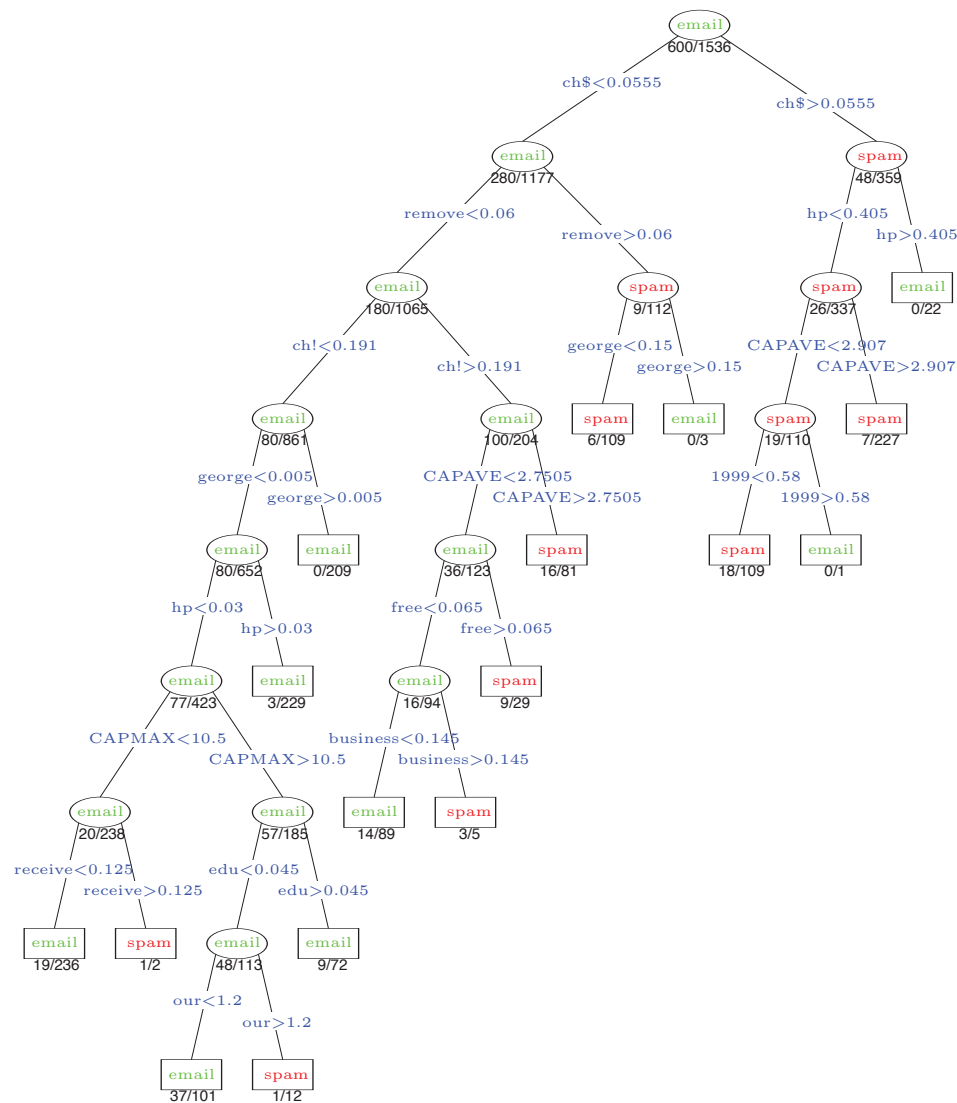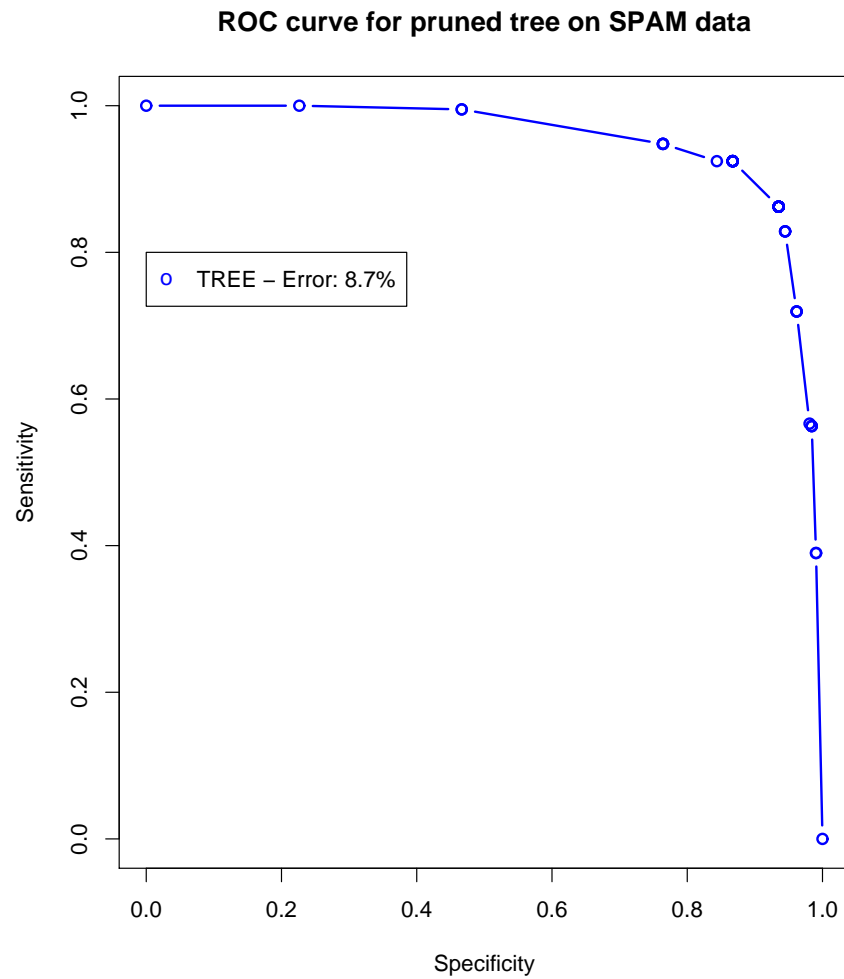# Ensemble Learners

- Classification Trees (Done)

- Bagging: Averaging Trees

- Random Forests: Cleverer Averaging of Trees

- Boosting: Cleverest Averaging of Trees

- Details in  chaps 10, 15 and 16

Methods for dramatically improving the performance of weak learners such as Trees. Classification trees are adaptive and robust, but do not make good prediction models. The techniques discussed here enhance their performance considerably.

# Properties of Trees

- [✔]Can handle huge datasets

- [✔]Can handle *mixed* predictors—quantitative and qualitative

- [✔]Easily ignore redundant variables

- [✔]Handle missing data elegantly through surrogate splits

- [✔]Small trees are easy to interpret

- [✘]Large trees are hard to interpret

- [✘]Prediction performance is often poor

**ROC curve for pruned tree on SPAM data**



○ TREE – Error: 8.7%

<div style="text-align:center">

**SPAM Data**

</div>

Overall error rate on test data: 8.7%.

ROC curve obtained by varying the *threshold* $c_0$ of the classifier:
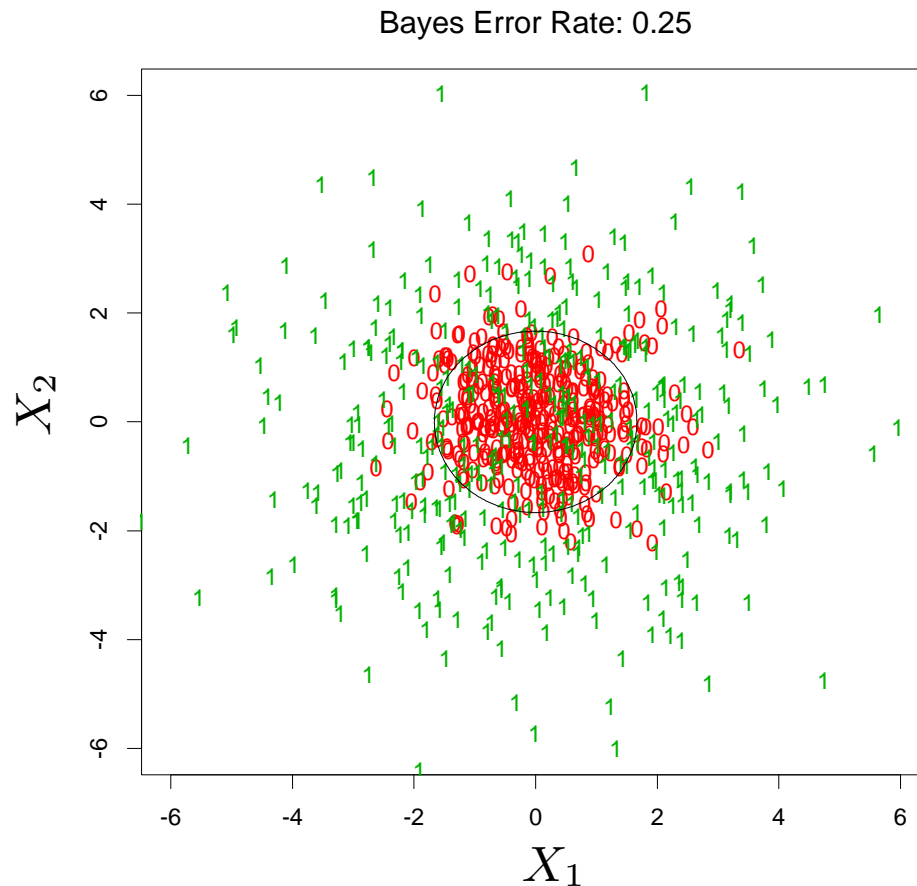
$C(X) = +1$ if $\widehat{\Pr}(+1|X) > c_0$.

Sensitivity: proportion of true spam identified

Specificity: proportion of true email identified.

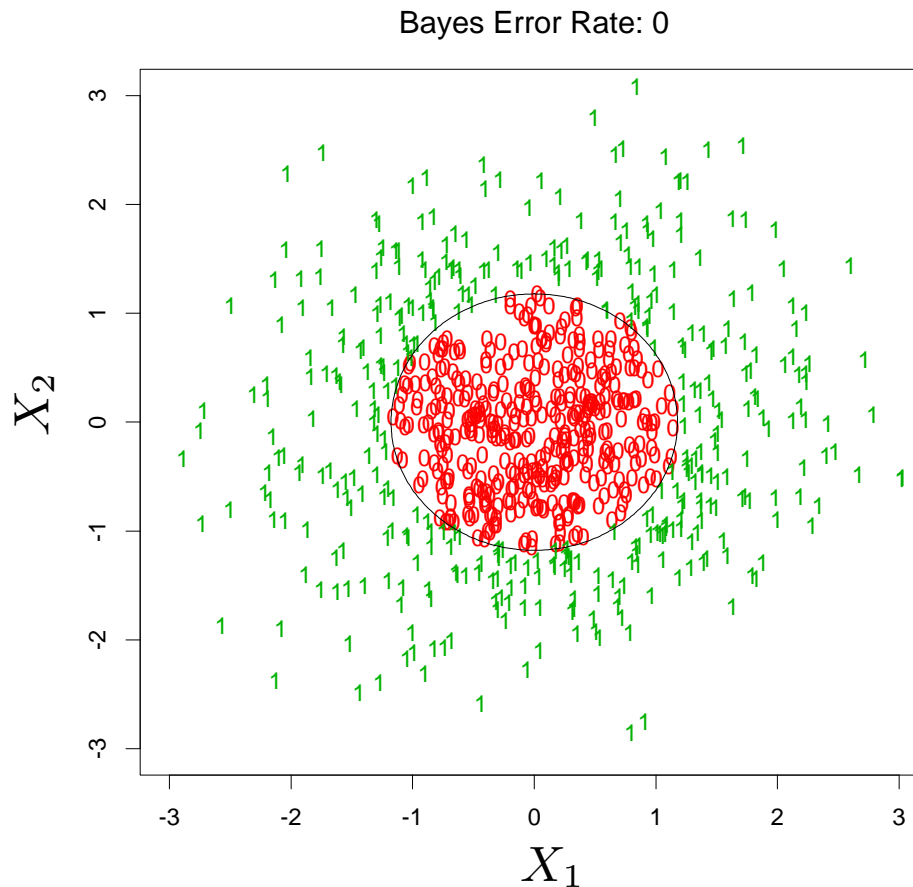We may want specificity to be high, and suffer some spam:

