

# An Empirical Comparison of Voting Classification Algorithms: Bagging, Boosting, and Variants

---

**Eric Bauer**

Computer Science Department  
Stanford University

**Ron Kohavi**

Data Mining and Visualization  
Silicon Graphics Inc.

Presented by:  
Nick Hamatake  
Karan Singh



# Motivation

---

- i Methods for voting classification algorithms show much promise
- i Review algorithms and present a large empirical study comparing several variants
- i Improve understanding of why and when these algorithms work
- i Explore practical problems



# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions



# Outline

---

- i **Introduction**
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions



# Introduction

---

- i Two types of voting algorithms
  - | Perturb and Combine
  - | Adaptively Resample and Combine
- i Authors describe a large empirical study to understand why and when algorithms affect classification error



# Outline

---

- i Introduction
- i **The base inducers**
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions



# The base inducers

---

- i 4 base inducers, implemented in *MLC++*
- i Decision trees
  - | MC4 (MLC++ C4.5)
  - | MC4(1) (decision stump)
  - | MC4(1)-disc (entropy discretization)
- i Naïve-Bayes



# Outline

---

- i Introduction
- i The base inducers
- i **The voting algorithms**
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions





# Bagging (Brieman 1996)

---

**Input:** training set  $S$ , Inducer  $\mathcal{I}$ , integer  $T$  (number of bootstrap samples).

1. for  $i = 1$  to  $T$  {
2.      $S' =$  bootstrap sample from  $S$  (i.i.d. sample with replacement).
3.      $C_i = \mathcal{I}(S')$
4. }
5.  $C^*(x) = \arg \max_{y \in Y} \sum_{i: C_i(x)=y} 1$  (the most often predicted label  $y$ )

**Output:** classifier  $C^*$ .



# AdaBoost.M1 (Freund & Schapire 1995)

---

**Input:** training set  $S$  of size  $m$ , Inducer  $\mathcal{I}$ , integer  $T$  (number of trials).

1.  $S' = S$  with instance weights assigned to be 1.
2. For  $i = 1$  to  $T$  {
3.      $C_i = \mathcal{I}(S')$
4.      $\epsilon_i = \frac{1}{m} \sum_{x_j \in S' : C_i(x_j) \neq y_j} \text{weight}(x_j)$      (weighted error on training set).
5.     If  $\epsilon_i > 1/2$ , set  $S'$  to a bootstrap sample from  $S$  with weight 1 for every instance and goto step 3 (this step is limited to 25 times after which we exit the loop).
6.      $\beta_i = \epsilon_i / (1 - \epsilon_i)$
7.     For-each  $x_j \in S'$ , if  $C_i(x_j) = y_j$  then  $\text{weight}(x_j) = \text{weight}(x_j) \cdot \beta_i$ .
8.     Normalize the weights of instances so the total weight of  $S'$  is  $m$ .
9. }
10.  $C^*(x) = \arg \max_{y \in Y} \sum_{i: C_i(x) = y} \log \frac{1}{\beta_i}$

**Output:** classifier  $C^*$ .



## Arc-x4 (Breiman 1996)

---

**Input:** training set  $S$  of size  $m$ , Inducer  $\mathcal{I}$ , integer  $T$  (number of trials).

1. For  $i = 1$  to  $T$  {
2.      $S' = S$  with instance weight for  $x$  set to  $(1 + e(x)^4)$  (normalized so total weight is  $m$ ), where  $e(x)$  is the number of misclassifications made on  $x$  by classifiers  $C_1$  to  $C_{i-1}$ .
3.      $C_i = \mathcal{I}(S')$
4. }
5.  $C^*(x) = \arg \max_{y \in Y} \sum_{i: C_i(x)=y} 1$  (the most often predicted label  $y$ )

**Output:** classifier  $C^*$ .



# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i **Bias/Variance decomposition**
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions



# Bias/Variance decomposition

---

**Intrinsic “target noise” ( $\sigma^2$ )** This quantity is a lower bound on the expected error of any learning algorithm. It is the expected error of the Bayes-optimal classifier.

**Squared “bias” ( $\text{bias}^2$ )** This quantity measures how closely the learning algorithm’s average guess (over all possible training sets of the given training set size) matches the target.

**“Variance” ( $\text{variance}$ )** This quantity measures how much the learning algorithm’s guess fluctuates for the different training sets of the given size.



# Bias/Variance decomposition

---

$$\text{Error} = \sum_x P(x) (\sigma_x^2 + \text{bias}_x^2 + \text{variance}_x)$$

where

$$\sigma_x^2 \equiv \frac{1}{2} \left( 1 - \sum_{y \in Y} P(Y_F = y|x)^2 \right)$$

$$\text{bias}_x^2 \equiv \frac{1}{2} \sum_{y \in Y} [P(Y_F = y|x) - P(Y_H = y|x)]^2$$

$$\text{variance}_x \equiv \frac{1}{2} \left( 1 - \sum_{y \in Y} P(Y_H = y|x)^2 \right) .$$



# Bias/Variance decomposition

---

- i Kohavi and Wolpert (1996)
  - | two-stage sampling procedure
    - 1) test set is split from training set
    - 2) remaining data,  $D$ , is sampled repeatedly to estimate bias/variance on the test set
  - | whole process can be repeated multiple times to improve estimates



# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i **Experimental design**
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions





