#)

[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 [11]. This set consists of 130 images with 507 labeled frontal **faces**. A ...
[Cited by 8422](#) [Related articles](#) [All 129 versions](#) [Cite](#) [Save](#) [More](#)


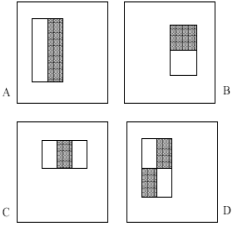
To give you some context of importance:

[The anatomy of a large-scale hypertextual Web search engine](#) 
S. Brin, L. Page - Computer networks and ISDN systems, 1998 - Elsevier
... This is largely because they all have high **PageRank**. ... However, once the system was running smoothly, S. Brin, L. Page/Computer Networks and ISDN Systems 30 ... Google employs a number of techniques to improve search quality including **page rank**, anchor text, and proximity ...
[Cited by 11070](#) [Related articles](#) [All 349 versions](#) [Cite](#) [Save](#)

or:

[Modeling word burstiness using the Dirichlet distribution](#)
RE Madsen, D. Kauchak, C. Elkan - Proceedings of the 22nd International ..., 2005 - dl.acm.org
Abstract 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 ...
[Cited by 148](#) [Related articles](#) [All 34 versions](#) [Cite](#) [Save](#)

“weak” learners

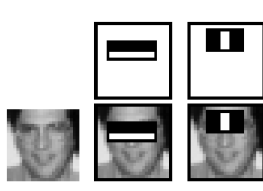



4 Types of “Rectangle filters” (Similar to Haar wavelet Papageorgiou, et al.)

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

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

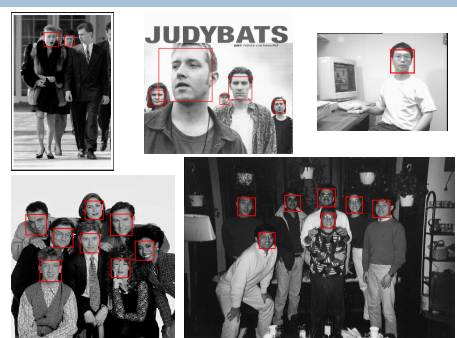
"weak" learners




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

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

Example output




Solving other "Face" Tasks



Facial Feature Localization Profile Detection

Demographic Analysis

"weak" classifiers learned



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>