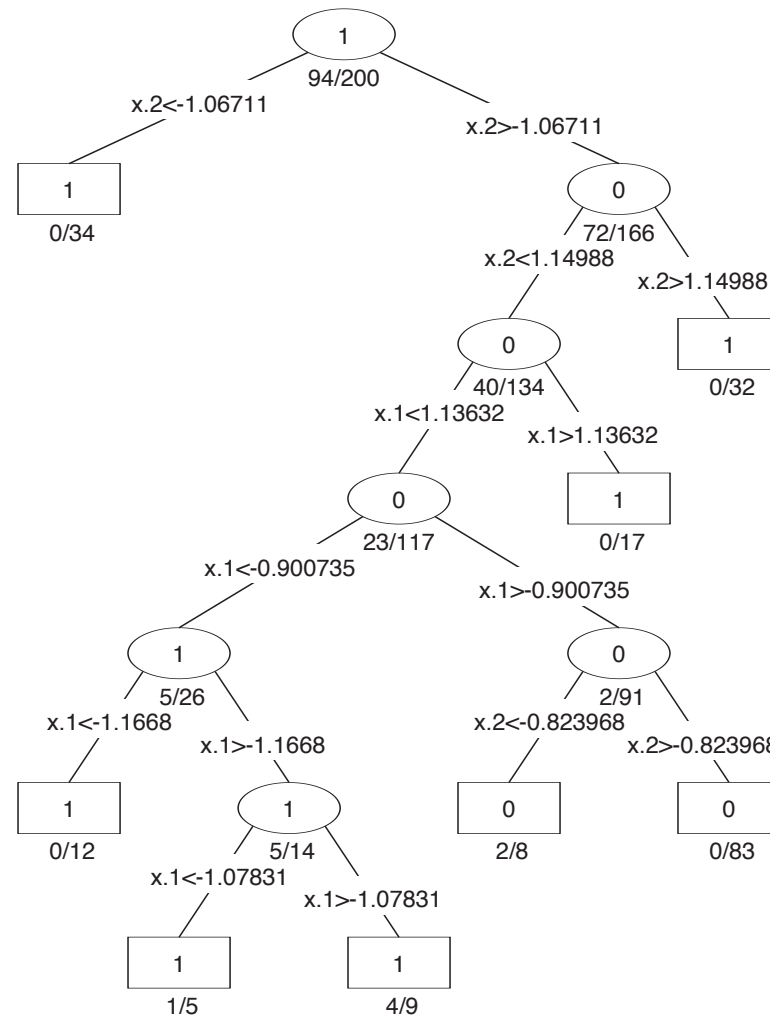Specificity : $95\% \implies$ Sensitivity : $79\%$

# Toy Classification Problem

Bayes Error Rate: 0.25



- Data $X$ and $Y$, with $Y$ taking values $+1$ or $-1$.

- Here $X = \in \mathbb{R}^2$

- The black boundary is the *Bayes Decision Boundary* - the best one can do. Here 25%

- Goal: given $N$ training pairs $(X_i, Y_i)$ produce a *classifier* $\hat{C}(X) \in \{-1, 1\}$

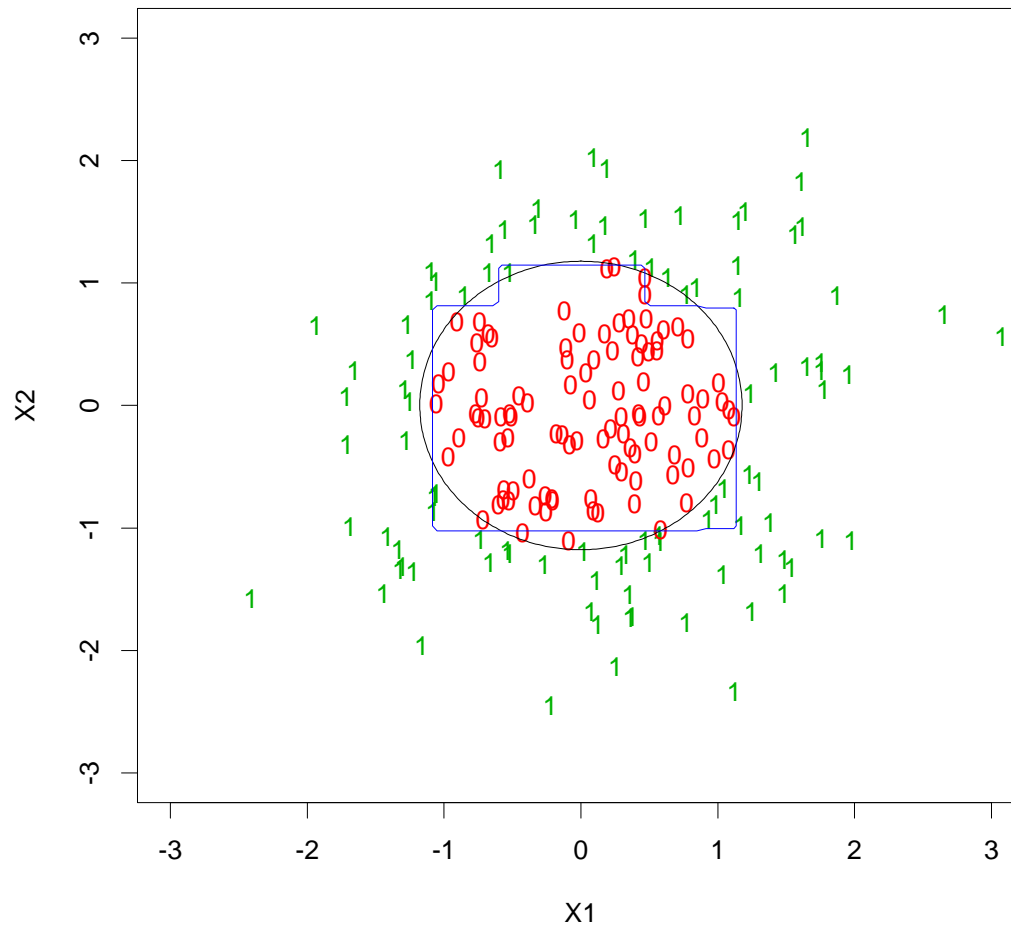- Also estimate the *probability* of the class labels $\Pr(Y = +1|X)$.

# Toy Example — No Noise

Bayes Error Rate: 0



- Deterministic problem; noise comes from sampling distribution of $X$.

- Use a training sample of size 200.

- Here *Bayes Error* is 0%.

# Classification Tree

# Decision Boundary: Tree

Error Rate: 0.073



- The eight-node tree is *boxy*, and has a test-error rate 0f 7.3%

- When the *nested spheres* are in 10-dimensions, a classification tree produces a rather noisy and inaccurate rule $\hat{C}(X)$, with error rates around 30%.
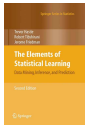
# Model Averaging

Classification trees can be simple, but often produce noisy (bushy) or weak (stunted) classifiers.

- *Bagging (Breiman, 1996):* Fit many large trees to bootstrap-resampled versions of the training data, and classify by majority vote.

- *Boosting (Freund & Shapire, 1996):* Fit many large or small trees to *reweighted* versions of the training data. Classify by weighted majority vote.

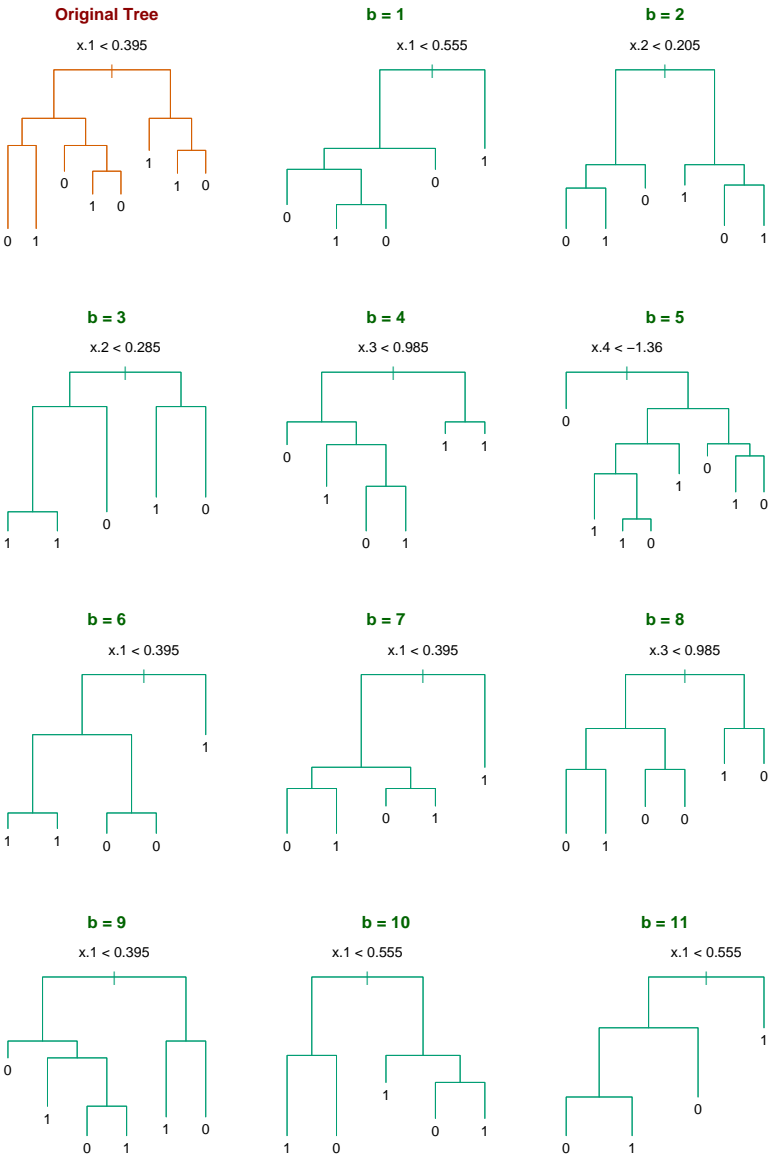- *Random Forests (Breiman 1999):* Fancier version of bagging.

In general *Boosting ≻ Random Forests ≻ Bagging ≻ Single Tree*.

# Bagging

- Bagging or *bootstrap aggregation* averages a noisy fitted function refit to many bootstrap samples, to reduce its variance. Bagging is a poor man's Bayes; See  pp 282.

- Suppose $C(\mathcal{S}, x)$ is a fitted classifier, such as a tree, fit to our training data $\mathcal{S}$, producing a predicted class label at input point $x$.