# Experimental design

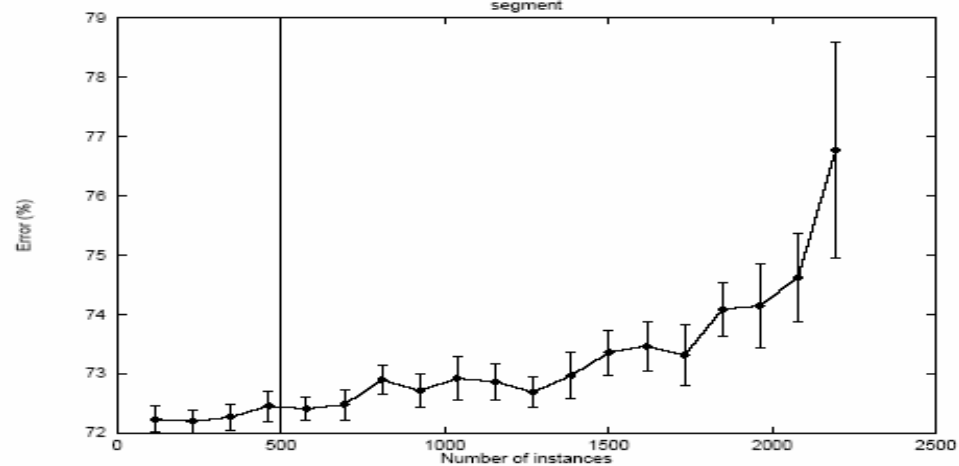
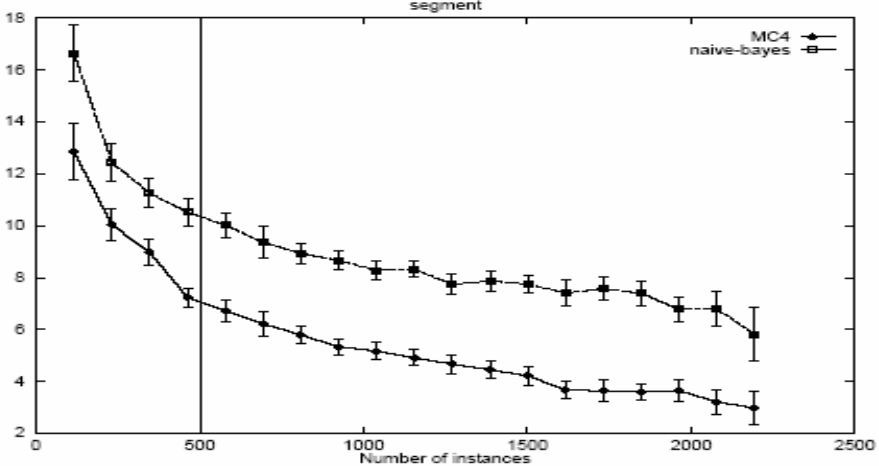
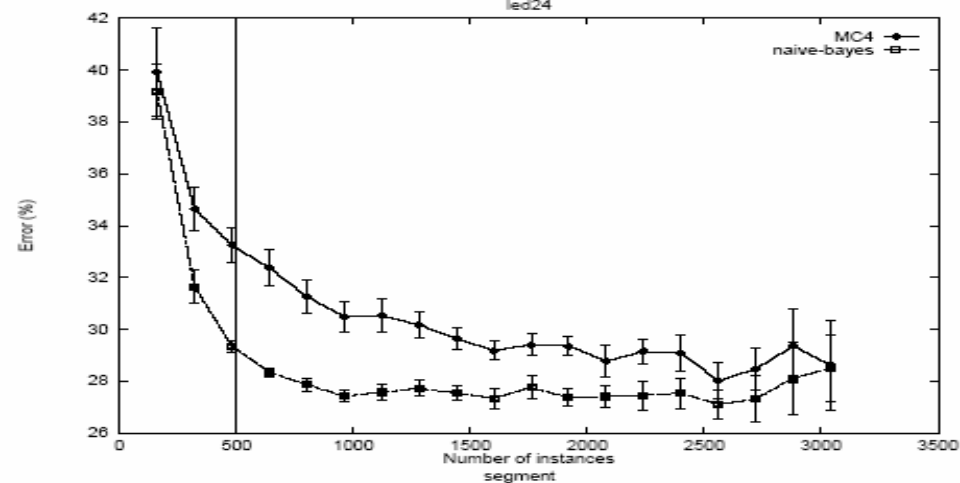
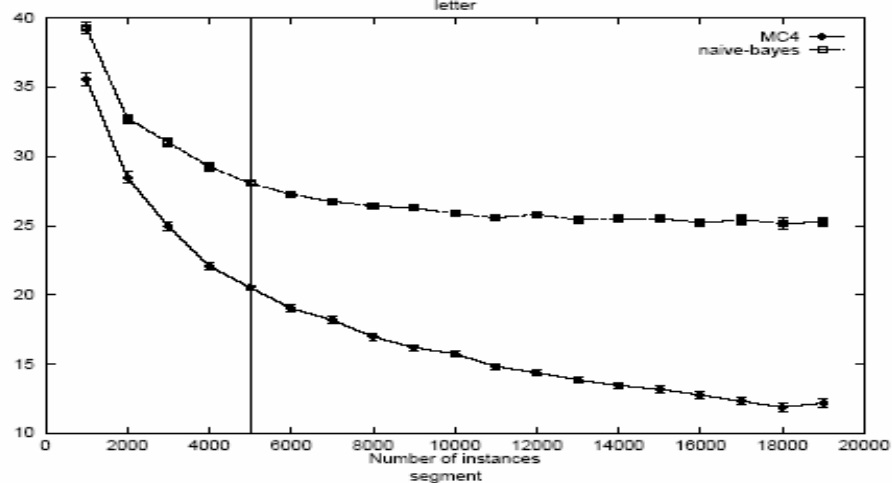
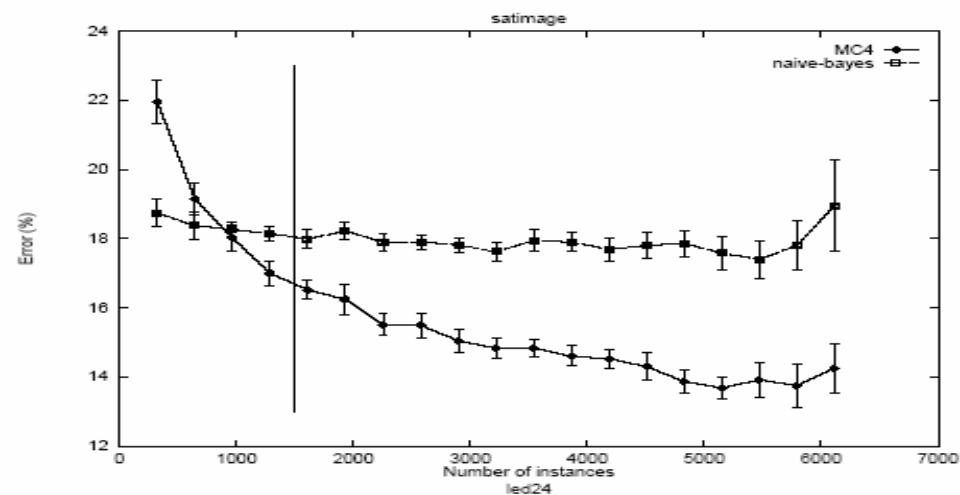
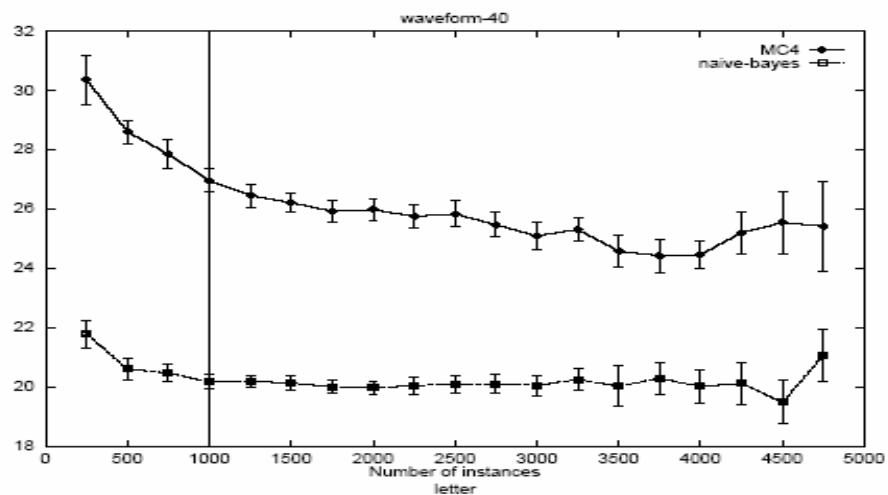
---

