

Ensemble Methods in Data Mining

Andrew Fast, Ph.D.
Chief Scientist
Elder Research, Inc.
fast@datamininglab.com

John F. Elder, IV, Ph.D.
President and CEO
Elder Research, Inc.
elder@datamininglab.com



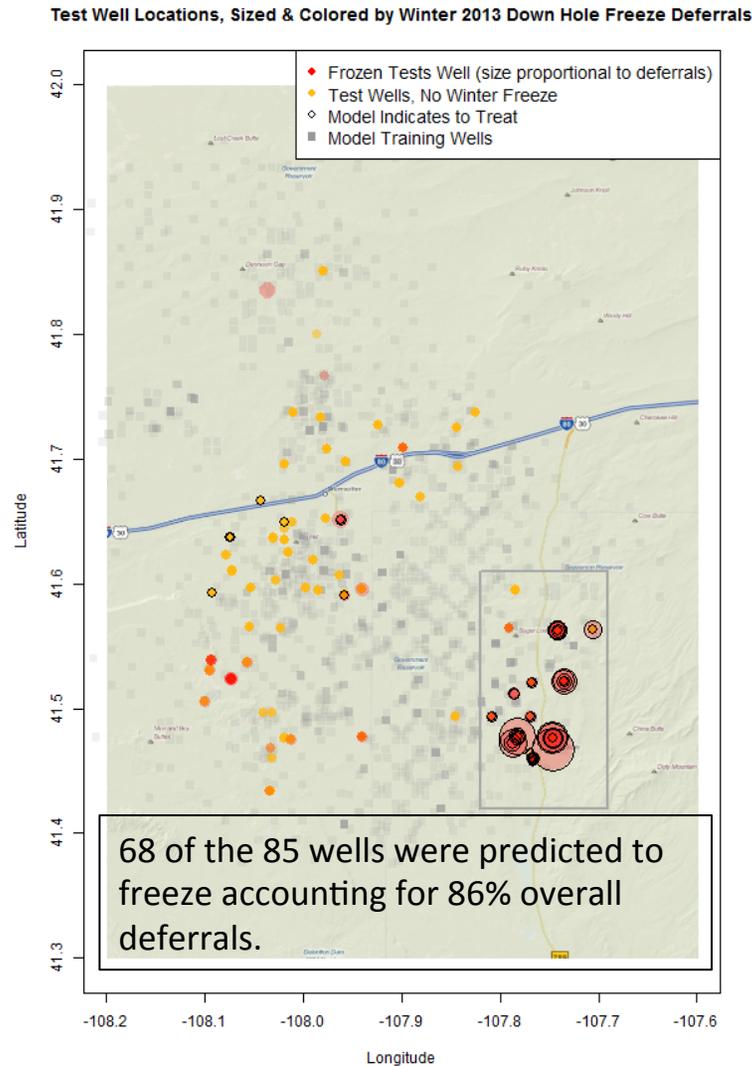
*Conference on Applied Statistics in Defense (CASD)
Washington, D.C., October 2014*

Occam's Razor and Induction

- “Entities should not be multiplied beyond necessity”
- Prefer parsimonious models when learning from data to increase performance on new, unseen data
- *Paradox of ensembles:* Combined models perform better than single models



