# BOOSTING

David Kauchak
CS451 – Fall 2013

---

## Admin

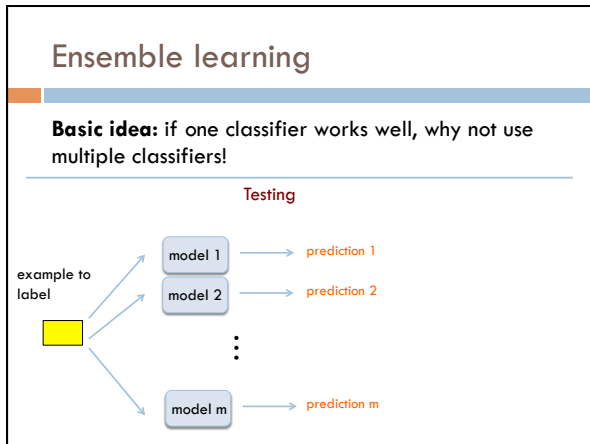Final project

---

## 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!

Training



Training Data → learning alg → model 1

learning alg → model 2

⋮

learning alg → model m

## Ensemble learning

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

Testing

example to label

model 1 → prediction 1

model 2 → prediction 2

⋮

model m → prediction m

## 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 |

## "Strong" learner

Given

☐ a reasonable amount of training data

☐ a target error rate $\varepsilon$

☐ a failure probability $p$

A **strong learning algorithm** will produce a classifier with error rate $< \varepsilon$ 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

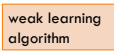| Data | Label | Weight |
|------|-------|--------|
|  | 0 | 0.2 |
|  | 0 | 0.2 |
|  | 1 | 0.2 |
|  | 1 | 0.2 |
|  | 0 | 0.2 |

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

## 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?

11/24/13

# 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**

# Boosting

Weights:

Examples:  E1  E2  E3  E4  E5

**weak 2**
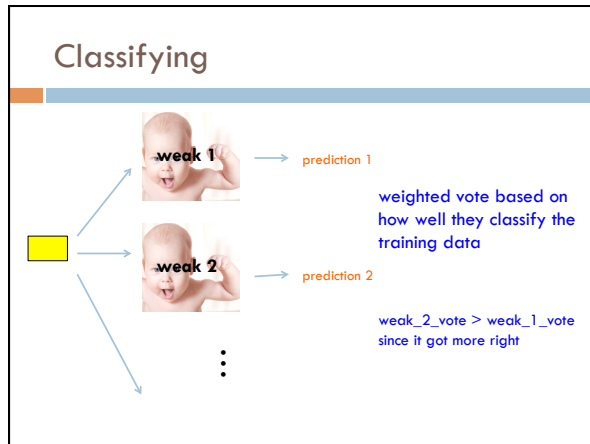
# Boosting

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

4

## Classifying



weak 1 → prediction 1

weighted vote based on how well they classify the training data

weak 2 → prediction 2

weak_2_vote > weak_1_vote since it got more right

## 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

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

- calculate "score" for this classifier:

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

- change the example weights

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

## 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 does this say?

## 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?

prediction

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

prediction

did we get the example wrong

weighted sum of the errors/mistakes

## 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 1 to -1
- errors of 50% = 0

## AdaBoost: classify

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

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 $\alpha$ (remember $\alpha$ varies from 1 to -1 training error)

What happens if a classifier has error >50%

## 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 $\alpha$ (remember $\alpha$ varies from 1 to -1 training error)

We actually vote the opposite!

## AdaBoost: train, updating the weights

update the example weights

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

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

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

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" $w_i$):

$$Z = \sum_{i=1}^{n} w_i \exp\left(-\alpha_k * label_i * classifier_k(x_i)\right)$$

## AdaBoost: train

update the example weights

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

What does this do?

## AdaBoost: train

update the example weights

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

correct       positive
incorrect    negative

correct       ?
incorrect

## AdaBoost: train

update the example weights

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

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

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

What does the $\alpha$ do?