- i all data sets have at least 1000 instances
- i training set size selected at a point where error was still decreasing
- i at least half of the data used for test set
- i only 25 sub-classifiers voted for each algorithm
- i authors verified their implementations of voting algorithms against original papers



# Experimental design: datasets and training set sizes used

Data set	Dataset size	Training set size	Attributes Cont/nominal	Classes
Credit (German)	1,000	300	7/13	2
Image Segmentation (segment)	2,310	500	19/0	7
Hypothyroid	3,163	1,000	7/18	2
Sick-euthyroid	3,163	800	7/18	2
DNA	3,186	500	0/60	3
Chess	3,196	500	0/36	2
LED-24	3,200	500	0/24	10
Waveform-40	5,000	1,000	40/0	3
Satellite image (satimage)	6,435	1,500	36/0	7
Mushroom	8,124	1,000	0/22	2
Nursery	12,960	3,000	0/8	5
Letter	20,000	5,000	16/0	26
Adult	48,842	11,000	6/8	2
Shuttle	58,000	5,000	9/0	7

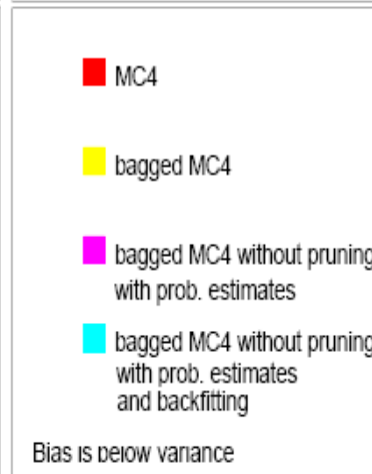
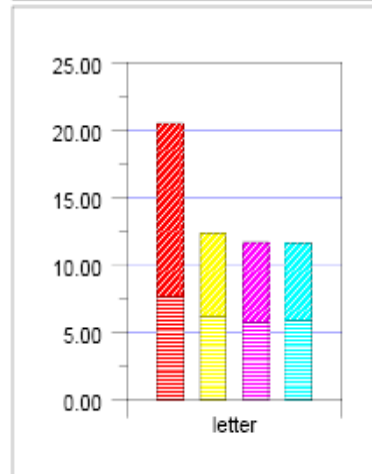
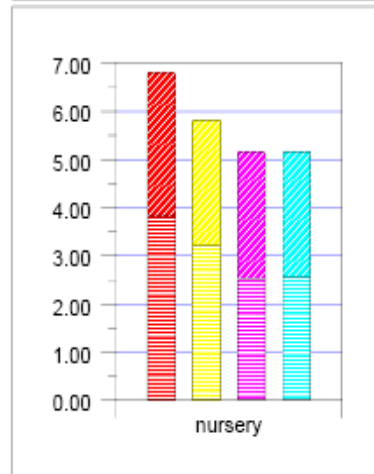
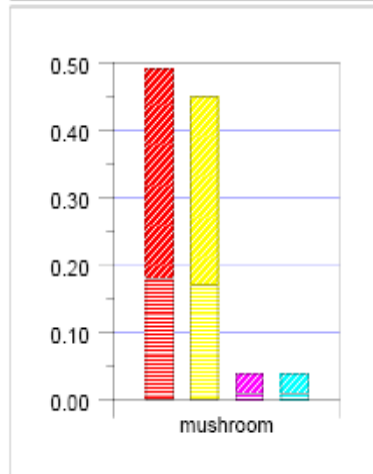
