

Tutorial 5. Combining Estimators to Improve Performance



Abstract.

Despite the diverse pedigrees of Data Mining methods, the underlying algorithms fall into a handful of families, whose properties suggest their likely performance on a given dataset. One typically selects an algorithm by matching its strengths to the properties of one's data. Yet, performance surprises, where competing models rank differently than expected, are common; model inference, even when semi-automated, seems to yet be as much art as science.

Recently however, researchers in several fields have discovered that a simple technique — combining competing models — almost always improves classification accuracy. (Such “bundling” is a natural outgrowth of Data Mining, since much of the model search process is automated, and candidate models abound.)

This tutorial will describe an interdisciplinary collection of powerful model combination methods — including bundling, bagging, boosting, and Bayesian model averaging — and briefly demonstrate their positive effects on scientific, medical, and marketing case studies. The instructors will show why this simple, new idea will often improve a model's accuracy and stability (robustness).

About the Tutor.

John Elder is chief scientist of a data mining consulting firm in Charlottesville, Virginia (<http://www.datamininglab.com>). For fifteen years he has developed and applied adaptive, data-driven techniques to practical problems — at an engineering consulting firm, for an investment management company, at Rice University, and the University of Virginia. Dr. Elder has written and spoken widely on pattern discovery topics, is active on statistical and engineering journals and boards, and has authored some influential data mining tools. His practical experience with commercial applications — including credit scoring, direct marketing, sales forecasting, market timing, and fraud detection — helps illustrate the tutorial concepts.

Greg Ridgeway (<http://www.stat.washington.edu/greg>) is a statistician finishing his Ph.D. studies at the University of Washington. His research has focused on boosting algorithms, Monte Carlo methods, and Bayesian inference in massive datasets. His work on boosting has produced new models for survival analysis, interpretable classifiers, and systems for medical diagnosis.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD-99 Tutorial Notes San Diego CA USA
Copyright ACM 1999 1-58113-171-2/99/08...\$5.00

Combining Estimators to Improve Performance

A survey of “model bundling” techniques --
from boosting and bagging, to Bayesian model averaging
-- creating a breakthrough in the practice of Data Mining.

John F. Elder IV, Ph.D.
Elder Research, Charlottesville, Virginia
www.datamininglab.com

Greg Ridgeway
University of Washington, Dept. of Statistics
www.stat.washington.edu/greg

© 1999 Elder & Ridgeway

KDD99 T5-1

Outline

- Why combine? A motivating example
- Hidden dangers of model selection
- Reducing modeling uncertainty through *Bayesian Model Averaging*
- Stabilizing predictors through *bagging*
- Improving performance through *boosting*
- Emerging theory illuminates empirical success
- Bundling, in general
- Latest algorithms
- Closing Examples & Summary

© 1999 Elder & Ridgeway

KDD99 T5-2

Reasons to combine estimators

- Decreases variability in the predictions.
- Accounts for uncertainty in the model class.
- ☆→ Improved accuracy on new data.

© 1999 Elder & Ridgeway

KDD99 T5-3

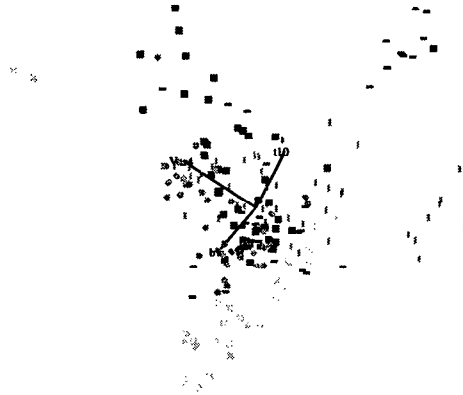
A Motivating Example: Classifying a bat's species from its chirp

- Goal: Use time-frequency features of echolocation signals to classify bats by species in the field (avoiding capture and physical inspection).
- U. Illinois biologists gathered data: 98 signals from 19 bats representing 6 species: Southeastern, Grey, Little Brown, Indiana, Pipistrelle, Big-Eared.
- ~35 data features (dimensions) calculated from signals, such as low frequency at the 3db level, time position of the signal peak, and amplitude ratio of 1st and 2nd harmonics.
- Turned out to have a nice level of difficulty for comparing methods: overlap in classes, but some separability.