- To bag $C$, we draw bootstrap samples $\mathcal{S}^{*1}, \ldots \mathcal{S}^{*B}$ each of size $N$ with replacement from the training data. Then
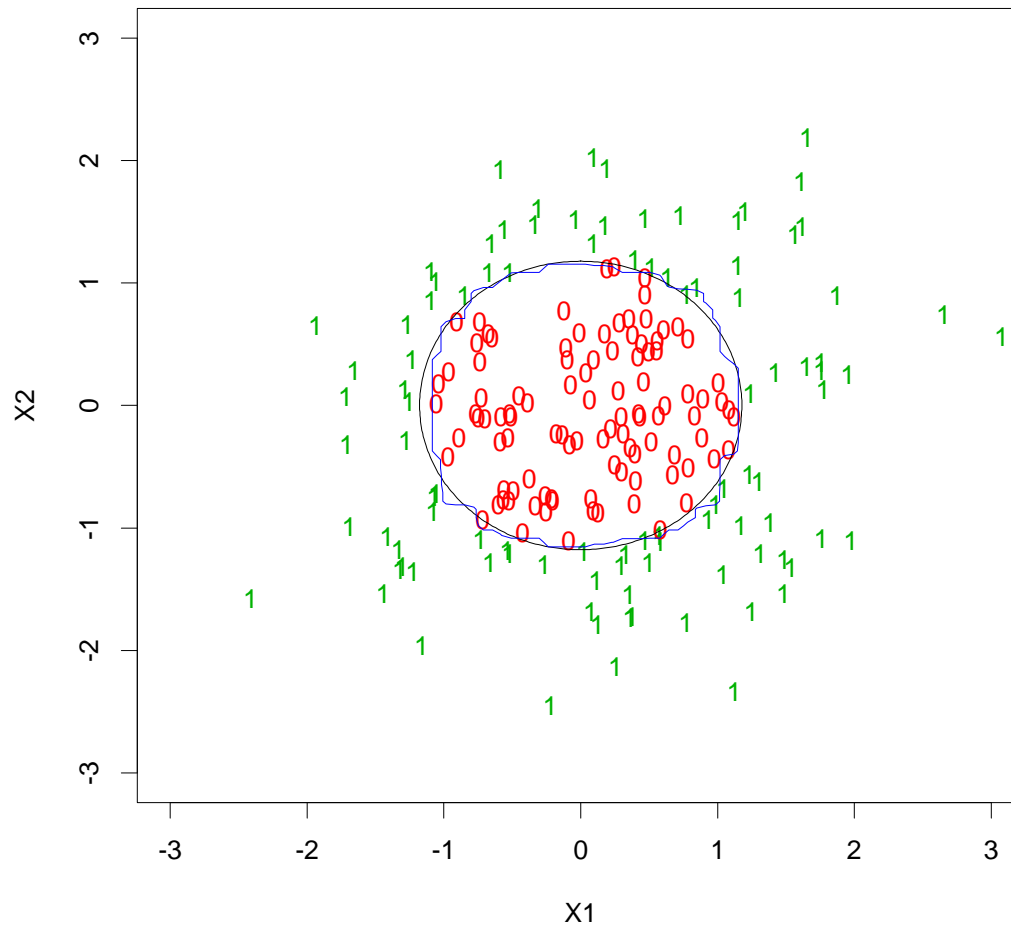
$$C_{bag}(x) = \text{Majority Vote } \{C(\mathcal{S}^{*b}, x)\}_{b=1}^{B}.$$

- Bagging can dramatically reduce the variance of unstable procedures (like trees), leading to improved prediction. However any simple structure in $C$ (e.g a tree) is lost.
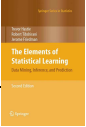
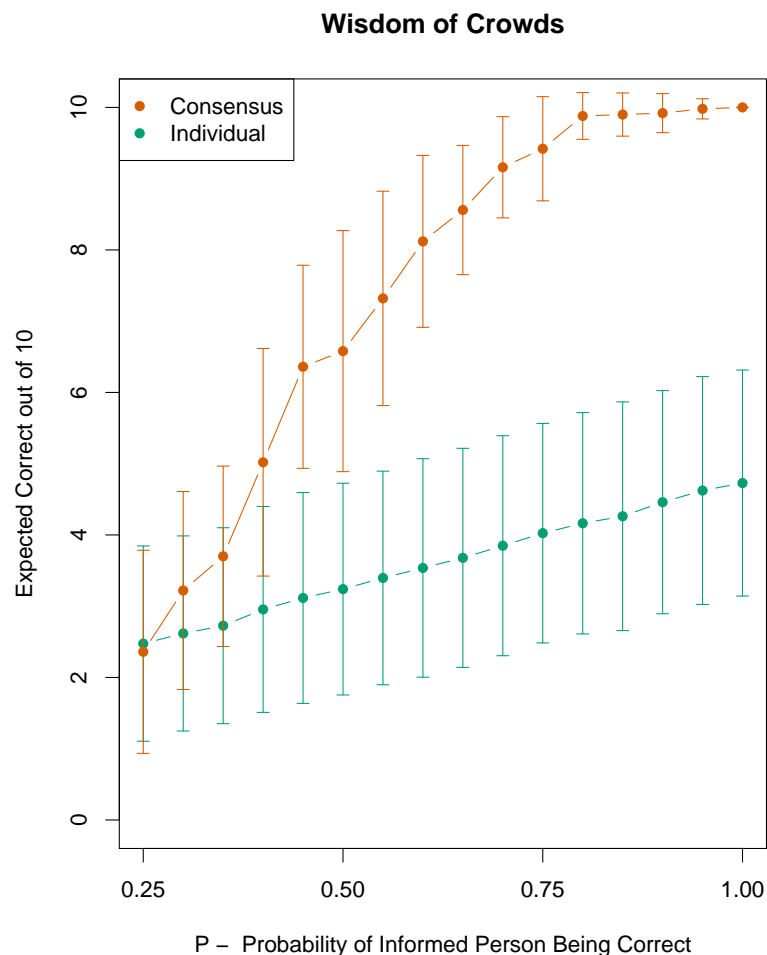# Decision Boundary: Bagging

Error Rate: 0.032



- Bagging averages many trees, and produces *smoother* decision boundaries.

- Bagging reduces error by variance reduction.

- Bias is slightly increased because bagged trees are slightly shallower.

# Random forests

- refinement of bagged trees; very popular—chapter 15.

- at each tree split, a random sample of $m$ features is drawn, and only those $m$ features are considered for splitting. Typically $m = \sqrt{p}$ or $\log_2 p$, where $p$ is the number of features

- For each tree grown on a bootstrap sample, the error rate for observations left out of the bootstrap sample is monitored. This is called the "out-of-bag" or OOB error rate.

- random forests tries to improve on bagging by "de-correlating" the trees, and reduce the variance.

- Each tree has the same expectation, so increasing the number of trees does not alter the bias of bagging or random forests.

## Wisdom of Crowds



**Wisdom of Crowds**

- 10 movie categories, 4 choices in each category

- 50 people, for each question 15 *informed* and 35 *guessers*

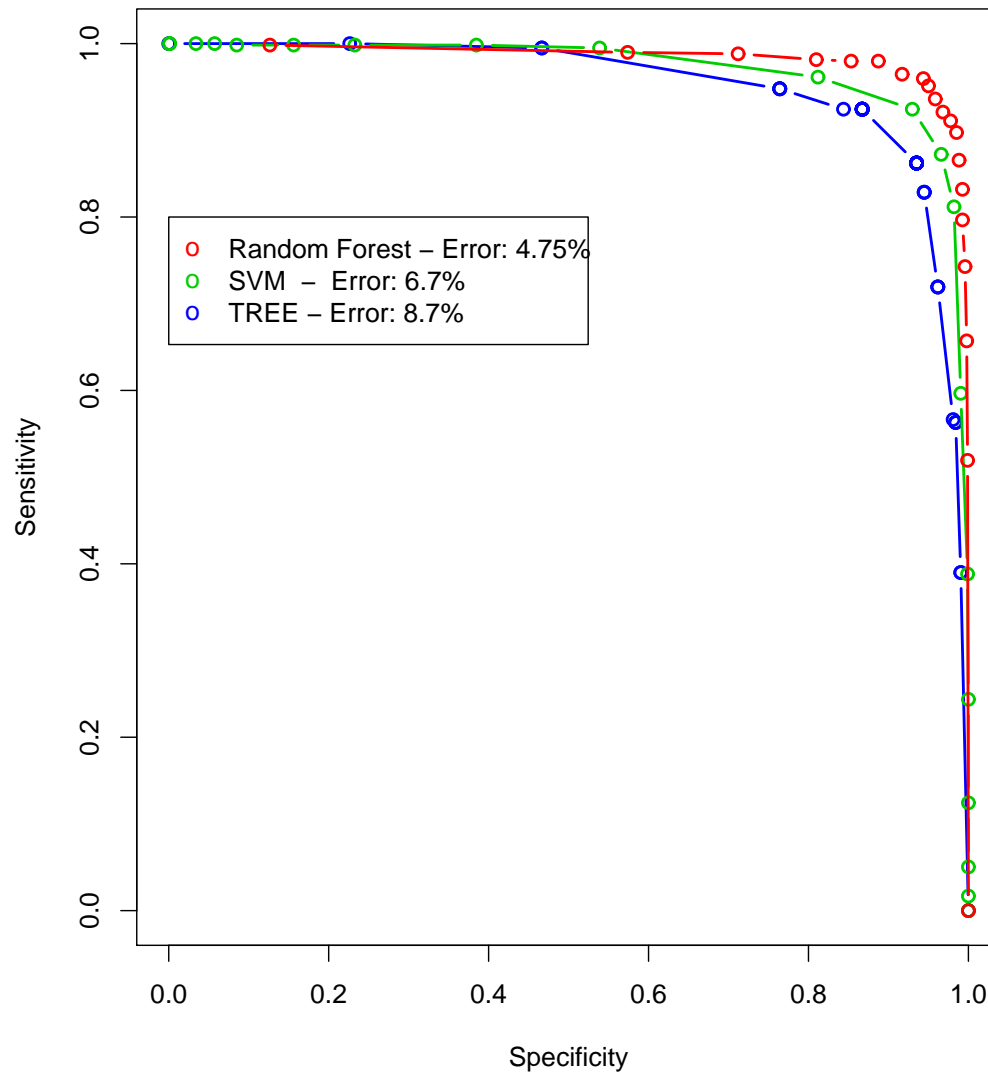- For informed people,

$$\mathrm{Pr}(\text{correct}) \quad = \quad p$$
$$\mathrm{Pr}(\text{incorrect}) \quad = \quad (1 - p)/3$$

- For guessers,
  $\mathrm{Pr}(\text{all choices}) = 1/4$.