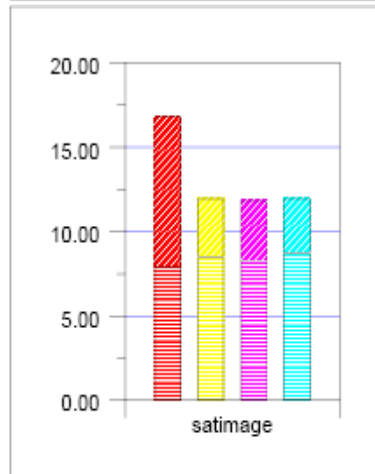
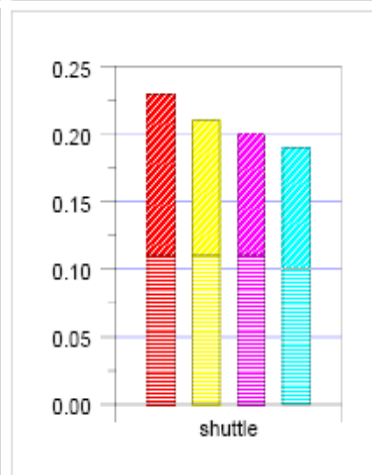
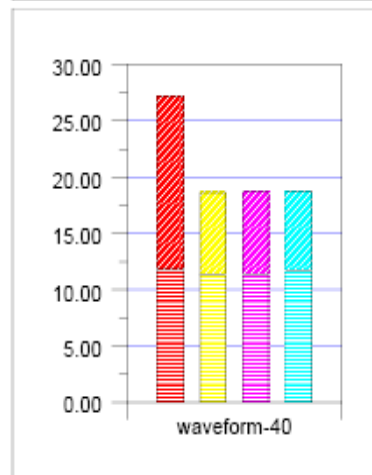
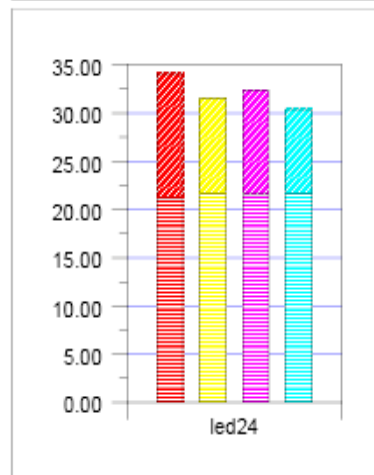
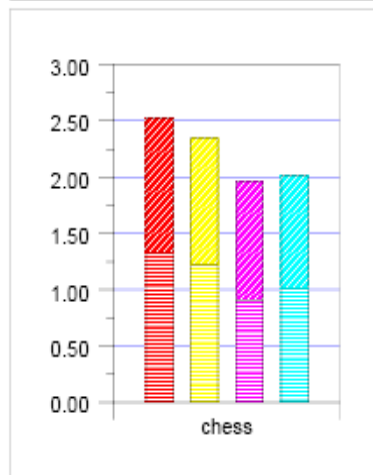
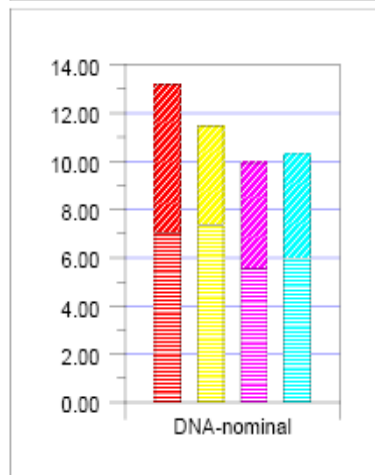
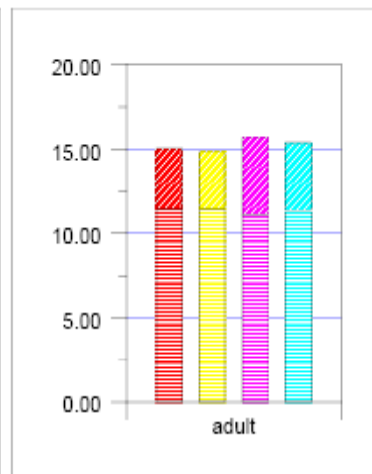
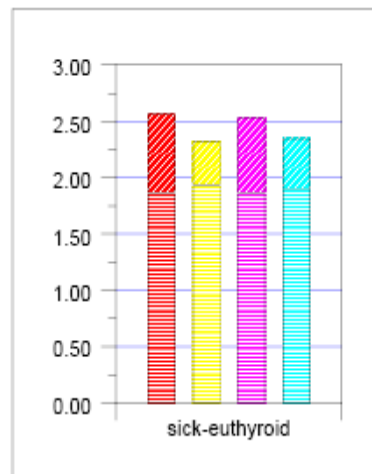
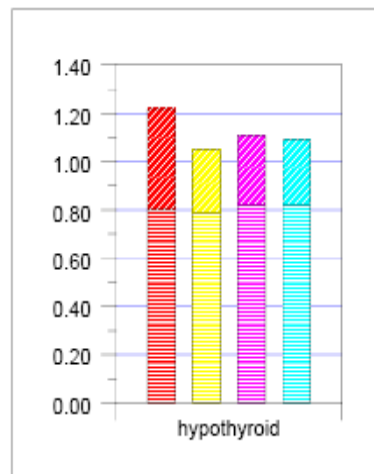
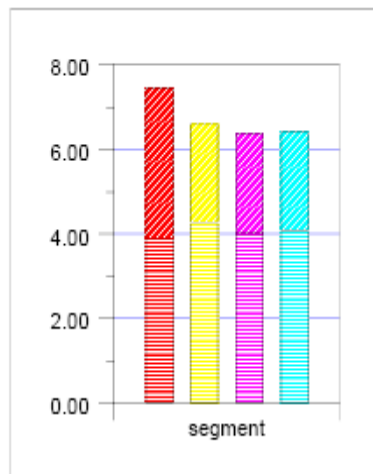
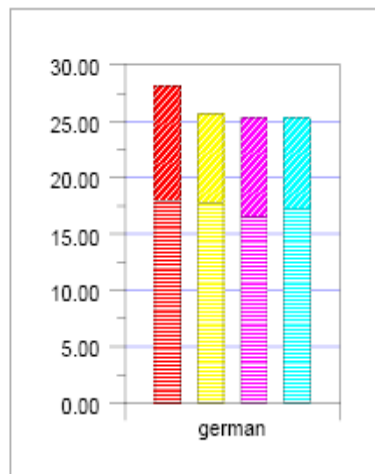




# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i **Bagging and variants**
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i Conclusions

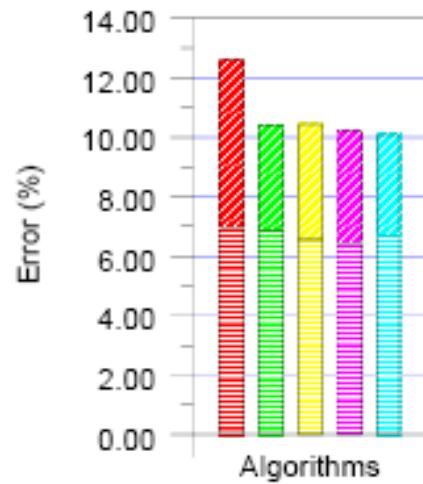




# Bagging Results

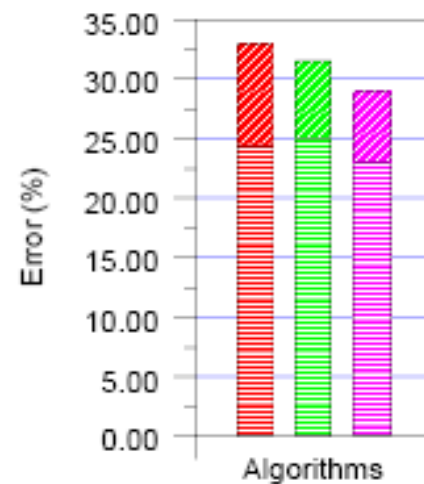
---

- i Performance boost across the board
- i Most gain due to reduced variance (but some bias)
- i Biggest effect on many-attr sets
- i Little effect on large training sets
- i Increased tree size (198 to 240)
- i Naïve Bayes - little to no change
  - | 3% average relative decrease in error
  - | Already highly-stable technique