© 1999 Elder & Ridgeway

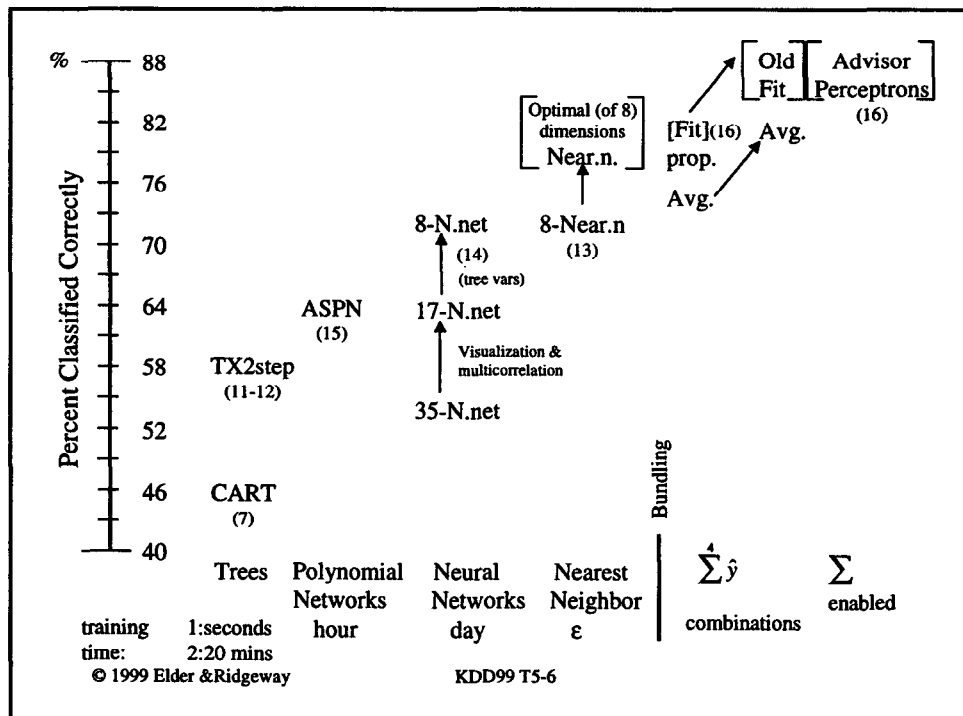
KDD99 T5-4

Sample Projection



© 1999 Elder & Ridgeway

KDD99 T5-5



What is model uncertainty?

- Suppose we wish to predict y from predictors x .
- Given a dataset of observations, D , for a new observation with predictors x^* we want to derive the predictive distribution of y^* given x^* and D .

$$P(y^* | x^*, D)$$

In practice...

- Although we want to use all the information in D to make the best estimate of y^* for an individual with covariates x^* ...

$$P(y^* | x^*, D)$$

- In practice, however, we always use

$$P(y^* | x^*, M)$$

where M is a model constructed from D .

Selecting M

- The process of selecting a model usually involves
 - Model class selection
 - Linear regression, tree regression, neural network
 - Variable selection
 - variable exclusion, transformation, smoothing
 - Parameter estimation
- We tend to choose the one model that fits the data or performs best as *the* model.

© 1999 Elder & Ridgeway

KDD99 T5-9

What's wrong with that?

- Two models may equally fit a dataset (with respect to some loss) but have different predictions.
- Competing interpretable models with equivalent performance offer ambiguous conclusions.
- Model search dilutes the evidence. “Part of the evidence is spent to specify the model.”

© 1999 Elder & Ridgeway

KDD99 T5-10

Bayesian Model Averaging

Goal: Account for model uncertainty

Method: Use Bayes' Theorem and average the models by their posterior probabilities

Properties:

- Improves predictive performance
- Theoretically elegant
- Computationally costly

© 1999 Elder & Ridgeway

KDD99 T5-11

Averaging the models

Consider a set containing the K candidate models — M_1, \dots, M_K .

With a few probability manipulations we can make predictions using all of them.

$$P(y^* | x^*, D) = \sum_k P(y^* | x^*, M_k) P(M_k | D)$$

The probability mass for a particular prediction value of y is a weighted average of the probability mass that each model places on that value of y . The weight is based on the posterior probability of that model given the data.