## AdaBoost: train

update the example weights

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

What does the $\alpha$ do?

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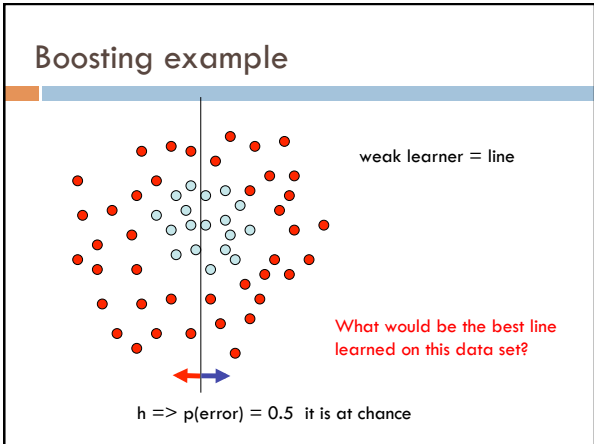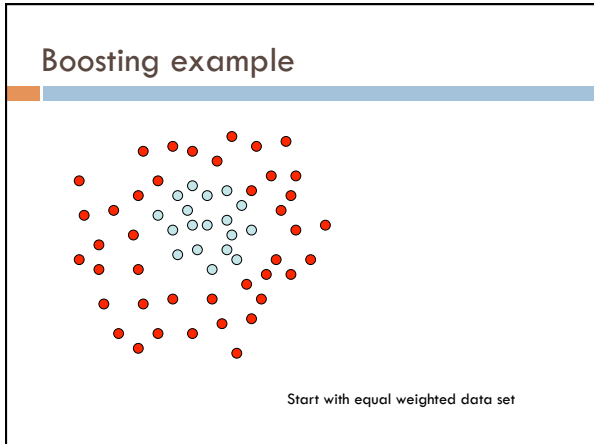    classifier is so bad, consider opposite prediction of
    classifier

## AdaBoost: train

update the example weights

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

correct       positive
incorrect    negative

correct        small value
incorrect     large value

If the classifier was good (<50% error) $\alpha$ is positive
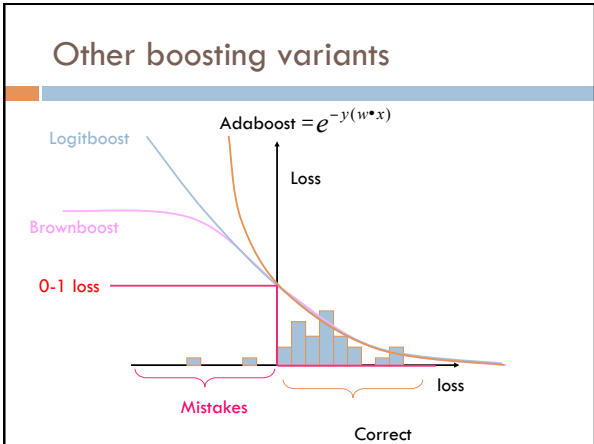If the classifier was back (>50% error) $\alpha$ is negative

## AdaBoost justification

update the example weights

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

Does this look like anything we've seen before?