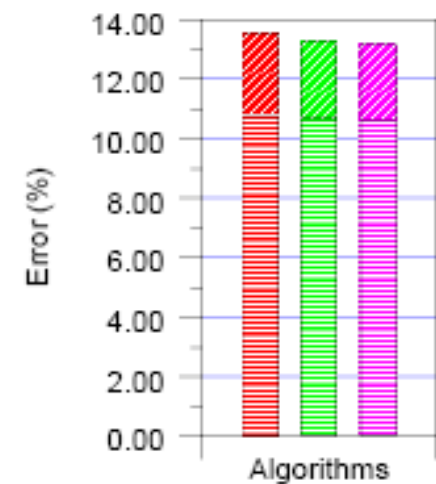
■ MC4  
■ Bagged MC4  
■ Above without pruning  
■ Above with p-Bagging  
■ Above with backfitting

Bias is below variance



■ MC4(1)-disc  
■ Bagged MC4(1)-disc  
■ p-Bagged MC4(1)-disc

Bias is below variance



■ Naive-Bayes  
■ Bagged Naive-Bayes  
■ p-Bagged Naive-Bayes

Bias is below variance



## Why the increased tree size?

---

- i Theory: Less pruning due to duplicate data
- i Test: Bagged trees actually *smaller (25%) before pruning*
- i *Idea: Skip pruning*
  - | *Pruning dampens variance while increasing bias*
  - | *Bagging does well with high variance*
  - | *Results: 11% increase in rel. error for un-bagged trees; slight decrease in error for bagged*





# Bagging Probabilities (p-Bagging)

---

- i Motivation: MC4 can generate probabilities; why not use them?
- i Skipping pruning process should help predict prob.'s even better (why?)
- i Worked well
  - | 2%, 4%, 8% rel. reduction in rel. error for MC4, MC4(1), MC4(1)-disc.
  - | No significant change for Naive-Bayes (when alg. tuned to predict prob.'s)



# Wagging and Backfitting

---

- i Wagging - Weight Aggregation
  - | Different method of perturbing train set -- noise vs. re-sampling
  - | Tunable: higher variance in noise leads to reduced variance (higher bias)
  - | Performs on-par with bagging when tuned optimally
- i Backfitting - Use whole training set to calculate leaf probabilities
  - | Avg. rel. var. decrease of 11% (3% total error); improvement (or no change) in *all models*



# Bagging: Summary

---

- i Three new variants: disabling pruning, p-Bagging, backfitting
- i Each variant improved performance; 4% rel. decrease in error at the end
- i Surprising reduce in bias for stumps
- i Naïve Bayes not helped much
- i MSE improved dramatically with p-Bagging



# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i **Boosting: AdaBoost and Arc-x4**
- i Future work
- i Conclusions



# Boosting: AdaBoost and Arc-x4

---

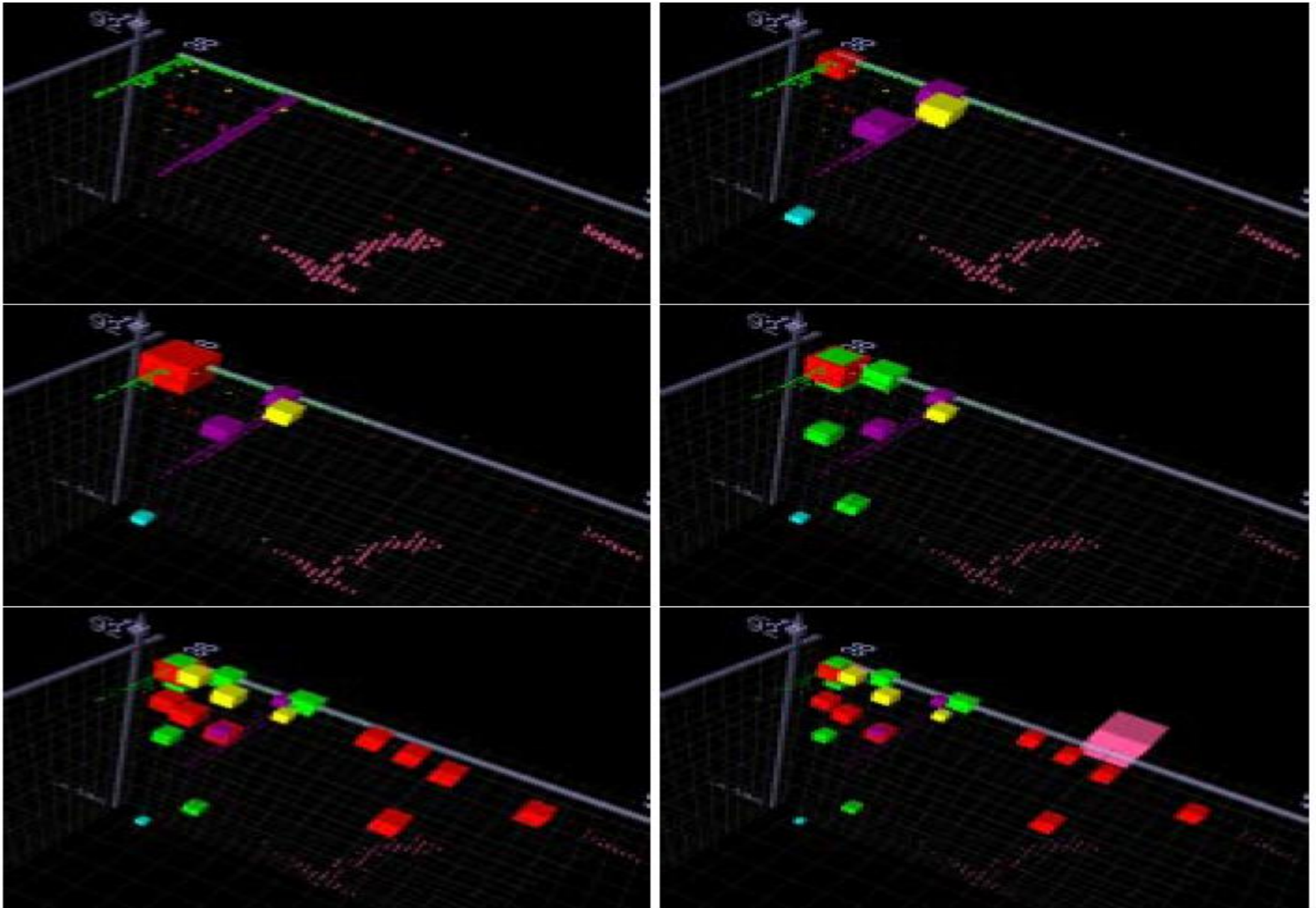
- i but first, an example:
  - | to get a better understanding of the process of boosting
  - | highlight the issue of numerical instabilities and underflows