© 1999 Elder & Ridgeway

KDD99 T5-12

Bayes' Theorem

$$P(M_k | D) = \frac{P(D | M_k)P(M_k)}{\sum_{l=1}^K P(D | M_l)P(M_l)}$$

- M_k - model
- D - data
- $P(D|M_k)$ - integrated likelihood of M_k
- $P(M_k)$ - prior model probability

© 1999 Elder & Ridgeway

KDD99 T5-13

Challenges

- The size of the model set may cause exhaustive summation to be impossible.
- The integrated likelihood of each model is usually complex.
- Specifying a prior distribution (even a non-informative one) across the space of models is non-trivial.
- Proposed solutions to these challenges often involve MCMC, BIC approximation, MLE approximation, Occam's window, Occam's razor.

© 1999 Elder & Ridgeway

KDD99 T5-14

Performance

- Survival model: Primary biliary cirrhosis
 - BMA vs. Stepwise regression — 2% improvement
 - BMA vs. expert selected model — 10% improvement
- Linear regression: Body fat prediction
 - BMA provides best 90% predictive coverage.
- Graphical models
 - BMA yields an improvement

© 1999 Elder & Ridgeway

KDD99 T5-15

BMA References

- Chris Volinsky's BMA homepage
www.research.att.com/~volinsky/bma.html
- J. Hoeting, D. Madigan, A. Raftery, C. Volinsky (1999). "Bayesian Model Averaging: A Practical Tutorial" (to appear in *Statistical Science*),
www.stat.colostate.edu/~jah/documents/bma2.ps

© 1999 Elder & Ridgeway

KDD99 T5-16

Unstable predictors

We can always assume

$$y = f(x) + \varepsilon, \text{ where } E(\varepsilon | x) = 0$$

Assume that we have a way of constructing a predictor, $\hat{f}_D(x)$, from a dataset D .

We want to choose the estimator of f that minimizes J , squared loss for example.

$$J(\hat{f}, D) = E_{y,x} (y - \hat{f}_D(x))^2$$

© 1999 Elder & Ridgeway

KDD99 T5-17

Bias-variance decomposition

If we could average over all possible datasets, let the average prediction be

$$\bar{f}(x) = E_D \hat{f}_D(x)$$

The average prediction error over all datasets that we might see is decomposable

$$\begin{aligned} E_D J(\hat{f}, D) &= E \varepsilon^2 + E_x (f(x) - \bar{f}(x))^2 + E_{x,D} (\hat{f}_D(x) - \bar{f}(x))^2 \\ &= \text{noise} + \text{bias} + \text{variance} \end{aligned}$$

© 1999 Elder & Ridgeway

KDD99 T5-18

Bias-variance decomposition (cont.)

$$E_D J(\hat{f}, D) = E \varepsilon^2 + E_X (f(x) - \bar{f}(x))^2 + E_{X,D} (\hat{f}_D(x) - \bar{f}(x))^2$$

= noise + bias + variance

- The noise cannot be reduced.
- The squared-bias term might be reducible
- The variance term is 0 if we use

$$\hat{f}_D(x) = \bar{f}(x)$$

But this requires having an infinite number of datasets

© 1999 Elder & Ridgeway

KDD99 T5-19

Bagging (Bootstrap Aggregating)

Goal: Variance reduction.

Method: Create bootstrap replicates of the dataset and fit a model to each. Average the predictions of each model.

Properties:

- Stabilizes “unstable” methods
- Easy to implement, parallelizable
- Theory is not fully explained

© 1999 Elder & Ridgeway

KDD99 T5-20

Bagging algorithm

1. Create K bootstrap replicates of the dataset.
2. Fit a model to each of the replicates.
3. Average (or vote) the predictions of the K models.

Bootstrapping simulates the stream of infinite datasets in the bias-variance decomposition.