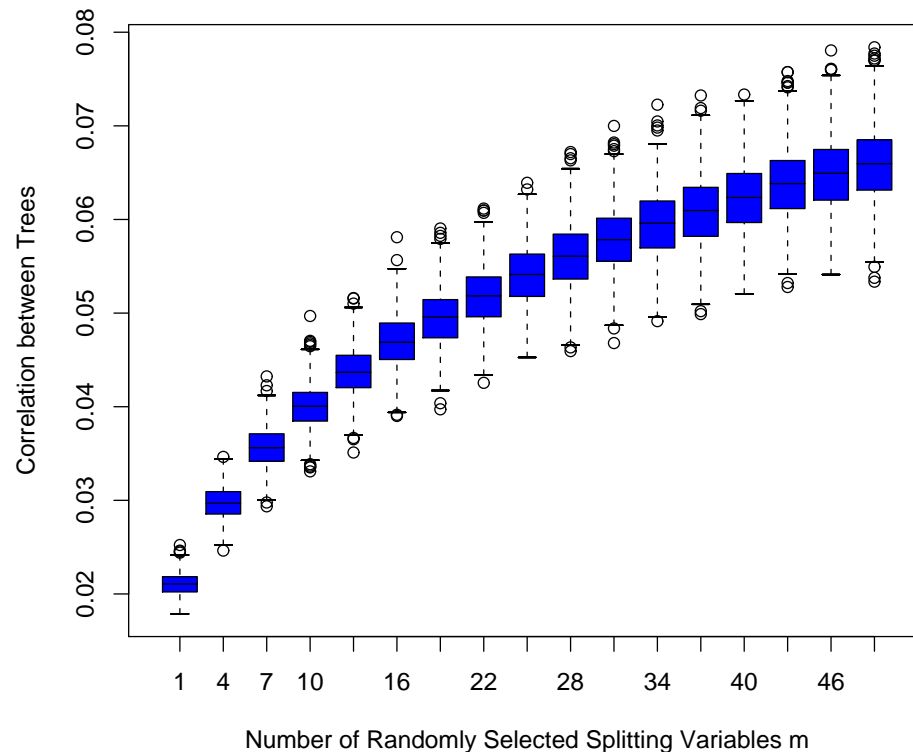Consensus (orange) improve over individuals (teal) even for moderate p.

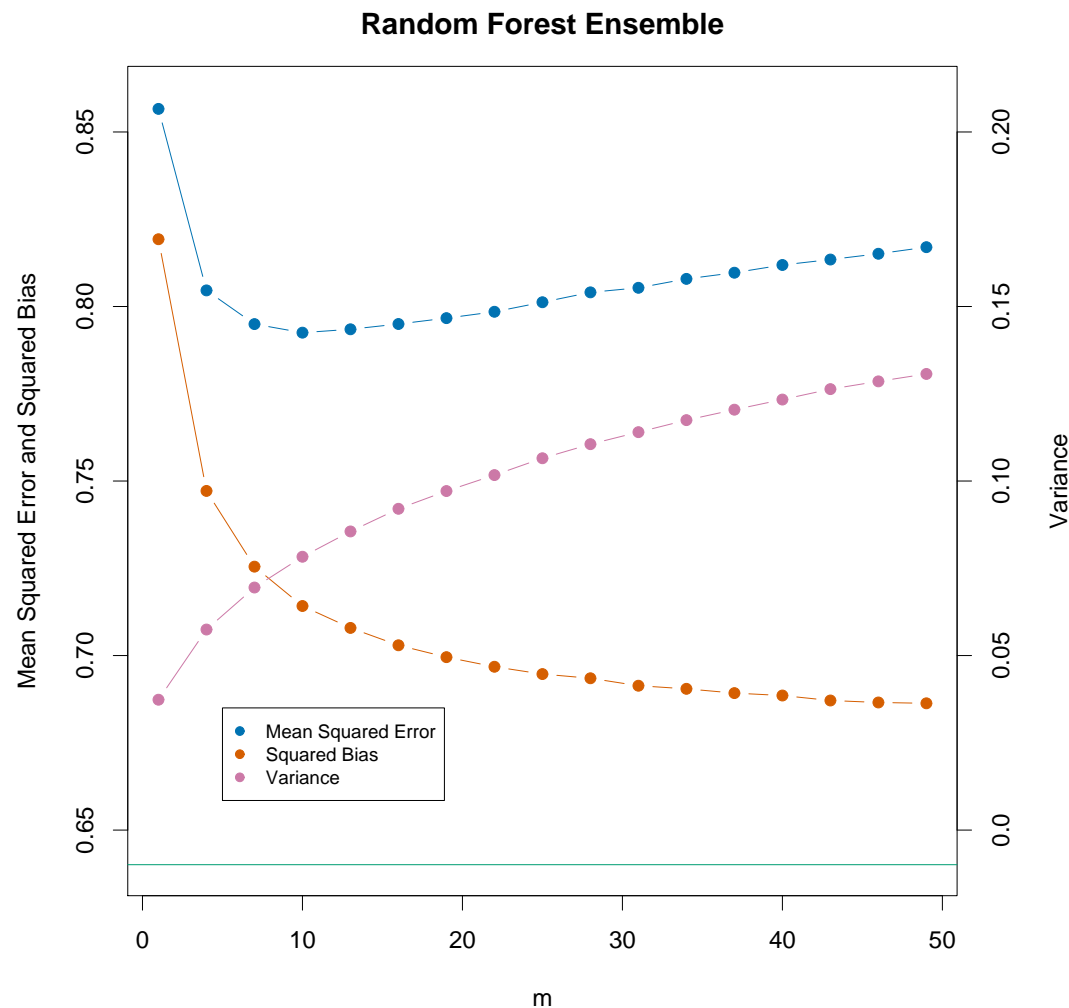**ROC curve for TREE, SVM and Random Forest on SPAM data**



## TREE, SVM and RF

Random Forest dominates both other methods on the SPAM data — 4.75% error. Used 500 trees with default settings for `random Forest` package in `R`.

Correlation between Trees

Number of Randomly Selected Splitting Variables m

## RF: decorrelation

- $\mathrm{Var}\hat{f}_{RF}(x) = \rho(x)\sigma^2(x)$

- The correlation is because we are growing trees to the *same data*.

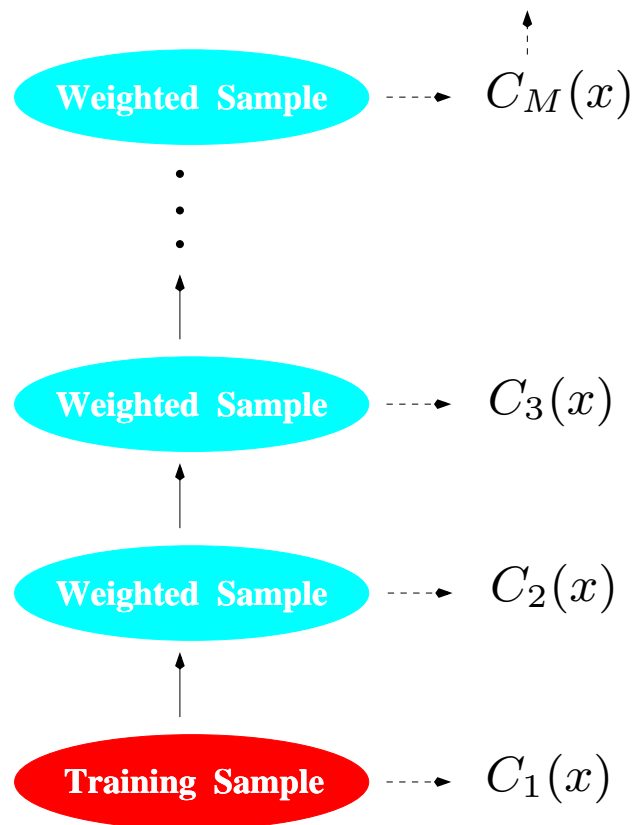- The correlation decreases as $m$ decreases—$m$ is the number of variables sampled at each split.

Based on a simple experiment with random forest regression models.

**Random Forest Ensemble**

## Bias & Variance

- The smaller $m$, the lower the variance of the random forest ensemble

- Small $m$ is also associated with higher bias, because important variables can be missed by the sampling.
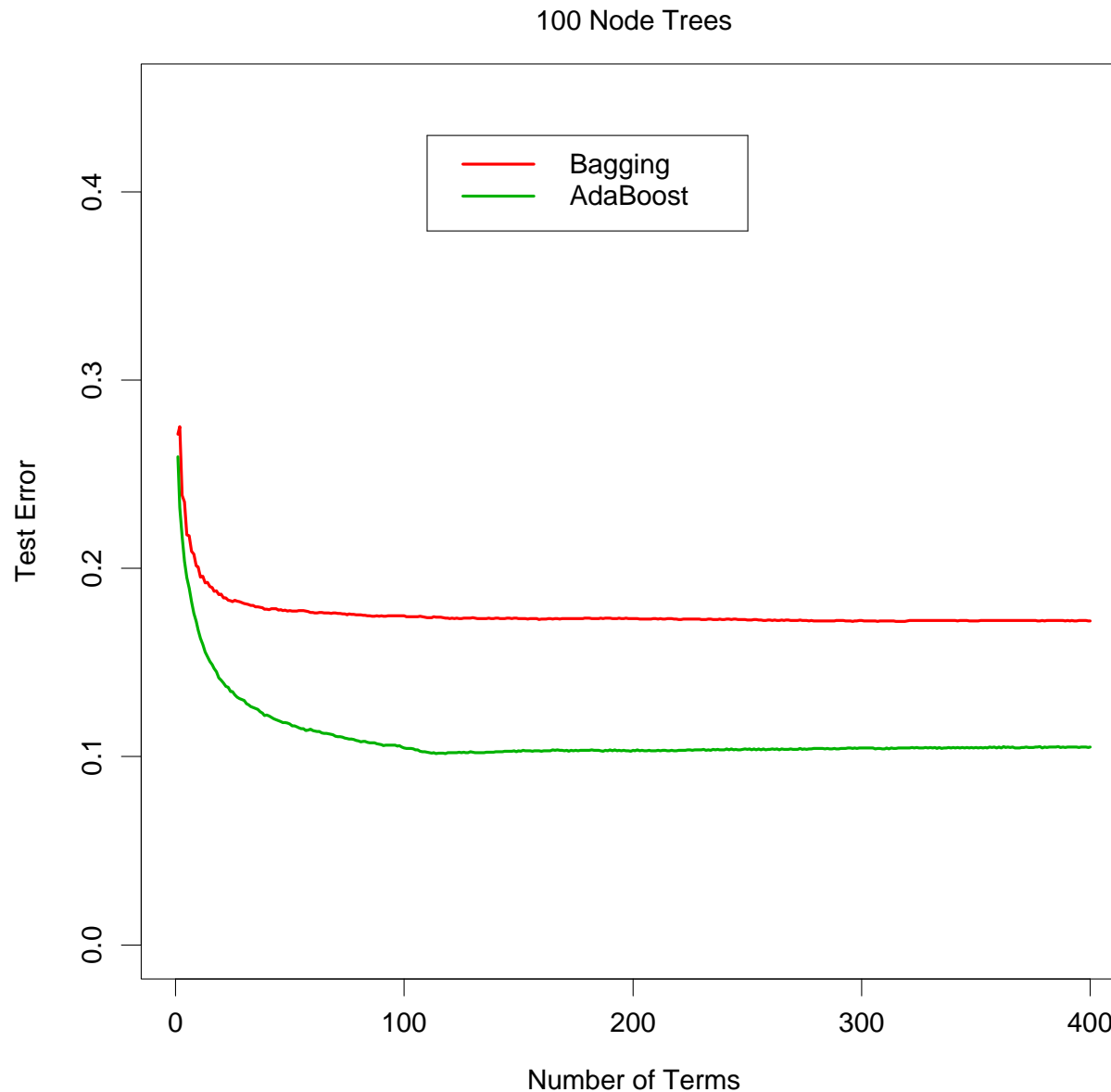
# Boosting

Weighted Sample $\dashrightarrow C_M(x)$

Weighted Sample $\dashrightarrow C_3(x)$

Weighted Sample $\dashrightarrow C_2(x)$

Training Sample $\dashrightarrow C_1(x)$

- Breakthrough invention of Freund & Schapire (1997)

- Average many trees, each grown to re-weighted versions of the training data.

- Weighting *decorrelates* the trees, by focussing on regions missed by past trees.