## Example

---

- i "A training set of size 5000 was used with the shuttle dataset. The MC4 algorithm already has relatively small test error of 0.38%."
- i Update rule:  
For each  $x_j$ , divide  $\text{weight}(x_j)$  by  $2\epsilon_i$  if  $C_j(x_j) \neq y_j$  and  $2(1-\epsilon_i)$  otherwise



*Figure 8.* The shuttle dataset projected on the three most discriminatory axes. Color/shading denotes the class of instances and the cube sizes correspond to the instance weights. Each picture shows one AdaBoost trial, where progression is from the top left to top right, then middle left, etc.



## Example: lesson learned ...

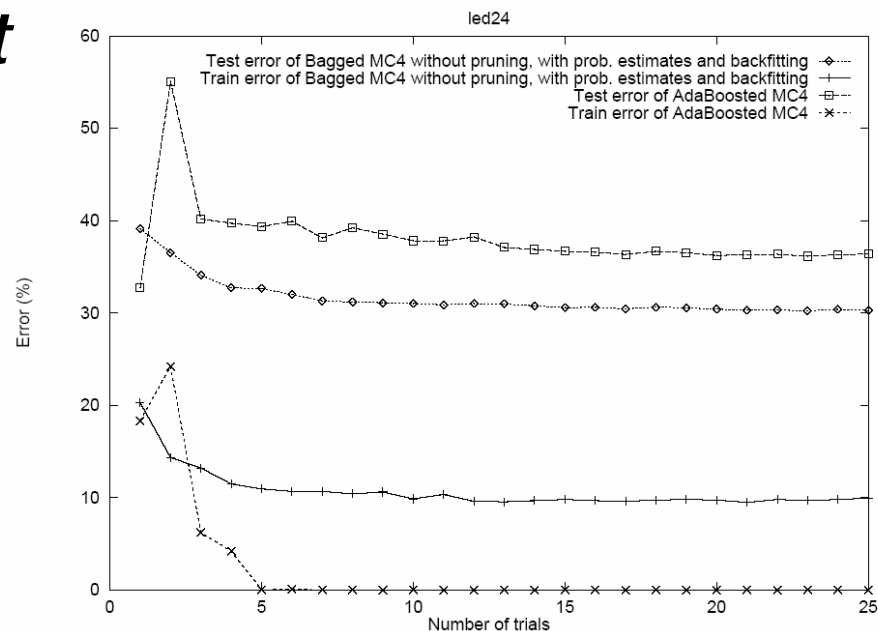
---

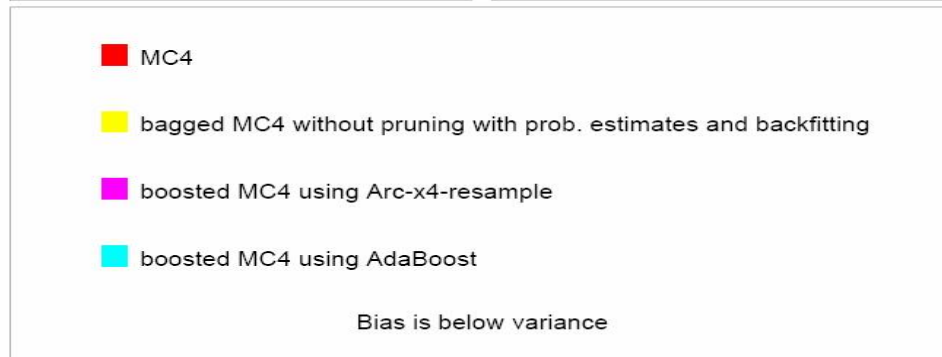
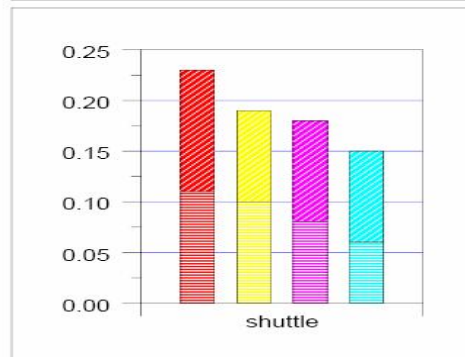