© 1999 Elder & Ridgeway

KDD99 T5-21

Regression results

Squared error loss

	CART	Bagging	% Reduction
Boston Housing	1.1	11.7	39%
Ozone	23.1	18.0	22%
Friedman #1	11.4	6.2	46%
Friedman #2 +3	31.8	21.7	30%
Friedman #3 -3	40.3	24.9	38%

© 1999 Elder & Ridgeway

KDD99 T5-22

Classification results

Misclassification rates

	CART	Bagging	% Reduction
Diabetes	23.4	18.8	20%
Breast	6.0	4.2	30%
Ionosphere	11.2	8.6	23%
Heart	10.0	5.3	47%
Soybean	14.5	10.6	27%
Glass	32.0	24.9	22%
Waveform	29.0	19.4	33%

© 1999 Elder & Ridgeway

KDD99 T5-23

Bagging References

- Leo Breiman's homepage
www.stat.berkeley.edu/users/breiman/
- Breiman, L. (1996) "Bagging Predictors,"
Machine Learning, 26:2, 123-140.

© 1999 Elder & Ridgeway

KDD99 T5-24

Boosting

Goal: Improve misclassification rates

Method: Sequentially fit models, each more heavily weighting those observations poorly predicted by the previous model

Properties:

- Bias and variance reduction
- Easy to implement
- Theory is not fully (but almost) explained

© 1999 Elder & Ridgeway

KDD99 T5-25

Origin of Boosting

Classification problems

$$\{\underline{X}, Y\}_i, i = 1, \dots, n$$

$$Y \in \{0, 1\}$$

The task - construct a function,

$$h(\underline{X}) : \underline{X} \rightarrow \{0, 1\}$$

so that h minimizes misclassification error.

© 1999 Elder & Ridgeway

KDD99 T5-26

Generic boosting algorithm

Equally weight the observations $(\underline{X}, Y)_i$

For t in $1, \dots, T$

Using the weights, fit a classifier $h_t(\underline{X}) \rightarrow Y$

Upweight the poorly predicted observations

Downweight the well-predicted observations

Merge h_1, \dots, h_T to form the boosted classifier

© 1999 Elder & Ridgeway

KDD99 T5-27

AdaBoost algorithm

Freund & Schapire 1996

$(X, Y)_i$ where $Y_i \in \{0, 1\}$, $w_i^{(1)} = \frac{1}{N}$

- With weights, fit the model $H_t(x_i) : X \rightarrow [0, 1]$.
- Compute the error $\varepsilon_t = \sum_{i=1}^N w_i^{(t)} |y_i - H_t(x_i)|$
- Reweight

$$w_i^{(t+1)} = w_i^{(t)} \beta_t^{1-|y_i - H_t(x_i)|} \quad \beta_t = \frac{\varepsilon_t}{1 - \varepsilon_t}$$

Lastly, predict

$$H(x) = \frac{1}{1 + \prod_{i=1}^T \beta_i^{2r(x)-1}} \quad r(x) = \frac{\sum_{i=1}^T (\log \frac{1}{\beta_i}) H_i(x)}{\sum_{i=1}^T (\log \frac{1}{\beta_i})}$$

© 1999 Elder & Ridgeway

KDD99 T5-28

AdaBoost's Performance

Freund & Schapire [1996]

- Leo Breiman - AdaBoost with trees is the “best off-the-shelf classifier in the world.”
- Performs well with many base classifiers and in a variety of problem domains.
- AdaBoost is generally slow to overfit.
- Boosted naïve Bayes tied for first place in the 1997 KDD Cup. (Elkan [1997])
- Boosted naïve Bayes is a scalable, interpretable classifier (Ridgeway, *et al* [1998]).

© 1999 Elder & Ridgeway

KDD99 T5-29

Boosting as optimization

- Friedman, Hastie, Tibshirani [1998] - AdaBoost is an optimization method for finding a classifier.
- Let $y \in \{-1, 1\}$, $F(x) \in (-\infty, \infty)$

$$J(F) = E\left(e^{-yF(x)} \mid x\right)$$

© 1999 Elder & Ridgeway

KDD99 T5-30

Criterion

- $E(e^{-yF(x)})$ bounds the misclassification rate.

$$I(yF(x) < 0) < e^{-yF(x)}$$

- The minimizer of $E(e^{-yF(x)})$ coincides with the maximizer of the expected Bernoulli likelihood.

