

Admin

Assignment 9

Midterm 2

Final project proposals due today!

2

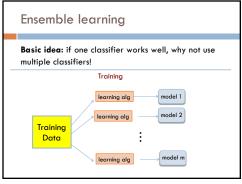
4

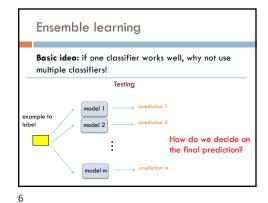
Write down on the paper (don't write your name): | Something you're happy about right now | Something you're worried about right now | Fold the piece of paper | I'll collect them, redistribute them and we'll read them out loud If you don't want to participate, just leave the paper blank

3

Ensemble learning

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





Ensemble learning

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

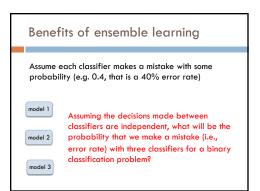
Testing

prediction 1

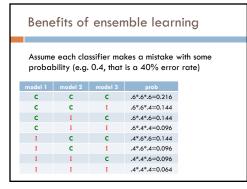
prediction 2

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

How does having multiple classifiers help?

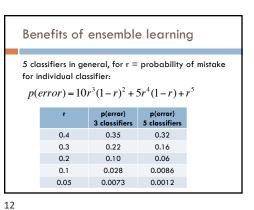


8



Bene	efits of	ensen	nble learr	ning
			kes a mistake v is a 40% erroi	
model 1	model 2	model 3	prob	
С	С	С	.6*.6*.6=0.216	
С	С	I	.6*.6*.4=0.144	0.096+ 0.096+ 0.096+ 0.064 = 35% error!
С	I	С	.6*.4*.6=0.144	
С	I	I	.6*.4*.4=0.096	
I	С	С	.4*.6*.6=0.144	
I	С	I	.4*.6*.4=0.096	
	I	С	.4*.4*.6=0.096	
1				

Benefits of ensemble learning 3 classifiers in general, for $\mathbf{r} = \mathbf{probability}$ of mistake for individual classifier: $p(error) = 3r^2(1-r) + r^3$ binomial distribution $\begin{array}{c|cccc} \mathbf{r} & \mathbf{p(error)} \\ 0.4 & 0.35 \\ 0.3 & 0.22 \\ 0.2 & 0.10 \\ 0.1 & 0.028 \\ 0.05 & 0.0073 \\ \end{array}$



Benefits of ensemble learning

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

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

(cumulative probability distribution for the binomial distribution)

Given enough classifiers... $p(error) = \sum_{i=(m\pi)/2}^{\infty} {m \choose i} p^i (1-r)^{m-i} \qquad r = 0.4$

13

14

What's 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

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?

Assuming the decisions made between

model 3

15

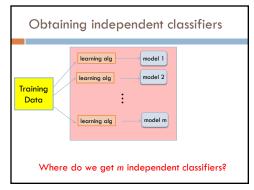
What's the catch?

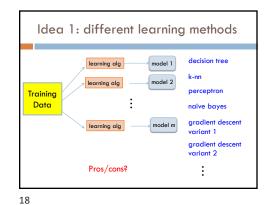
Assume each classifier makes a mistake with some probability (e.g. 0.4, that is a 40% error rate)

model 1

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?

model 3





19

Idea 1: different learning methods

Pros:

Lots of existing classifiers already
Can work well for some problems

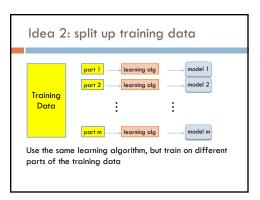
Cons/concerns:

Often, classifiers are not independent, that is, they

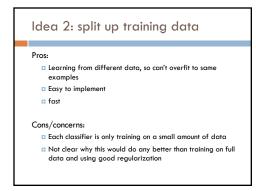
voting won't help us if they're making the same mistakes

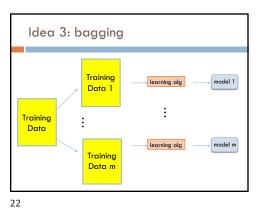
e.g. many of these classifiers are linear models

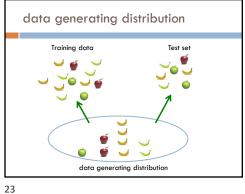
make the same mistakes!

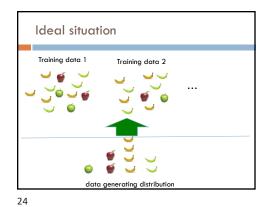


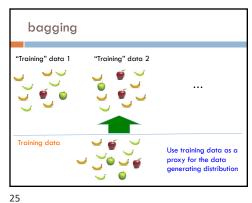
20







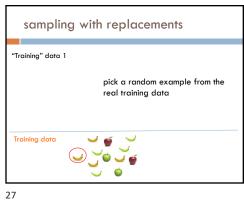




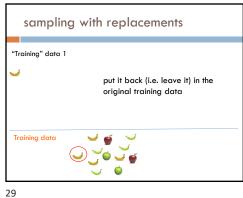
sampling with replacements "Training" data 1 Training data