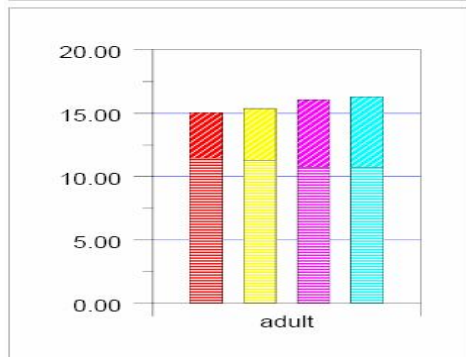
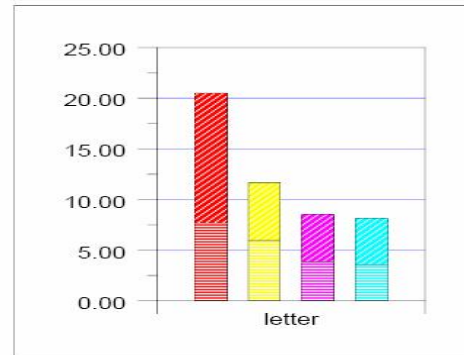
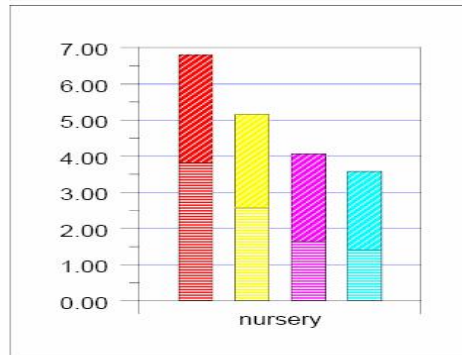
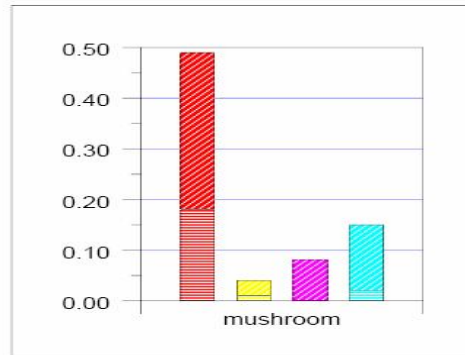
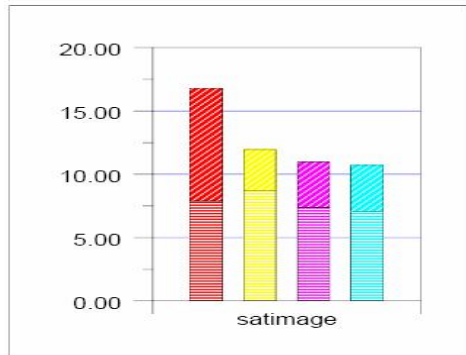
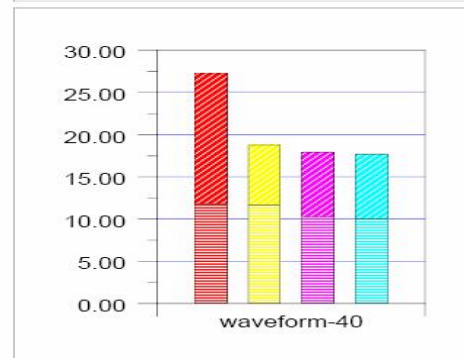
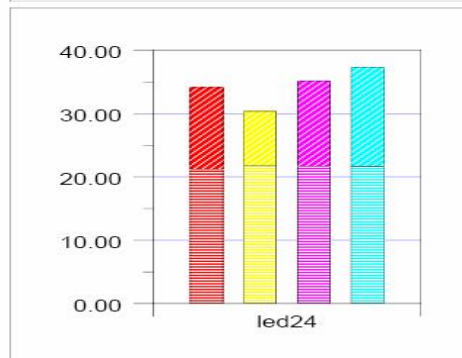
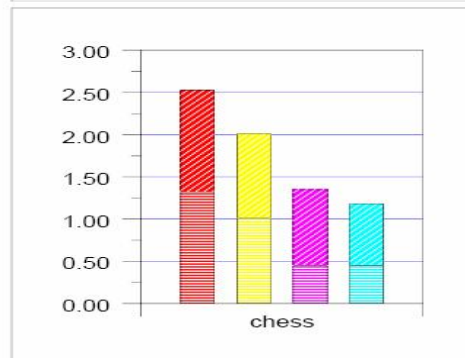
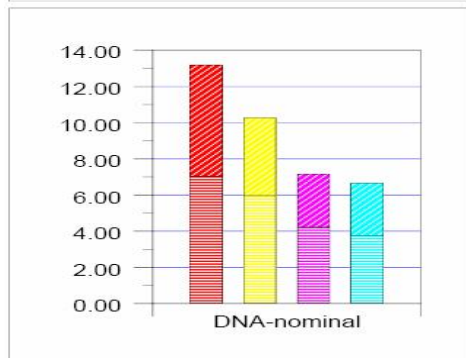
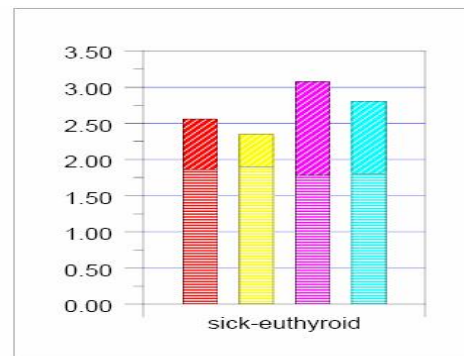
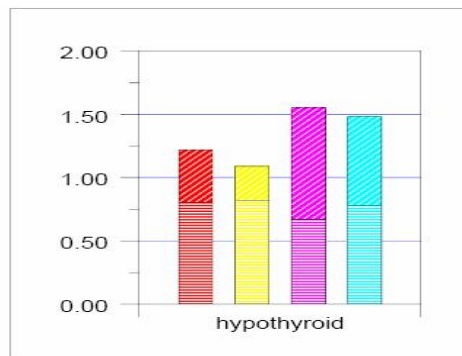
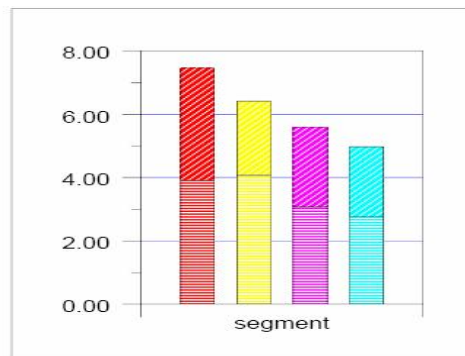
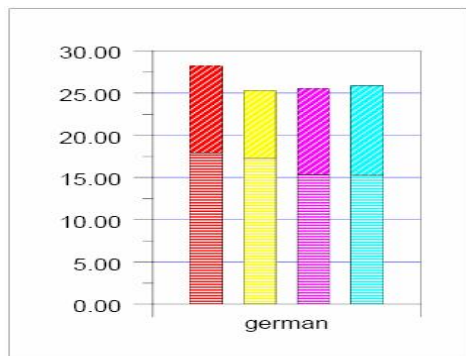
- i using the update rule that avoids the normalization step circumvents underflows early in the process, but they can still happen
- i authors set instances with weights falling below the minimum weight to have the minimum weight