$$E(\ell(p(x), y)) = -E \log(1 + e^{-2yF(x)})$$

Optimization step

$$J(F + f) = E(e^{-y(F(x) + f(x))} | x)$$

- Select f to minimize $J...$

$$F^{(t+1)} \leftarrow F^{(t)} + \frac{1}{2} \log \frac{E_w[I(y=1) | x]}{1 - E_w[I(y=1) | x]}$$

$$w(x, y) = e^{-yF^{(t)}(x)}$$

LogitBoost

Friedman, Hastie, Tibshirani [1998]

- Logistic regression

$$y = \begin{cases} 1 & \text{with probability } p(x) \\ 0 & \text{with probability } 1 - p(x) \end{cases}$$

$$p(x) = \frac{1}{1 + e^{-F(x)}}$$

- Expected log-likelihood of a regressor, $F(x)$

$$E \ell(F) = E(yF(x) - \log(1 + e^{F(x)}) | x)$$

© 1999 Elder & Ridgeway

KDD99 T5-33

Newton steps

$$J(F + f) = E(y(F(x) + f(x)) - \log(1 + e^{F(x) + f(x)}) | x)$$

- Iterate to optimize expected log-likelihood.

$$F^{(t+1)}(x) \leftarrow F^{(t)}(x) - \frac{\frac{\partial}{\partial f} J(F^{(t)} + f) \Big|_{f=0}}{\frac{\partial^2}{\partial f^2} J(F^{(t)} + f) \Big|_{f=0}}$$

© 1999 Elder & Ridgeway

KDD99 T5-34

LogitBoost, continued

- Newton steps for Bernoulli likelihood

$$F(x) \leftarrow F(x) + E_w \left(\frac{y - p(x)}{p(x)(1 - p(x))} \middle| x \right)$$
$$w(x) = p(x)(1 - p(x))$$

- In practice the $E_w(\bullet|x)$ can be any regressor - trees, smoothers, etc.
- Trees are adaptive and work well for high dimensional data.

© 1999 Elder & Ridgeway

KDD99 T5-35

Classification results

Friedman, Hastie, Tibshirani [1998]

	CART	AdaBoost	LogitBoost
Breast	4.5%	4.0%	2.9%
Ion	7.6%	6.8%	7.1%
Glass	40.0%	25.7%	26.6%
Sonar	59.6%	20.2%	20.2%
Waveform	36.4%	19.5%	20.6%

© 1999 Elder & Ridgeway

KDD99 T5-36

Boosting References

- Rob Schapire's homepage
www.research.att.com/~schapire
- Freund, Y. and R. Schapire (1996). "Experiments with a new boosting algorithm," Machine Learning: Proceedings of the 13th International Conference, 148-156.
- Jerry Friedman's homepage
www.stat.stanford.edu/~jhf
- Friedman, J., T. Hastie, R. Tibshirani (1998). "Additive Logistic Regression: a statistical view of boosting," Technical report, Statistics Department, Stanford University.

© 1999 Elder & Ridgeway

KDD99 T5-37

In general, combining ("bundling") estimators consists of two steps:

- 1) Constructing varied models, and
- 2) Combining their estimates

Generate component models by varying:

- Case Weights
- Data Values
- Guiding Parameters
- Variable Subsets

Combine estimates using:

- Estimator Weights
- Voting
- Advisor Perceptrons
- Partitions of Design Space, X

© 1999 Elder & Ridgeway

KDD99 T5-38

Other Bundling Techniques

We've Examined:

- **Bayesian Model Averaging:** sum estimates of possible models, weighted by posterior evidence
- **Bagging** (Breiman 96) (*bootstrap aggregating*) -- bootstrap data (to build trees mostly); take majority vote or average
- **Boosting** (Freund & Shapire 96) -- weight error cases by $\beta\tau = (1-e(t))/e(t)$, iteratively re-model; average, weighing model t by $\ln(\beta\tau)$

Additional Example Techniques:

- **GMDH** (Ivakhenko 68) -- multiple layers of quadratic polynomials, using two inputs each, fit by Linear Regression
- **Stacking** (Wolpert 92) -- train a 2nd-level (LR) model using leave-1-out estimates of 1st-level (neural net) models
- **ARCing** (Breiman 96) (Adaptive Resampling and Combining) -- Bagging with reweighting of error cases; similar to boosting
- **Bumping** (Tibshirani 97) -- bootstrap, select single best
- **Crumpling** (Anderson & Elder 98) -- average cross-validations
- **Born-Again** (Breiman 98) -- invent new X data...