26

28

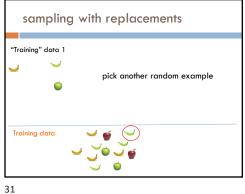


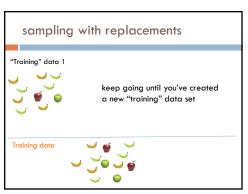
sampling with replacements "Training" data 1 add it to the new "training" data Training data



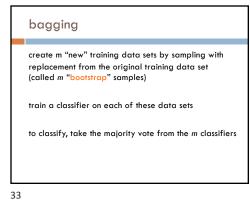
sampling with replacements "Training" data 1 pick another random example Training data

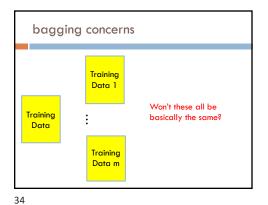
30

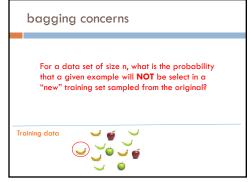


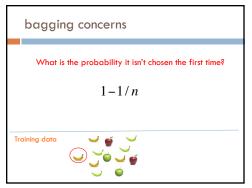


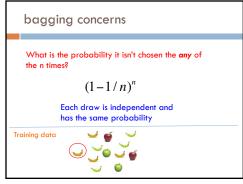
32











probability of overlap $\frac{(1-1/n)^n}{\sqrt{1-1/n}}$ Converges very quickly to 1-1/e \approx 37%

38

40

37

39

Training Data 1

Training Data 1

Training Data 1

Training Data

Training Data m

Training Data m

Training Data m

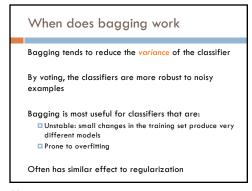
When does bagging work

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?

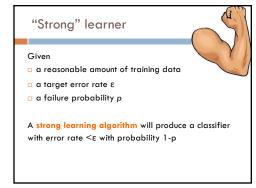
They'll still have ~10% of the examples as noisy
However, these examples will only represent about two-thirds of the original noisy examples

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



42

41



"Weak" learner

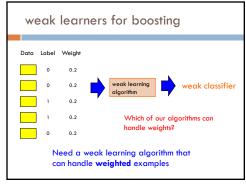
Given

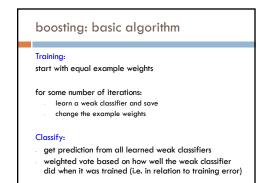
a reasonable amount of training data
a failure probability p

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

Weak learners are much easier to create!

43 44





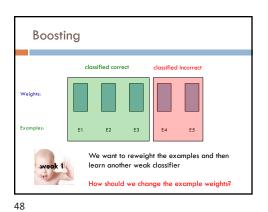
boosting basics

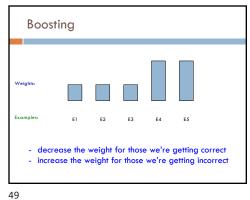
Start with equal weighted examples

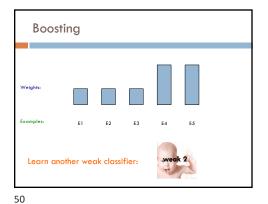
Weights:

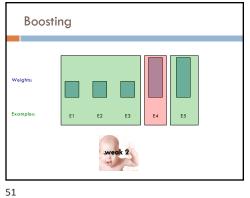
Examples:

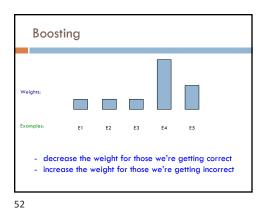
Exa

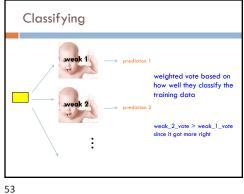












Notation example i in the training data weight for example i, we will enforce: $w_i \geq 0$ $\sum_{i=1}^{n} w_i = 1$ classifier $k(x_i)$ +1/-1 prediction of classifier k example i

AdaBoost: train

for k = 1 to iterations:

55

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

$$\varepsilon_k = \sum\nolimits_{i=1}^n w_i *1[label_i \neq classifier_k(x_i)]$$

calculate "score" for this classifier:
$$\alpha_k = \frac{1}{2} \log \left(\frac{1-\varepsilon_i}{\varepsilon_i} \right)$$
 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:

$$\varepsilon_k = \sum\nolimits_{i=1}^n w_i * 1[label_i \neq classifier_k(x_i)]$$

What does this say?