# Boosting: AdaBoost results

- i AdaBoosted MC4 gave average relative error reduction of 27%
- i AdaBoost beat backfit-p-Bagging (average absolute error: 9.8% vs 10.1%)
- i AdaBoost was *not* uniformly better for all datasets







## Boosting: AdaBoost results

---

- i MC4(1) with boosting fails
- i AdaBoosted MC4(1)-disc gave average relative error reduction of 31%
- i AdaBoosted Naïve-Bayes gave average relative error reduction of 24%

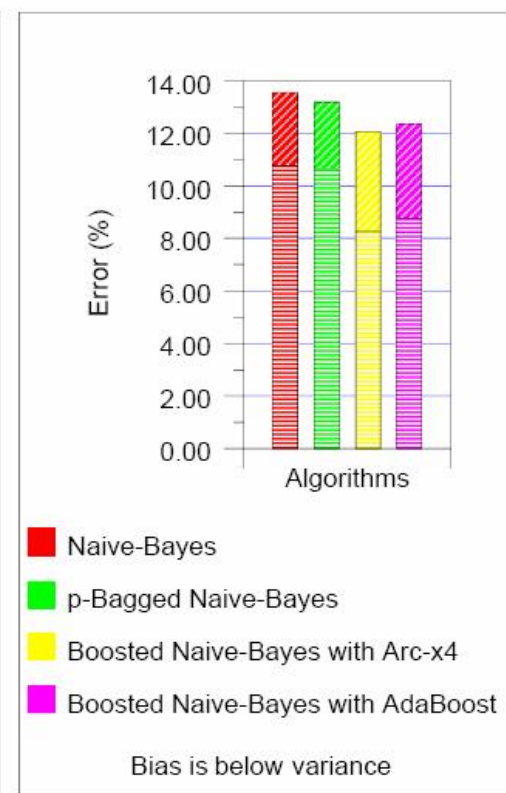
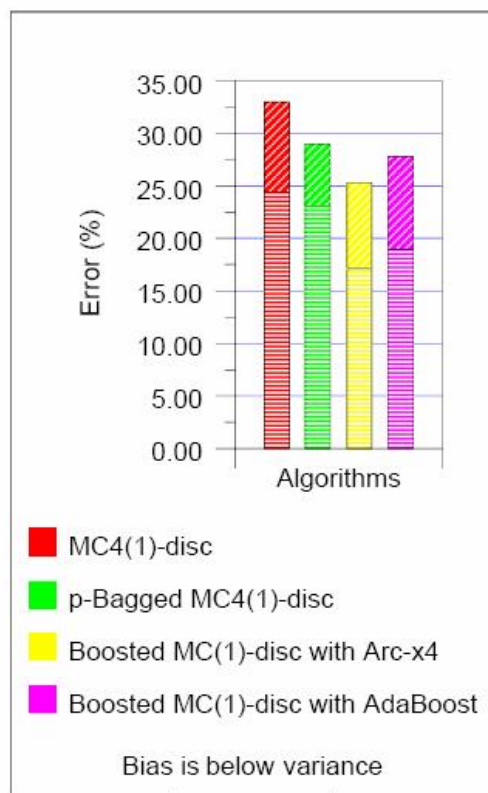
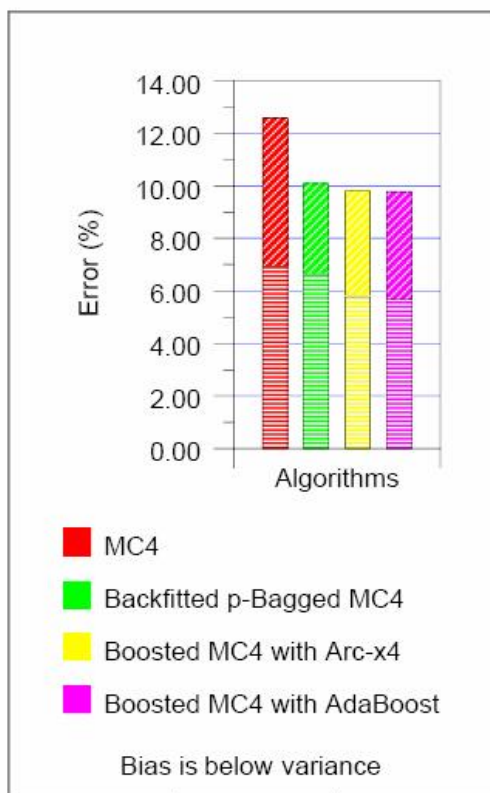


## Boosting: Arc-x4 results

---

- i two versions:
  - | Arc-x4-reweight
  - | Arc-x4-resample
- i Arc-x4-resample MC4 9.81% avg error
- i Arc-x4-reweight MC4 10.86% avg error
- i Arc-x4-reweight has higher variance
- i does worse on same datasets as AdaBoost, likely due to noise
- i Arc-x4 outperforms AdaBoost for MC4(1)-disc and Naïve-Bayes

# Boosting: summary





## Boosting: summary

---

- i *on average*, boosting is better than bagging for the given datasets
- i boosting is *not uniformly* better than bagging
- i AdaBoost does not deal well with noise
- i both algorithms reduced bias and variance for decision trees; increased variance for naïve-bayes



# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i **Future work**
- i Conclusions



## Future work

---

- i Improving robustness of boosting w.r.t. noise
- i Are there better ways to handle a “perfect” classifier when boosting?
- i Bagging & boosting seem to drift from Occam's razor; can we quantify this bias toward complex models?
- i Are there cases when pruning helps bagging?
- i Can stacking improve bagging & boosting even further?
- i Applying boosting to other methods: kNN
- i Other techniques of combining probabilities
- i Can boosting be parallelized?





# Outline

---

- i Introduction
- i The base inducers
- i The voting algorithms
- i Bias/Variance decomposition
- i Experimental design
- i Bagging and variants
- i Boosting: AdaBoost and Arc-x4
- i Future work
- i **Conclusions**



## Conclusions

---

- i Performance of bagging and boosting depends highly on the data set
- i Bagging very helpful w.r.t. MSE
- i Larger trees in AdaBoost trials correlated with greater success in reducing error
- i Arc-x4 works far better with re-sampling, unlike other boosting methods