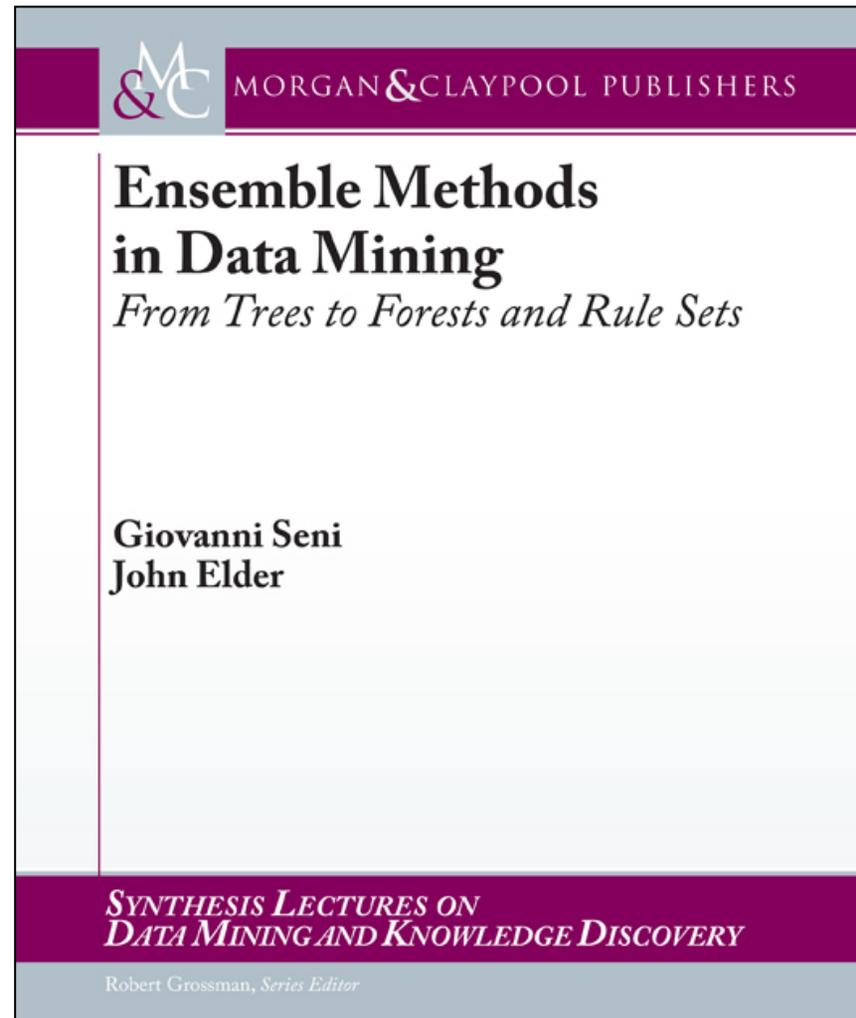
Ensembles in Action



- Work for a major petrochemical company
- Focusing on well-stoppages called “freezes”
 - Lost production called “deferrals”
- Can we provide a ranked list of wells most likely to “freeze”?

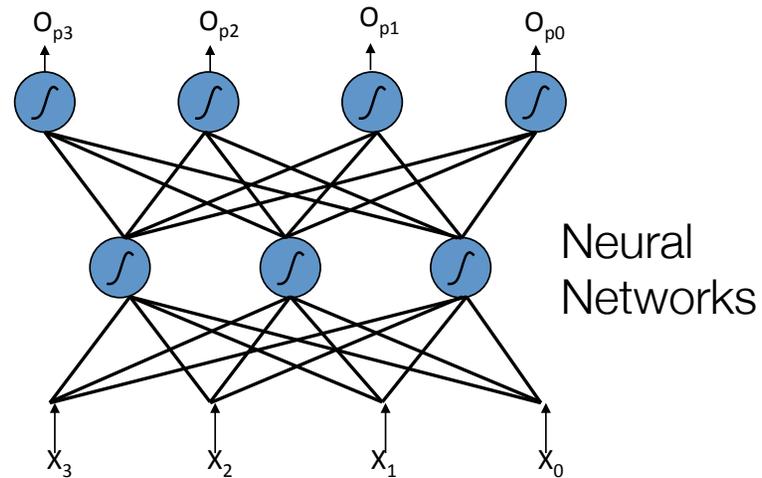
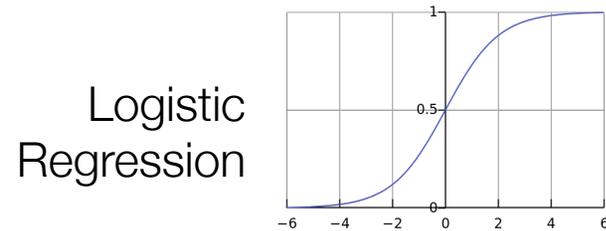
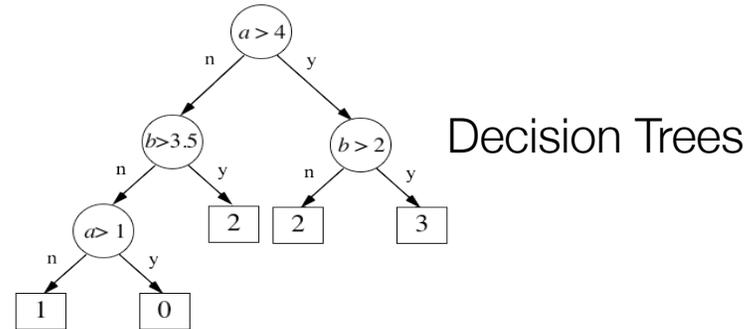
Ensemble Methods in Data Mining