© 1999 Elder & Ridgeway

KDD99 T5-39

When does Bundling work?

Hypotheses:

- Breiman (1996): when the prediction method is *unstable* (significantly different models are constructed)
- Ali & Pazzani (1996): when there is low noise, lots of irrelevant variables, and good individual predictors which make different errors
- when models are slightly overfit
- when models are from different families

© 1999 Elder & Ridgeway

KDD99 T5-40

Advanced techniques

- Stochastic gradient boosting
- Adaptive bagging
- Example regression and classification results

© 1999 Elder & Ridgeway

KDD99 T5-41

Stochastic Gradient Boosting

Goal: Non-parametric function estimation

Method: Cast the problem as optimization and use gradient ascent to obtain predictor

Properties:

- Bias and variance reduction
- Widely applicable
- Can make use of existing algorithms
- Many tuning parameters

© 1999 Elder & Ridgeway

KDD99 T5-42

Improving boosting

- Boosting usually has the form

$$F^{(t+1)}(x) \leftarrow F^{(t)}(x) + \lambda E_w(z(y, x)|x)$$

Improve by...

- Sub-sampling a fraction of the data at each step when computing the expectation.
- “Robustifying” the expectation.
- Trimming observations with small weights.

© 1999 Elder & Ridgeway

KDD99 T5-43

Stochastic gradient boosting offers...

- Application to likelihood based models (GLM, Cox models)
- Bias reduction - non-linear fitting
- Massive datasets - bagging, trimming
- Variance reduction - bagging
- Interpretability - additive models
- High-dimensional regression - trees
- Robust regression

© 1999 Elder & Ridgeway

KDD99 T5-44

SGB References

- Friedman, J. (1999). "Greedy function approximation: a gradient boosting machine," Technical report, Dept. of Statistics, Stanford University.
- Friedman, J. (1999). "Stochastic gradient boosting," Technical report, Dept. of Statistics, Stanford University.

Adaptive Bagging

Goal: Bias and variance reduction

Method: Sequentially fit *bagged* models,
where each fits the current residuals

Properties:

- Bias and variance reduction
- No tuning parameters

Adaptive bagging algorithm

1. Fit a bagged regressor to the dataset D .
2. Predict “out-of-bag” observations.
3. Fit a new bagged regressor to the bias (error) and repeat.

For a new observation, sum the predictions from each stage.

© 1999 Elder & Ridgeway

KDD99 T5-47

Regression results

Squared error loss

	Bagging	Debias	% Reduction
Boston Housing	12.7	10.8	14%
Ozone	17.8	17.8	0%
Servo -2	24.5	25.1	-3%
Abalone	4.9	4.9	0%
Robot arm -2	4.7	2.8	41%
Peak20	12.8	3.7	71%
Friedman #1	6.3	4.1	35%
Friedman #2 +3	21.5	21.5	0%
Friedman #3 -3	24.8	24.8	0%