9
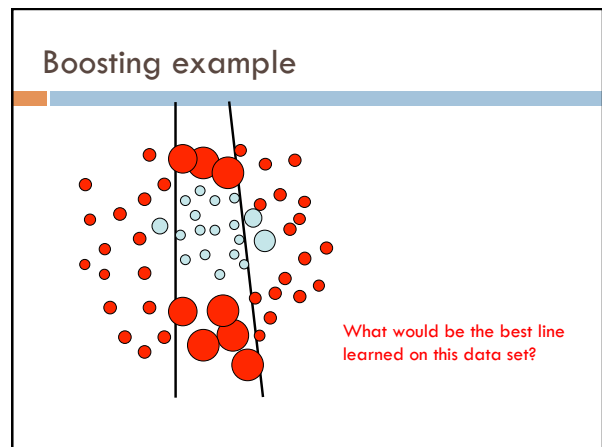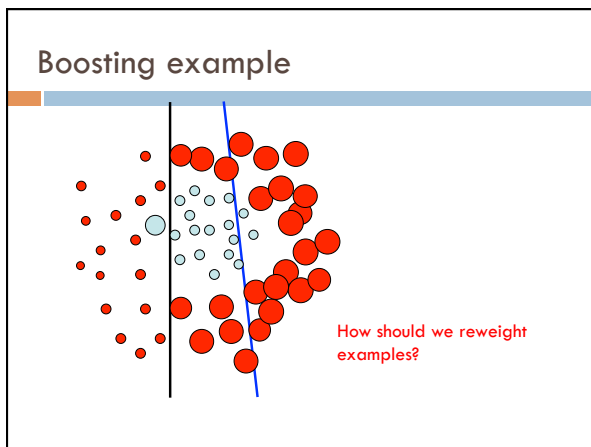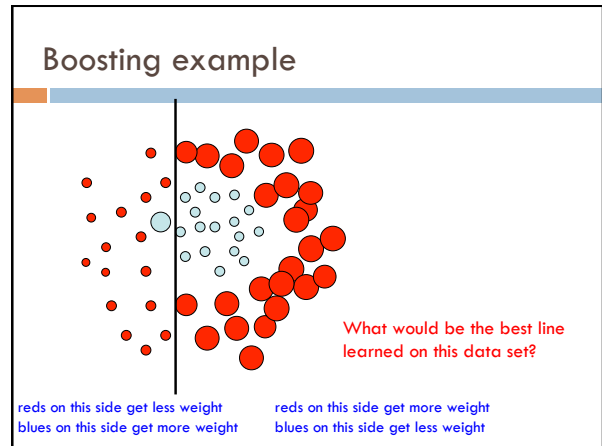
## AdaBoost justification
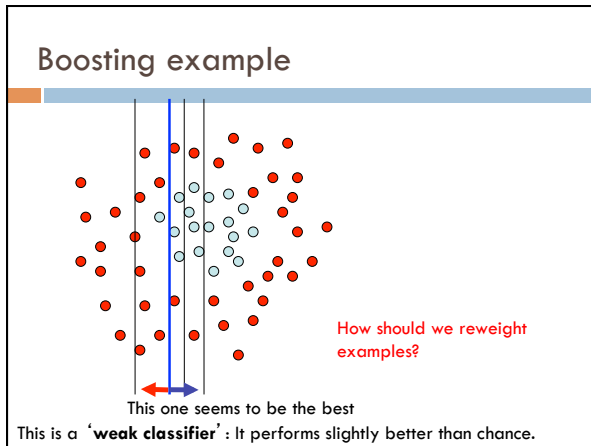
update the example weights

$$w_i = \frac{1}{Z} w_i \exp\left(-\alpha_k * label_i * classifier_k(x_i)\right)$$
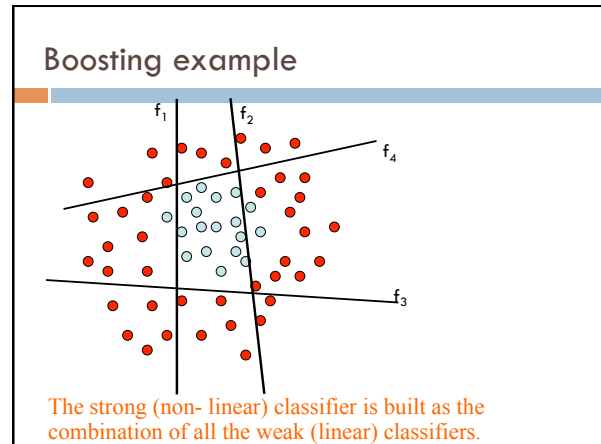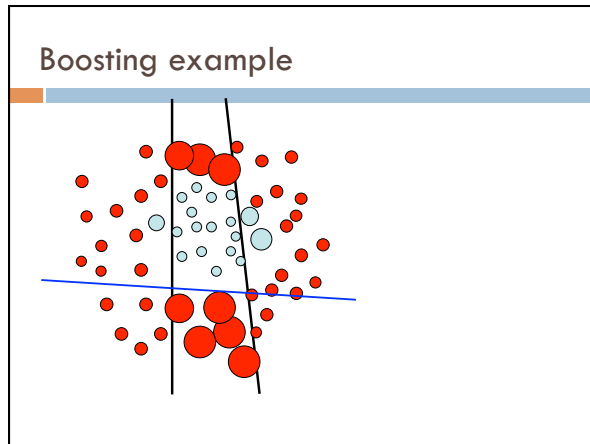
Exponential loss!