56

AdaBoost: train

classifier $_k$ = learn a weak classifier based on weights

weighted error for this classifier is: $\varepsilon_k = \sum_{i=1}^n w_i * 1[label_i \neq classifier_k(x_i)]$ What is the range of possible values?

did we get the example wrong

weighted sum of the errors/mistakes

AdaBoost: train

classifier $_k$ = learn a weak classifier based on weights

weighted error for this classifier is: $\varepsilon_k = \sum_{i=1}^n w_i * 1[label_i \neq classifier_k(x_i)]$ Between 0 (if we get all examples right) and 1 (if we get them all wrong)

weighted sum of the errors/mistakes

57

59

58

60

AdaBoost: train $\text{classifier}_k = \text{learn a weak classifier based on weights}$ "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
- error < 50% = positive error > 50% = negative

AdaBoost: classify

 $classify(x) = sign\left(\sum_{k=1}^{literations} \alpha_k * classifier_k(x)\right)$

What does this do?

AdaBoost: classify

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

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%

61

63

62

AdaBoost: classify

 $classify(x) = sign\left(\sum_{k=1}^{tterations} \alpha_k * classifier_k(x)\right)$

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

We vote the opposite!

AdaBoost: train, updating the weights

update the example weights

$$w_i = \frac{1}{7} w_i \exp(-label_i * \alpha_k * classifier_k(x_i))$$

Remember, we want to enforce:

$$w_i \ge 0$$

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

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

What should it be?

AdaBoost: train $w_{l} = \frac{1}{Z}w_{l} \exp(-label_{l}*\alpha_{k}*classifier_{k}(x_{l}))$ Remember, we want to enforce: $w_{l} \geq 0$ $\sum_{i=1}^{n}w_{i}=1$ normalizing constant (i.e. the sum of the "new" w_l): $Z = \sum_{i=1}^{n}w_{i} \exp(-\alpha_{k}*label_{i}*classifier_{k}(x_{i}))$

AdaBoost: train update the example weights $w_i = \frac{1}{Z} w_i \exp(-\frac{label_i * \alpha_k * classifier_k(x_i)}{})$ What does this do?

65

67

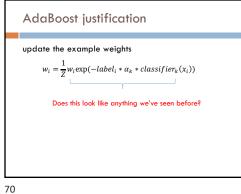
AdaBoost: train

update the example weights $w_i = \frac{1}{Z}w_i \exp(-label_i * \alpha_k * classifier_k(x_i))$ correct \rightarrow positive incorrect \rightarrow negative

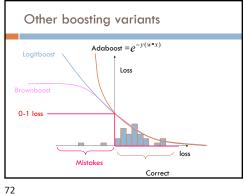
correct $\stackrel{?}{\underset{incorrect}{2}}$

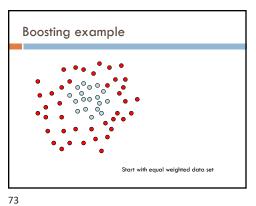
AdaBoost: train $w_i = \frac{1}{Z} w_i \exp(-label_i * \alpha_k * classifier_k(x_i))$ $\text{correct} \rightarrow \text{positive}$ $\text{incorrect} \rightarrow \text{negative}$ $\text{correct} \rightarrow \text{small value}$ $\text{incorrect} \rightarrow \text{large value}$ Note: only change weights based on current classifier (not all previous classifiers)

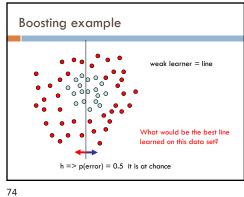
68

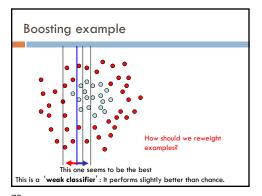


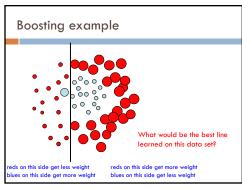
AdaBoost justification update the example weights $w_i = \frac{1}{7} w_i \exp(-label_i * \alpha_k * classifier_k(x_i))$ Exponential loss! $l(y, y') = \exp(-yy')$ AdaBoost turns out to be another approach for minimizing the exponential loss!

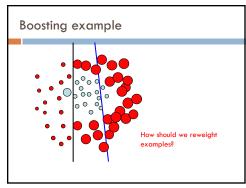


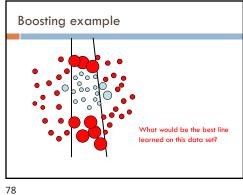


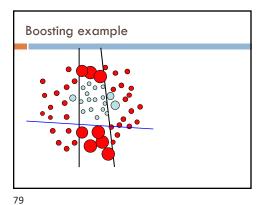


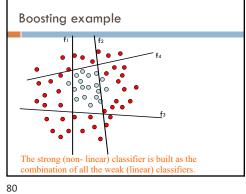












AdaBoost: train for k = 1 to iterations: $\mathsf{classifier}_k \equiv \mathsf{learn} \ \mathsf{a} \ \mathsf{weak} \ \mathsf{classifier} \ \mathsf{based} \ \mathsf{on} \ \mathsf{weights}$ weighted error for this classifier is: "score" or weight for this classifier is: change the example weights What can we use as a classifier?

81

AdaBoost: train

for k = 1 to iterations:

82

84

- $classifier_k = learn a weak classifier based on weights$
- 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

- 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?

85

83

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!

86

- Each stump defines the weight for that dimension
 - If you learn multiple stumps for that dimension then it's the weighted average

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

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

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!

87

Adaboost application example: face detection