- Final Classifier is weighted average of classifiers:

$$C(x) = \text{sign}\left[\sum_{m=1}^{M} \alpha_m C_m(x)\right]$$

Note: $C_m(x) \in \{-1, +1\}$.
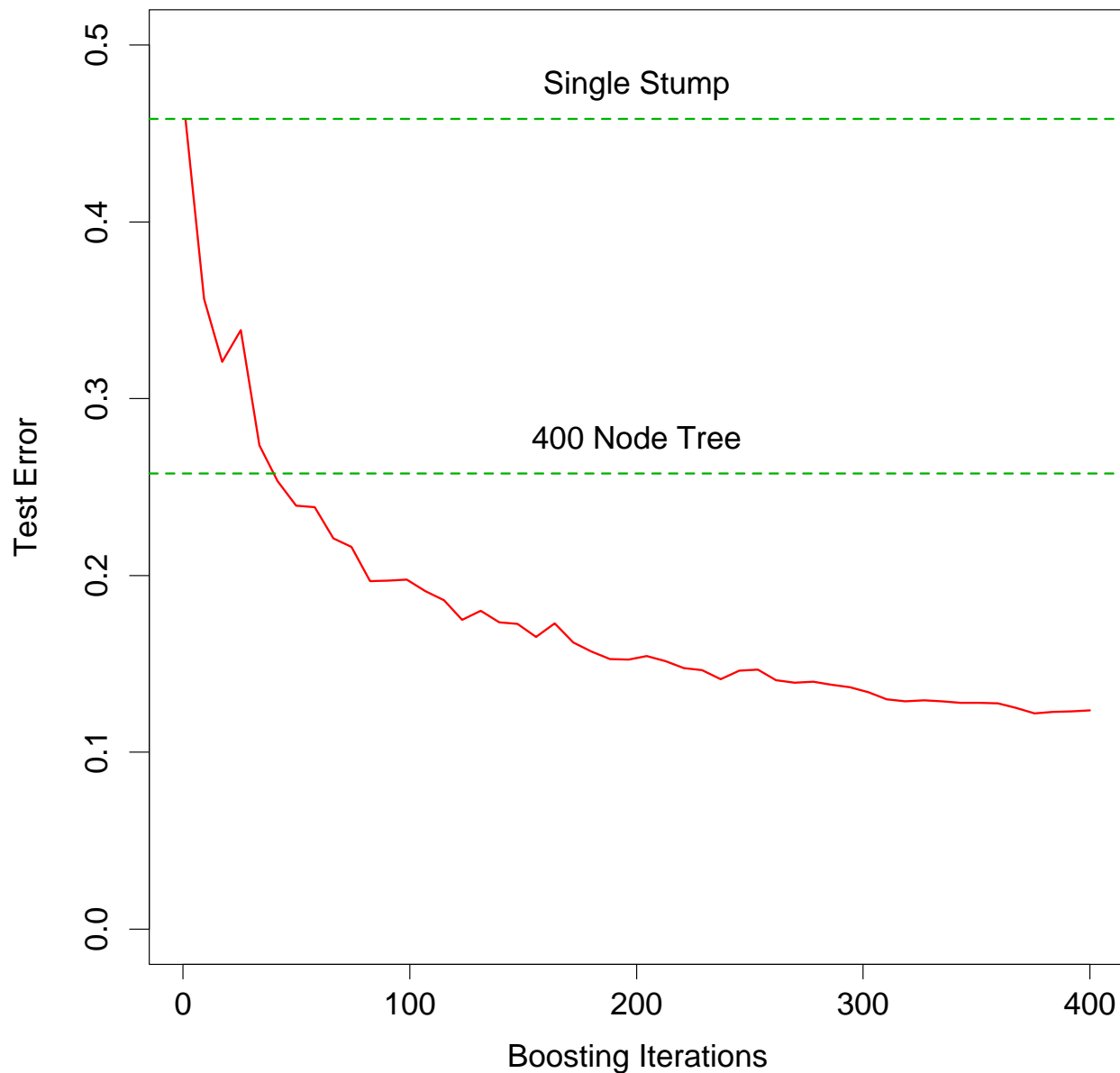
100 Node Trees



# Boosting vs Bagging

- 2000 points from Nested Spheres in $R^{10}$

- Bayes error rate is 0%.

- Trees are grown *best first* without pruning.

- Leftmost term is a single tree.

# AdaBoost (Freund & Schapire, 1996)

1. Initialize the observation weights $w_i = 1/N$, $i = 1, 2, \ldots, N$.

2. For $m = 1$ to $M$ repeat steps (a)–(d):

   (a) Fit a classifier $C_m(x)$ to the training data using weights $w_i$.

   (b) Compute weighted error of newest tree

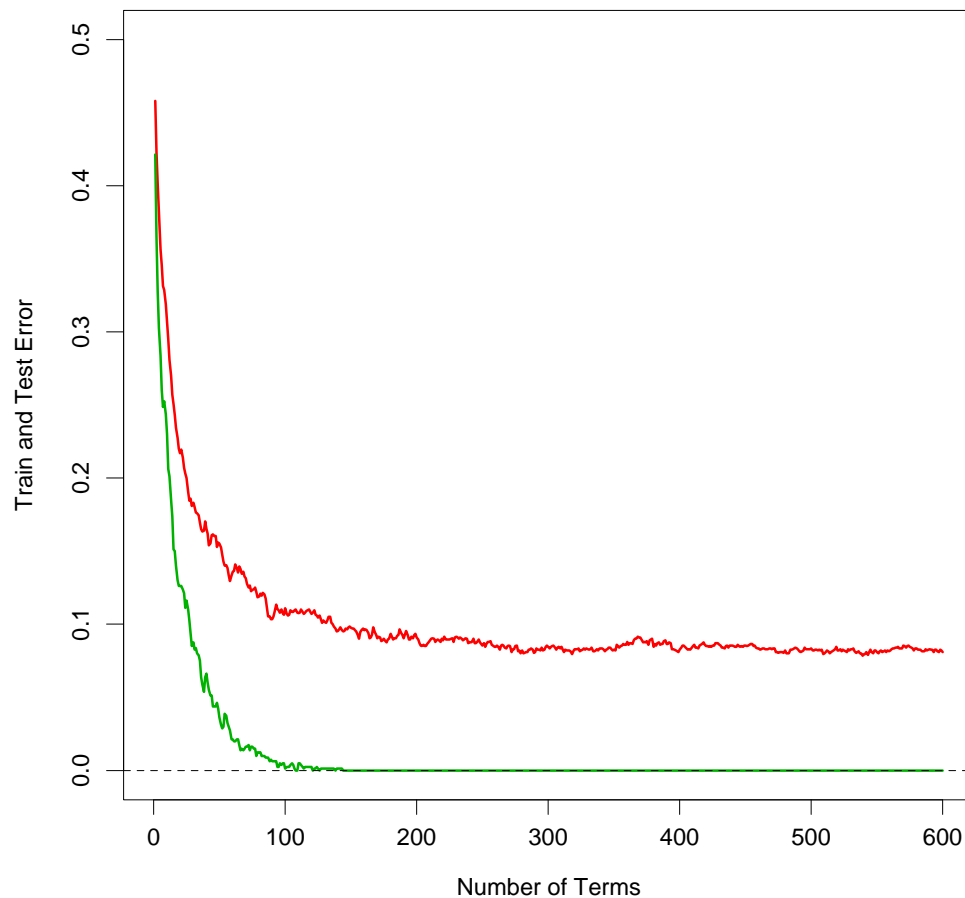   $$\text{err}_m = \frac{\sum_{i=1}^{N} w_i I(y_i \neq C_m(x_i))}{\sum_{i=1}^{N} w_i}.$$

   (c) Compute $\alpha_m = \log[(1 - \text{err}_m)/\text{err}_m]$.

   (d) Update weights for $i = 1, \ldots, N$:
   $$w_i \leftarrow w_i \cdot \exp[\alpha_m \cdot I(y_i \neq C_m(x_i))]$$
   and renormalize to $w_i$ to sum to 1.

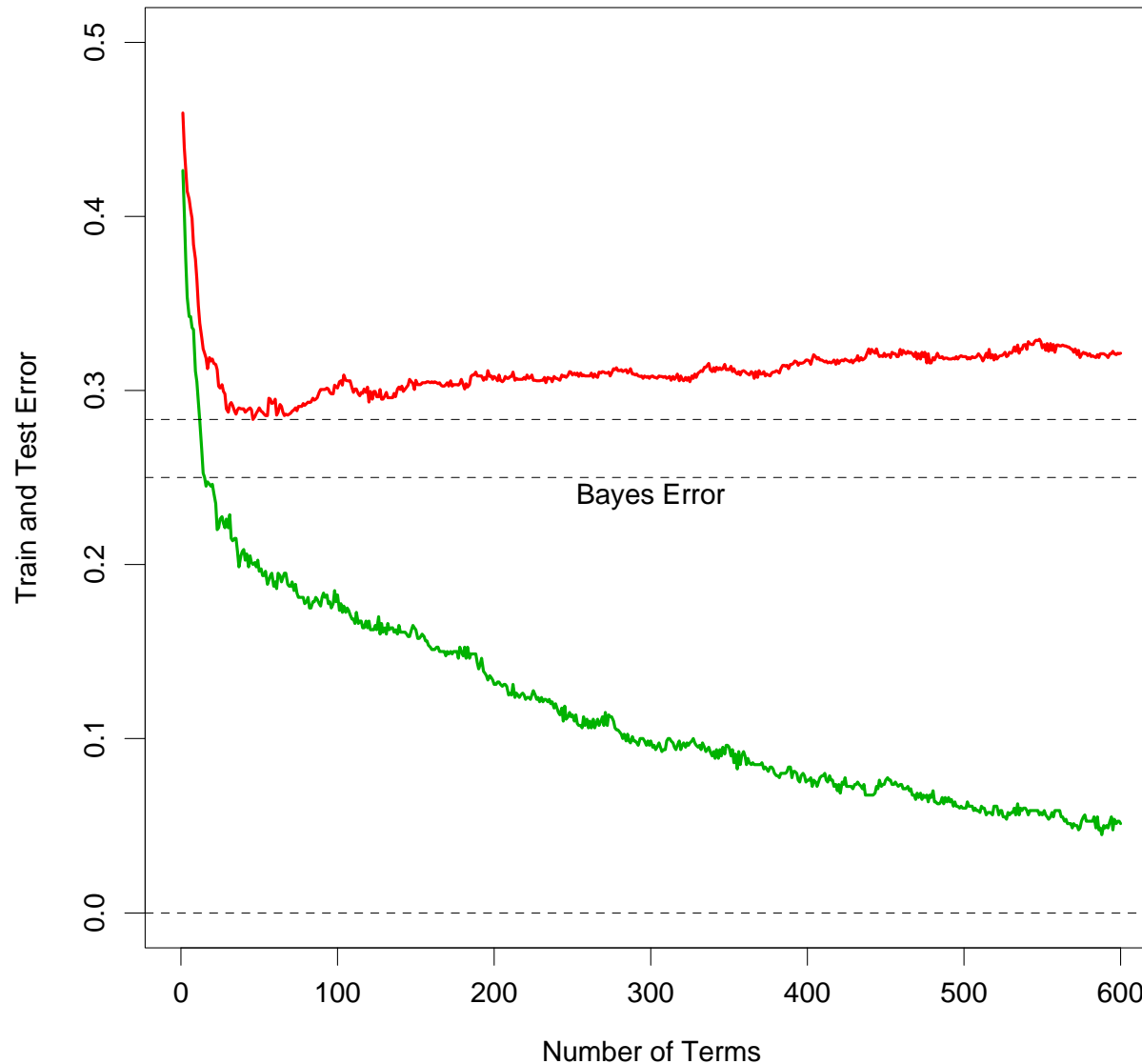3. Output $C(x) = \text{sign}\left[\sum_{m=1}^{M} \alpha_m C_m(x)\right]$.