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

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

## Other boosting variants

Adaboost $= e^{-y(w \bullet x)}$

Logitboost

Loss

Brownboost

0-1 loss

loss

Mistakes

Correct

## Boosting example

Start with equal weighted data set

## Boosting example

weak learner = line

What would be the best line learned on this data set?

h => p(error) = 0.5  it is at chance

## Boosting example

How should we reweight examples?

This one seems to be the best

This is a 'weak classifier': It performs slightly better than chance.

## Boosting example

What would be the best line learned on this data set?

reds on this side get less weight
blues on this side get more weight

reds on this side get more weight
blues on this side get less weight

## Boosting example

How should we reweight examples?

## Boosting example

What would be the best line learned on this data set?

## Boosting example



## Boosting example



The strong (non- linear) classifier is built as the combination of all the weak (linear) classifiers.

## 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

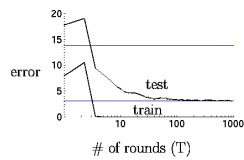One of the most common classifiers to use is a decision tree:

- 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
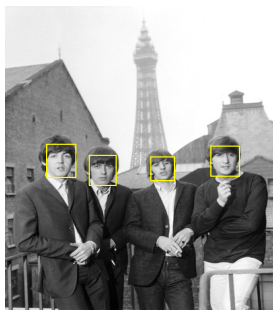


Using <10,000 training examples can fit >2,000,000 parameters!

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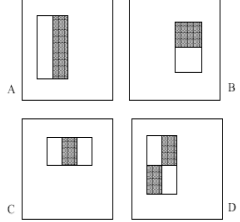Cited by 8422   Related articles   All 129 versions   Cite   Save   More ▾

## Slide 1

To give you some context of importance:
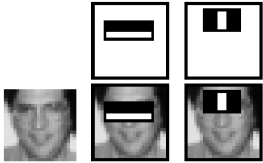
Google

or:

## Slide 2

# "weak" learners

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

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

A   B
C   D
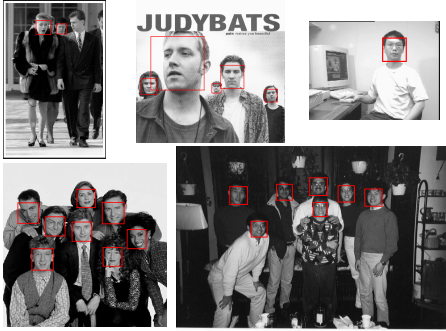
$$g(x) = sum(WhiteArea) - sum(BlackArea)$$

## Slide 3

# "weak" learners

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

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

## Slide 4

# Example output

JUDYBATS

## 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**                                                          OPITZ@CS.UMT.EDU
*Department of Computer Science*
*University of Montana*
*Missoula, MT 59812 USA*

**Richard Maclin**                                                       RMACLIN@D.UMN.EDU
*Computer Science Department*
*University of Minnesota*
*Duluth, MN 55812 USA*

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

Boosting Neural Networks



Change in error rate over standard classifier

Ada-Boosting
Arcing
Bagging

White bar represents 1 standard deviation

Boosting Decision Trees