© 1999 Elder & Ridgeway

KDD99 T5-48

Classification results

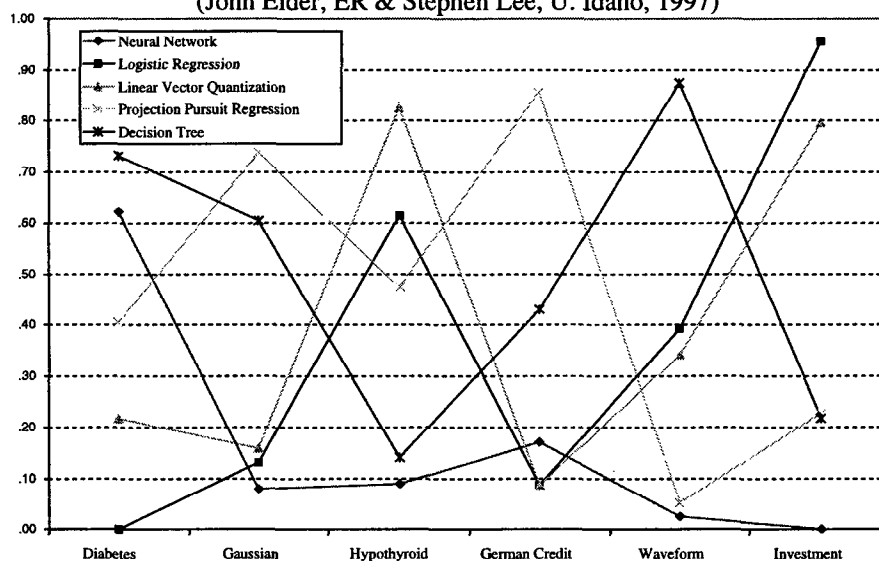
Misclassification rates

	Debias	AdaBoost
Diabetes	23.4	26.6
Breast	3.9	3.2
Ionosphere	6.6	6.4
Sonar	14.1	15.6
Heart	15.6	20.7
German credit	23.6	23.5
Votes	3.7	5.4
Liver	25.9	28.7

© 1999 Elder & Ridgeway

KDD99 T5-49

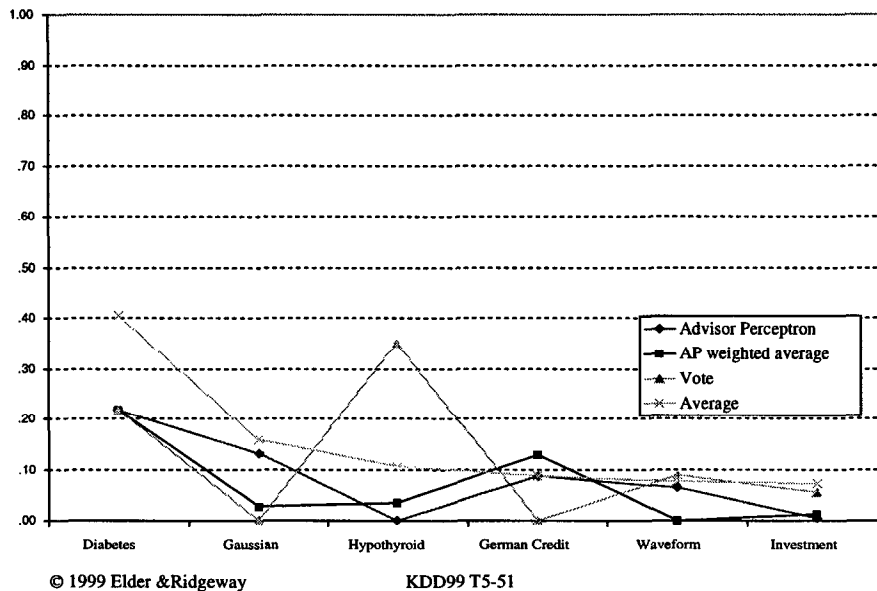
Relative Performance Examples: 5 Algorithms on 6 Datasets
(John Elder, ER & Stephen Lee, U. Idaho, 1997)