# Boosting Stumps

- A stump is a two-node tree, after a single split.

- Boosting stumps works remarkably well on the nested-spheres problem.

## Training Error



- Nested spheres in 10-Dimensions.

- Bayes error is 0%.

- Boosting drives the training error to zero (green curve).

- Further iterations continue to improve test error in many examples (red curve).

# Noisy Problems

- Nested Gaussians in 10-Dimensions.

- Bayes error is 25%.

- Boosting with stumps

- Here the test error does increase, but quite slowly.

# Stagewise Additive Modeling

Boosting builds an additive model

$$F(x) = \sum_{m=1}^{M} \beta_m b(x; \gamma_m).$$

Here $b(x, \gamma_m)$ is a tree, and $\gamma_m$ parametrizes the splits.

We do things like that in statistics all the time!

- GAMs: $F(x) = \sum_j f_j(x_j)$

- Basis expansions: $F(x) = \sum_{m=1}^{M} \theta_m h_m(x)$

Traditionally the parameters $f_m$, $\theta_m$ are fit *jointly* (i.e. least squares, maximum likelihood).

With boosting, the parameters $(\beta_m, \gamma_m)$ are fit in a *stagewise* fashion. This slows the process down, and overfits less quickly.
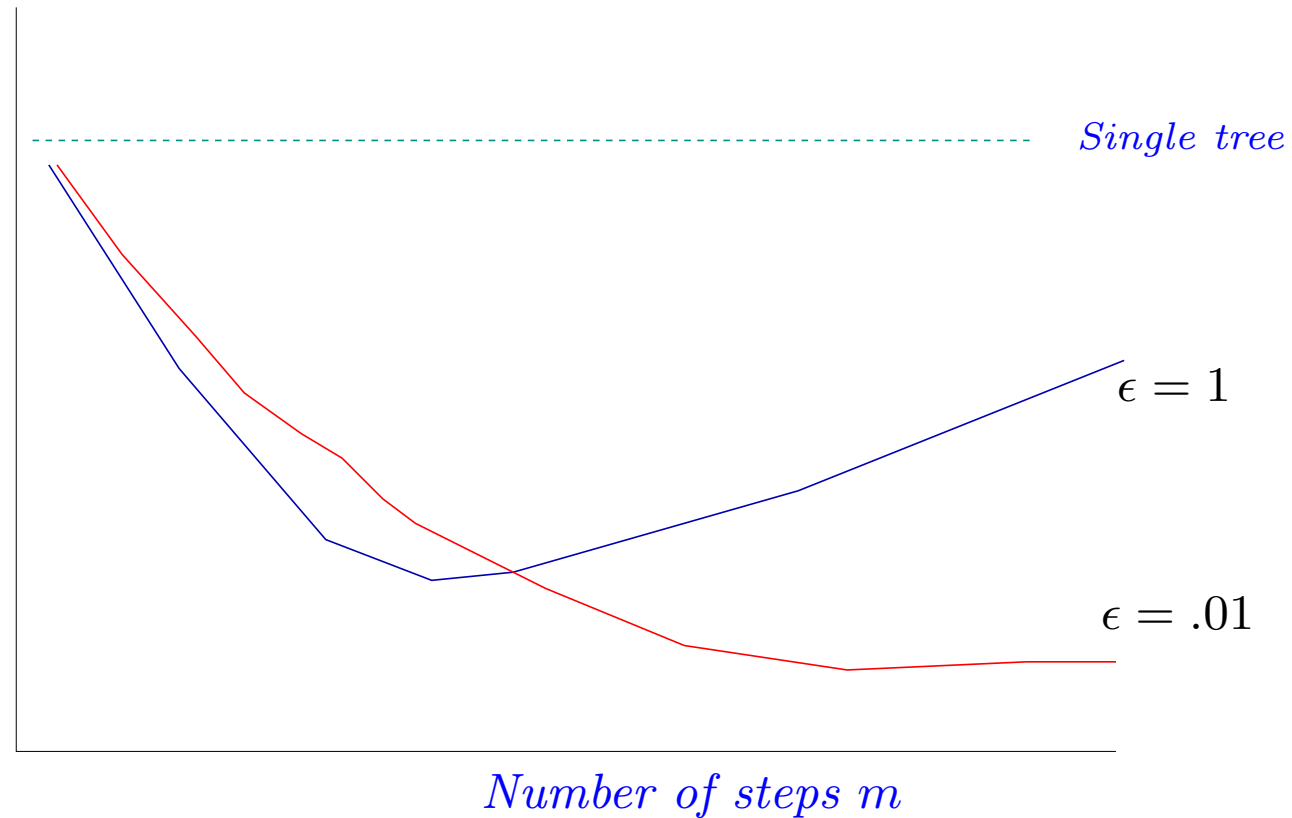
# Example: Least Squares Boosting

1. Start with function $F_0(x) = 0$ and residual $r = y$, $m = 0$

2. $m \leftarrow m + 1$

3. Fit a CART regression tree to $r$ giving $g(x)$

4. Set $f_m(x) \leftarrow \epsilon \cdot g(x)$

5. Set $F_m(x) \leftarrow F_{m-1}(x) + f_m(x)$, $r \leftarrow r - f_m(x)$ and repeat step 2–5 many times

- $g(x)$ typically shallow tree; e.g. constrained to have $k = 3$ splits only (depth).

- Stagewise fitting (no backfitting!). Slows down (over-) fitting process.

- $0 < \epsilon \leq 1$ is shrinkage parameter; slows overfitting even further.

# Least Squares Boosting and $\epsilon$

*Prediction Error*



*Single tree*

$\epsilon = 1$

$\epsilon = .01$

*Number of steps m*

- If instead of trees, we use linear regression terms, boosting with small $\epsilon$ very similar to lasso!

# Adaboost: Stagewise Modeling

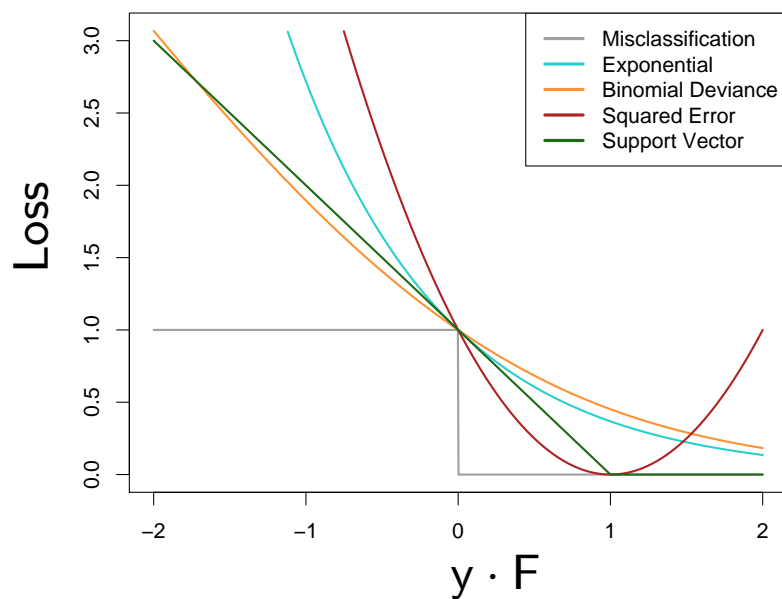- AdaBoost builds an *additive logistic regression model*

$$F(x) = \log \frac{\Pr(Y = 1|x)}{\Pr(Y = -1|x)} = \sum_{m=1}^{M} \alpha_m f_m(x)$$

by *stagewise fitting* using the loss function

$$L(y, F(x)) = \exp(-y\, F(x)).$$

- Instead of fitting trees to residuals, the special form of the exponential loss leads to fitting trees to weighted versions of the original data. See  pp 343.

# Why Exponential Loss?



- $y \cdot F$ is the *margin*; positive margin means correct classification.

- $e^{-yF(x)}$ is a monotone, smooth upper bound on misclassification loss at $x$.

- Leads to simple reweighting scheme.

- Has *logit* transform as population minimizer

$$F^*(x) = \frac{1}{2} \log \frac{\Pr(Y = 1|x)}{\Pr(Y = -1|x)}$$

- Other more robust loss functions, like *binomial deviance.*

# General Boosting Algorithms

The boosting paradigm can be extended to general loss functions
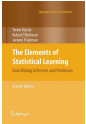— not only *squared error* or *exponential*.

1. Initialize $F_0(x) = 0$.

2. For $m = 1$ to $M$:

   (a) Compute
       $$(\beta_m, \gamma_m) = \arg\min_{\beta, \gamma} \sum_{i=1}^{N} L(y_i, F_{m-1}(x_i) + \beta b(x_i; \gamma)).$$
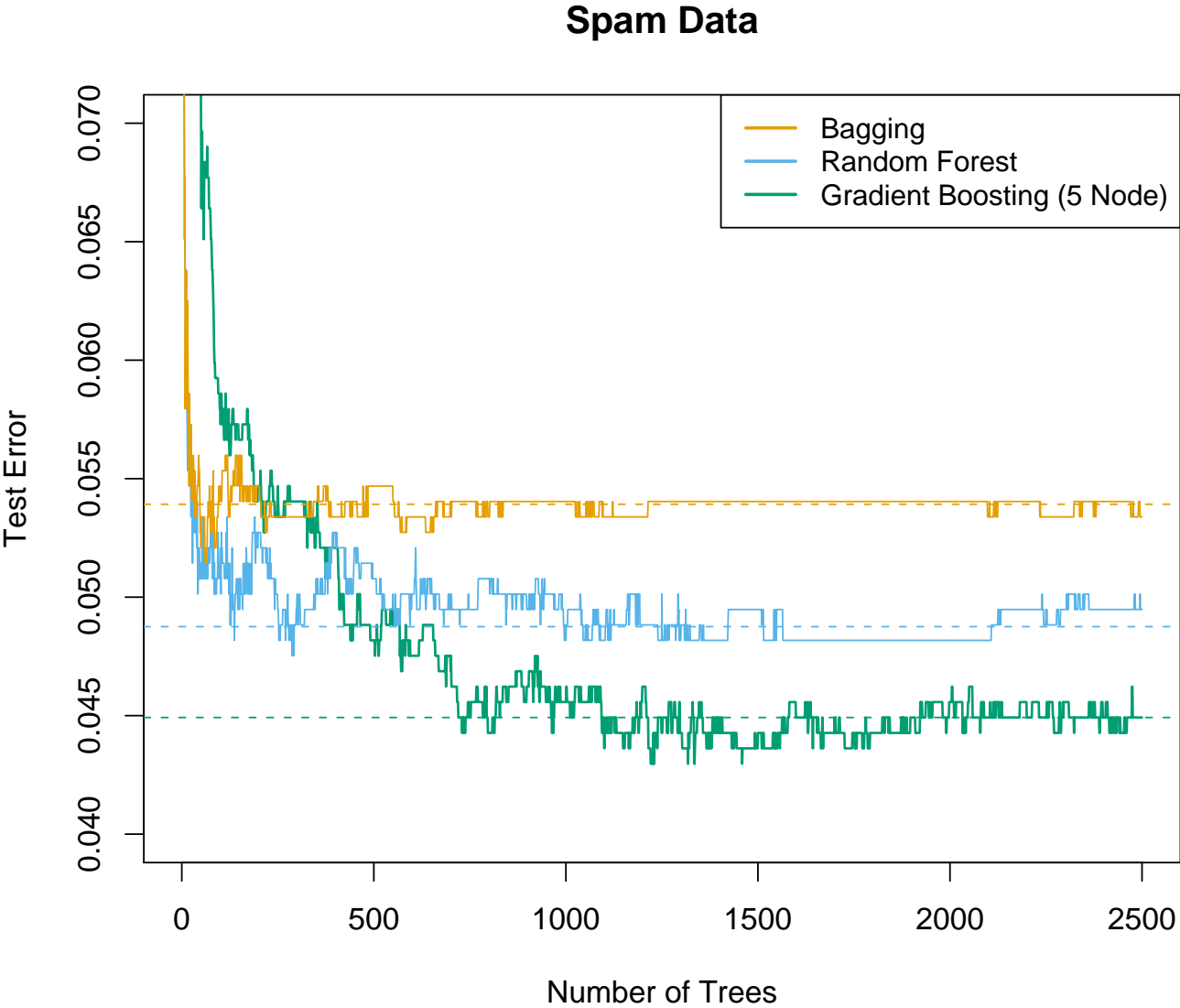
   (b) Set $F_m(x) = F_{m-1}(x) + \beta_m b(x; \gamma_m)$.

Sometimes we replace step (b) in item 2 by

   (b*) Set $F_m(x) = F_{m-1}(x) + \epsilon \beta_m b(x; \gamma_m)$

Here $\epsilon$ is a *shrinkage factor*; in `gbm` package in R, the default is
$\epsilon = 0.001$. Shrinkage slows the stagewise model-building even more,
and typically leads to better performance.

# Gradient Boosting

- General boosting algorithm that works with a variety of different loss functions. Models include regression, resistant regression, K-class classification and risk modeling.

- Gradient Boosting builds additive tree models, for example, for representing the logits in logistic regression.

- *Tree size* is a parameter that determines the order of interaction (next slide).

- Gradient Boosting inherits all the good features of trees (variable selection, missing data, mixed predictors), and improves on the weak features, such as prediction performance.

- Gradient Boosting is described in detail in , section 10.10, and is implemented in Greg Ridgeways `gbm` package in R.

**Spam Data**

## Spam Example Results

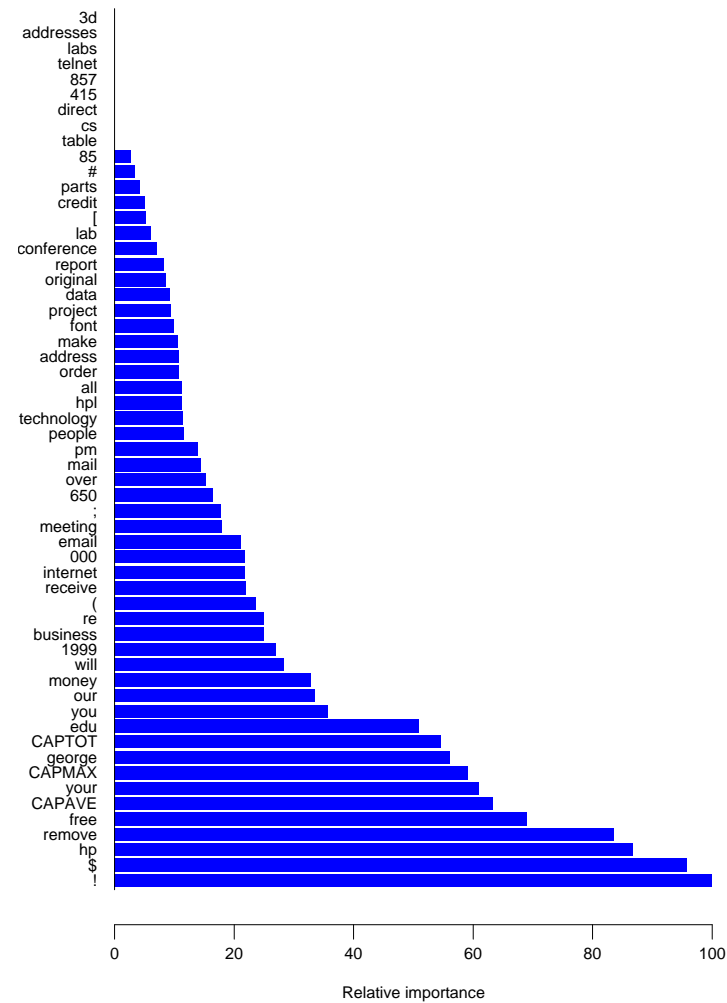With 3000 training and 1500 test observations, Gradient Boosting fits an additive logistic model
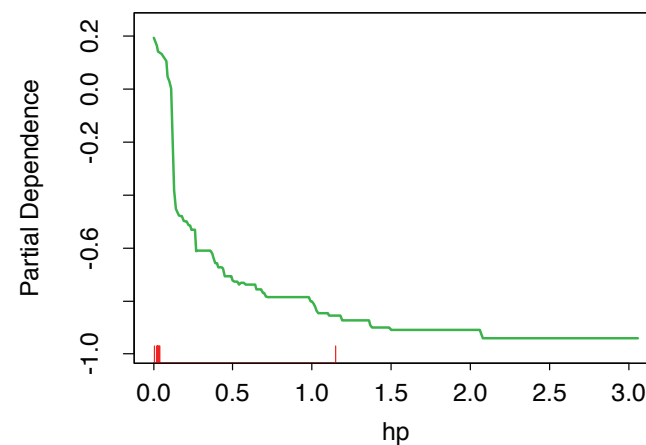
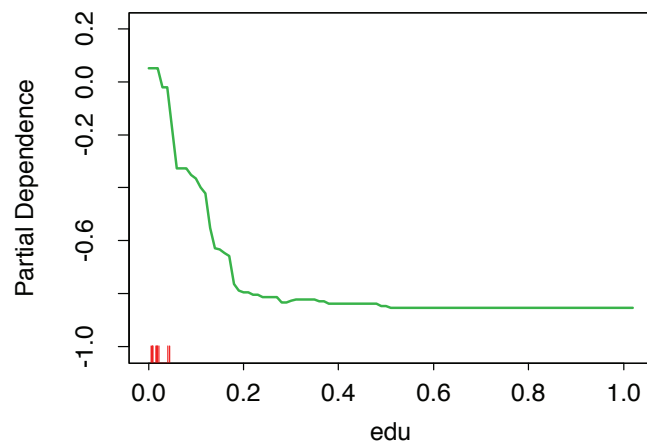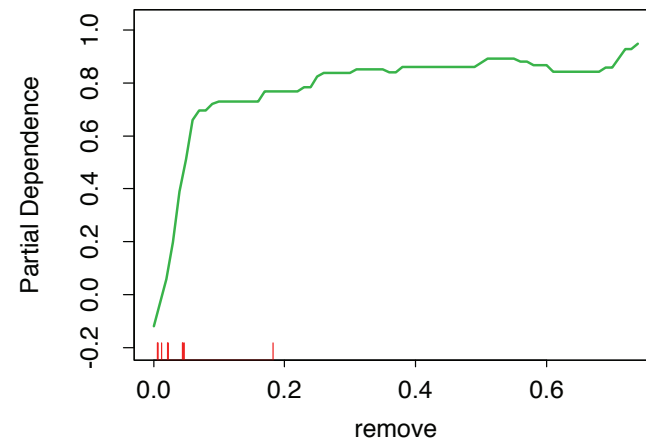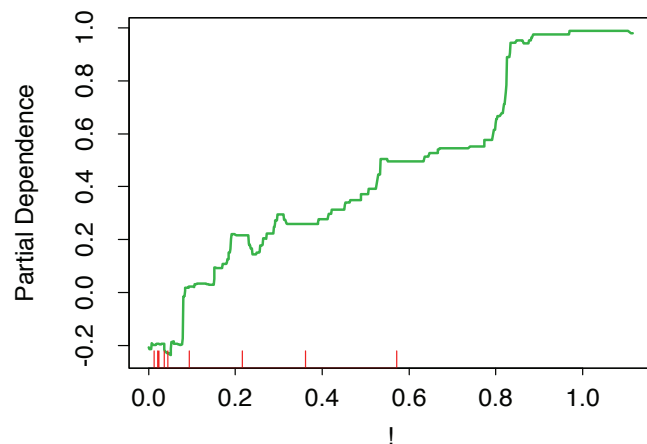$$f(x) = \log \frac{\Pr(\texttt{spam}|x)}{\Pr(\texttt{email}|x)}$$

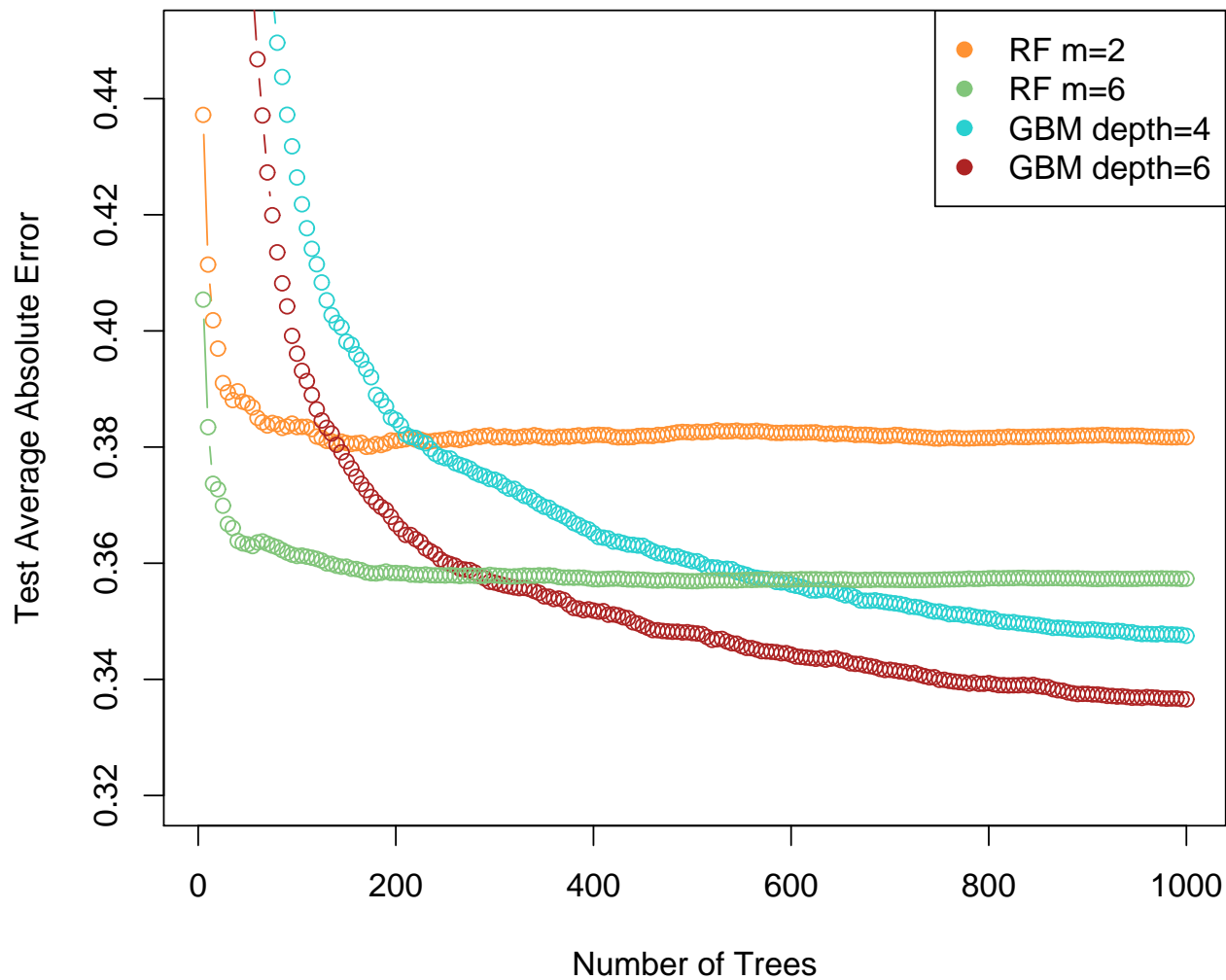using trees with $J = 6$ terminal-node trees.