© 1999 Elder & Ridgeway

KDD99 T5-50

Essentially every Bundling method improves performance

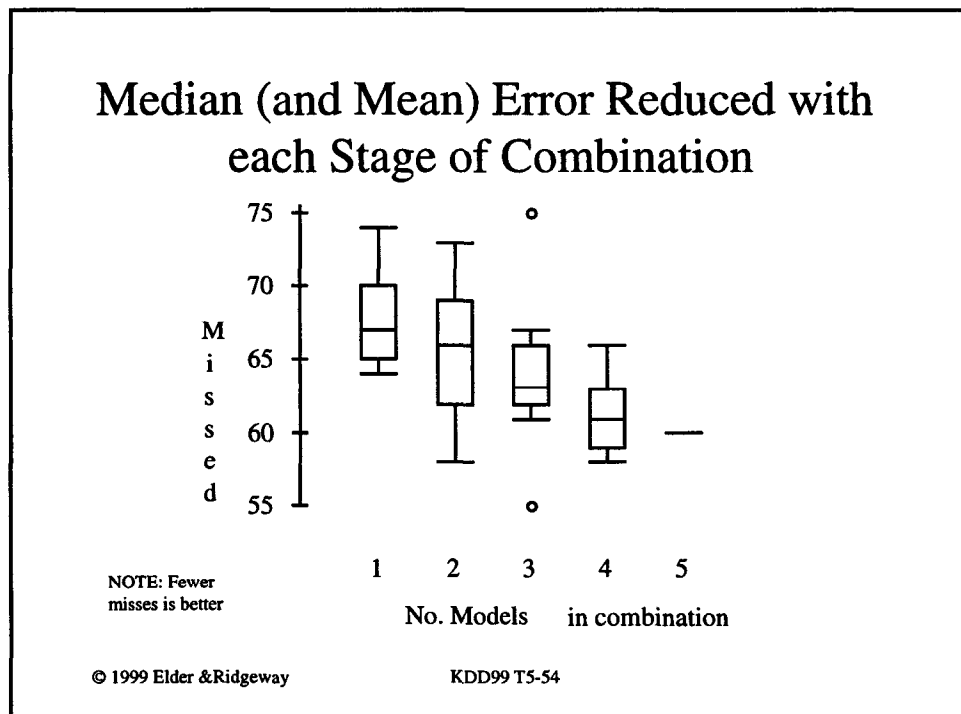
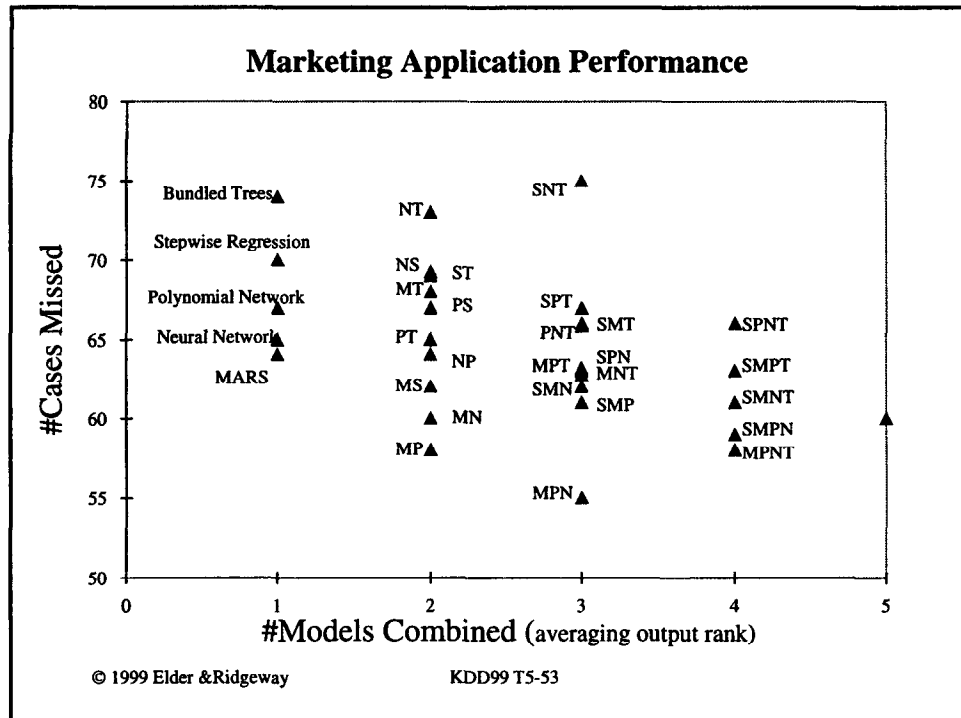


Application Ex.: Direct Marketing (Elder Research 1996-)

- Model respondents to direct marketing as binary variable: 0 (no response), 1 (response).
- Create models using several (here, 5) different algorithms, all employing the same candidate model inputs.
- Rank-order model responses:
 - Give highest-probability response value a rank of 1, second highest value 2, etc.
 - For bundling, combine model ranks (not estimates) into a new consensus estimate (which is again ranked).
- Report number of response cases missed (in top portion).

© 1999 Elder & Ridgeway

KDD99 T5-52



...and in a multitude of counselors there is safety.

Proverbs 24:6b

Why Bundling works

- (semi-) Independent Estimators
- Bayes Rule - weighing evidence
- Shrinking (ex.: stepwise LR)
- Smoothing (ex.: decision trees)
- Additive modeling and maximum likelihood
(Friedman, Hastie, & Tibshirani 8/20/98)

... Open research area.

Meanwhile, we recommend bundling competing candidate models both within, and between, model families.

© 1999 Elder & Ridgeway

KDD99 T5-55