Gradient Boosting achieves a test error of 4.5%, compared to 5.3% for an additive GAM, 4.75% for Random Forests, and 8.7% for CART.
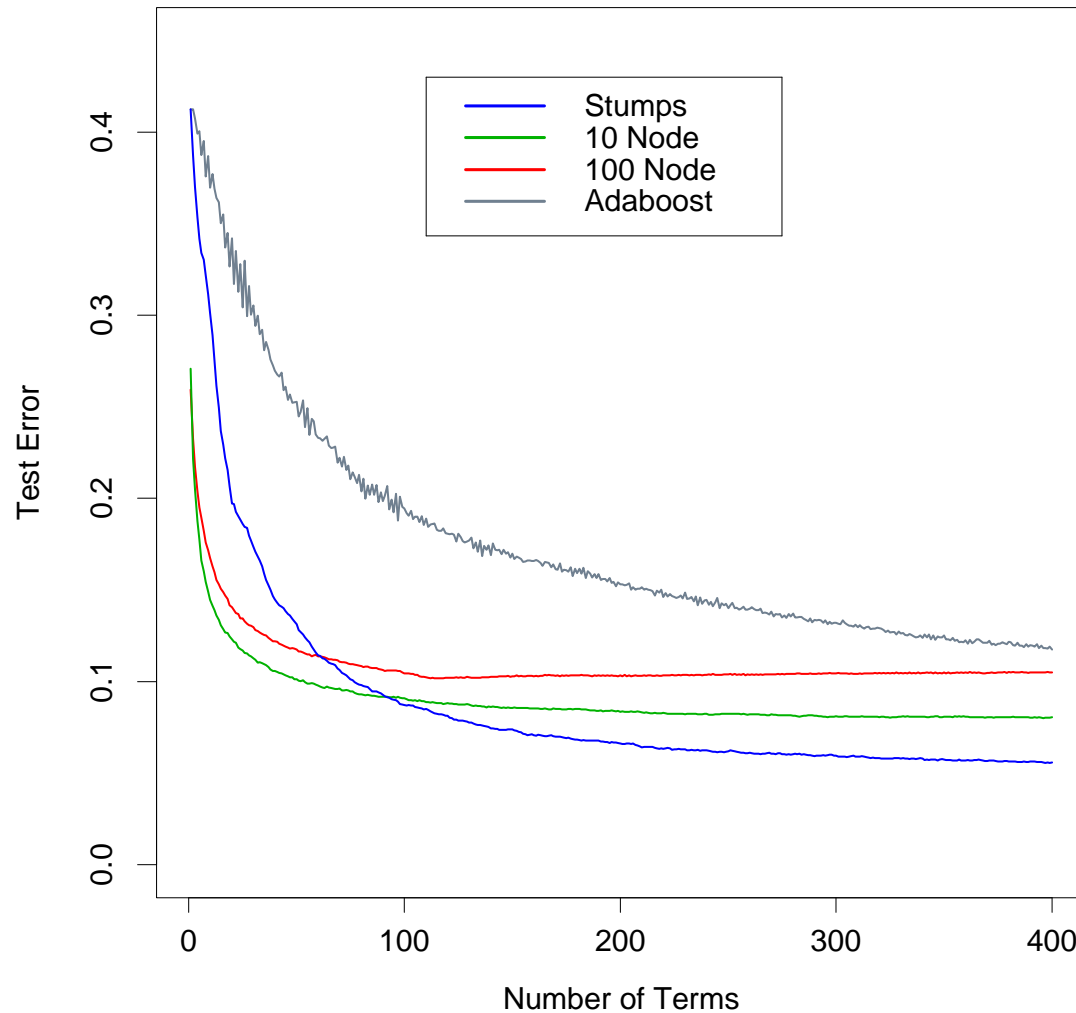
# Spam: Variable Importance

# Spam: Partial Dependence

**California Housing Data**

## Tree Size

The tree size $J$ determines the interaction order of the model:

$$
\begin{aligned}
\eta(X) \;=\; & \sum_j \eta_j(X_j) \\
& + \sum_{jk} \eta_{jk}(X_j, X_k) \\
& + \sum_{jkl} \eta_{jkl}(X_j, X_k, X_l) \\
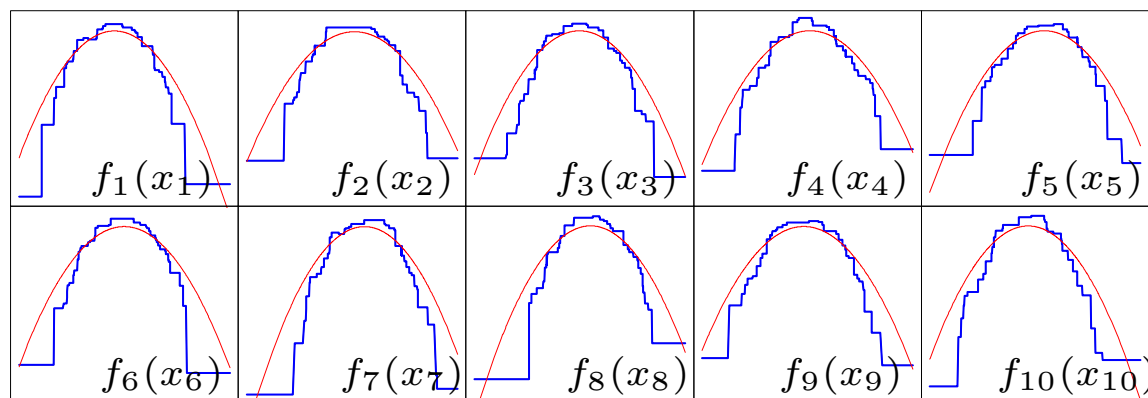& + \cdots
\end{aligned}
$$

# Stumps win!

Since the true decision boundary is the surface of a sphere, the function that describes it has the form

$$f(X) = X_1^2 + X_2^2 + \ldots + X_p^2 - c = 0.$$

Boosted stumps via Gradient Boosting returns reasonable approximations to these quadratic functions.

Coordinate Functions for Additive Logistic Trees

# Learning Ensembles

*Ensemble Learning* consists of two steps:

1. Construction of a *dictionary* $\mathcal{D} = \{T_1(X), T_2(X), \ldots, T_M(X)\}$ of basis elements (weak learners) $T_m(X)$.

2. Fitting a model $f(X) = \sum_{m \in \mathcal{D}} \alpha_m T_m(X)$.

Simple examples of ensembles

- Linear regression: The ensemble consists of the coordinate functions $T_m(X) = X_m$. The fitting is done by least squares.

- *Random Forests:* The ensemble consists of trees grown to booststrapped versions of the data, with additional randomization at each split. The fitting simply averages.

- *Gradient Boosting:* The ensemble is grown in an adaptive fashion, but then simply averaged at the end.

# Learning Ensembles and the Lasso

Proposed by Popescu and Friedman (2004):
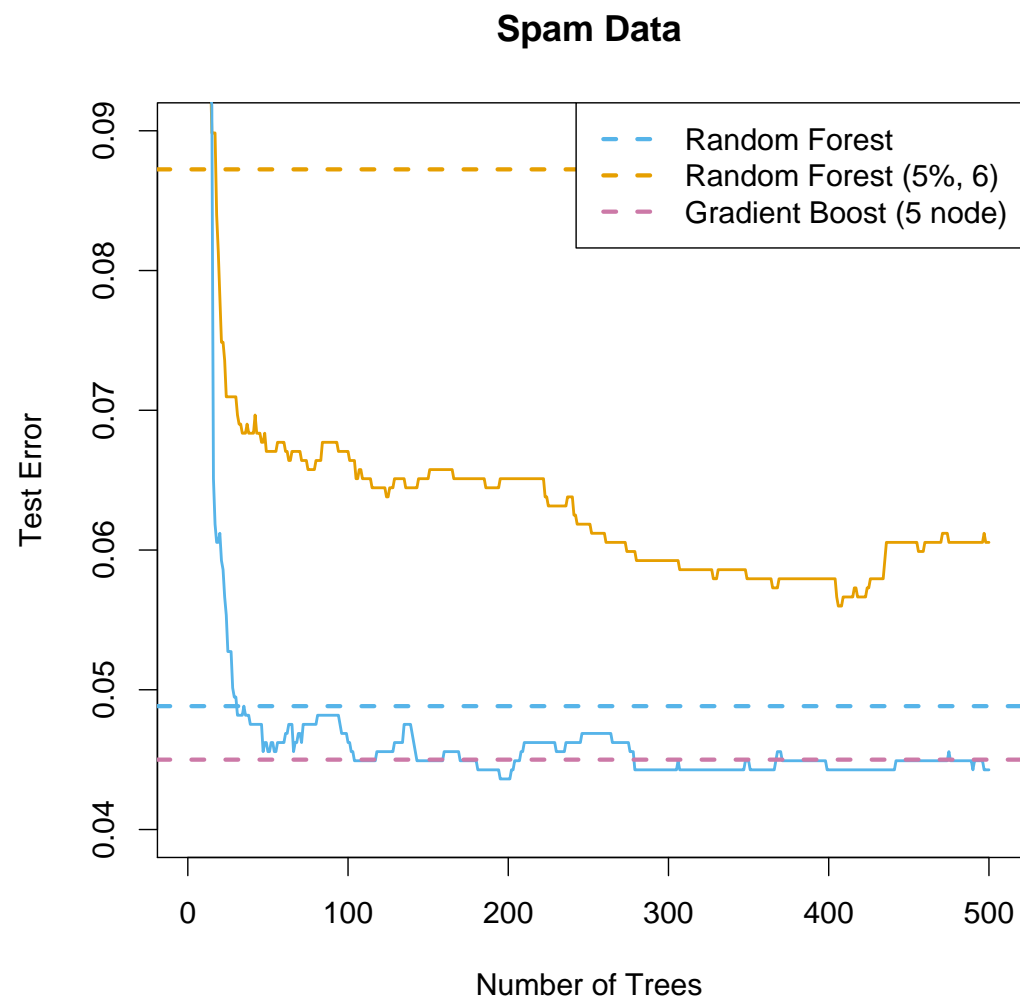
Fit the model $f(X) = \sum_{m \in \mathcal{D}} \alpha_m T_m(X)$ in step (2) using Lasso regularization:

$$\alpha(\lambda) = \arg\min_\alpha \sum_{i=1}^N L[y_i, \alpha_0 + \sum_{m=1}^M \alpha_m T_m(x_i)] + \lambda \sum_{m=1}^M |\alpha_m|.$$

Advantages:

- Often ensembles are large, involving thousands of small trees. The post-processing selects a smaller subset of these and combines them efficiently.

- Often the post-processing improves the process that generated the ensemble.

# Post-Processing Random Forest Ensembles

**Spam Data**

# Software

- *R*: free GPL statistical computing environment available from *CRAN*, implements the *S* language. Includes:

  - *randomForest*: implementation of Leo Breimans algorithms.

  - *rpart*: Terry Therneau's implementation of classification and regression trees.

  - *gbm*: Greg Ridgeway's implementation of Friedman's gradient boosting algorithm.

- *Salford Systems*: Commercial implementation of trees, random forests and gradient boosting.

- *Splus (Insightful)*: Commerical version of S.

- *Weka*: GPL software from University of Waikato, New Zealand. Includes Trees, Random Forests and many other procedures.