- Giovanni Seni & John Elder
- First published in 2010
 - 2nd Edition coming soon
- Complete exposition of the Importance Sampling Learning Ensembles (ISLE) approach



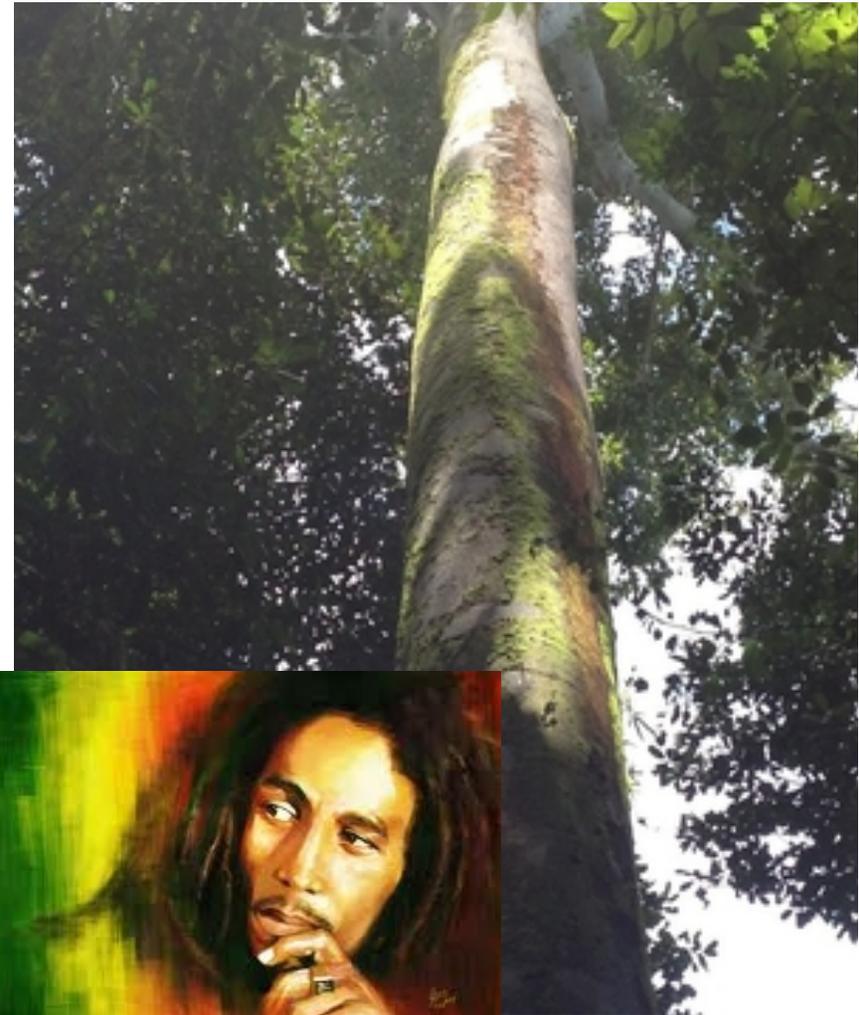
Distinctives of Inductive Modeling

- Observational data only, no experimental control
- Learning *structure* of model, not just fitting parameters
- Both Classification and Regression tasks



Challenge: Overfit

- Vast Search Effect
 - Many models considered
 - Can look good by chance alone
- Most inductive algorithms will learn structure from random data
- Overfit models do not generalize well to unseen data



<http://www.bypolaropposite.com/whats-hat-health/2014/10/8/pareidolia-that-thing-has-a-face-i-saw-bob-marley-in-a-tree>

Solution: Penalize Model Complexity

1. Explicit Complexity Constraints

- Ridge, Lasso, BIC, etc.

2. Incremental model building

- Stepwise Regression

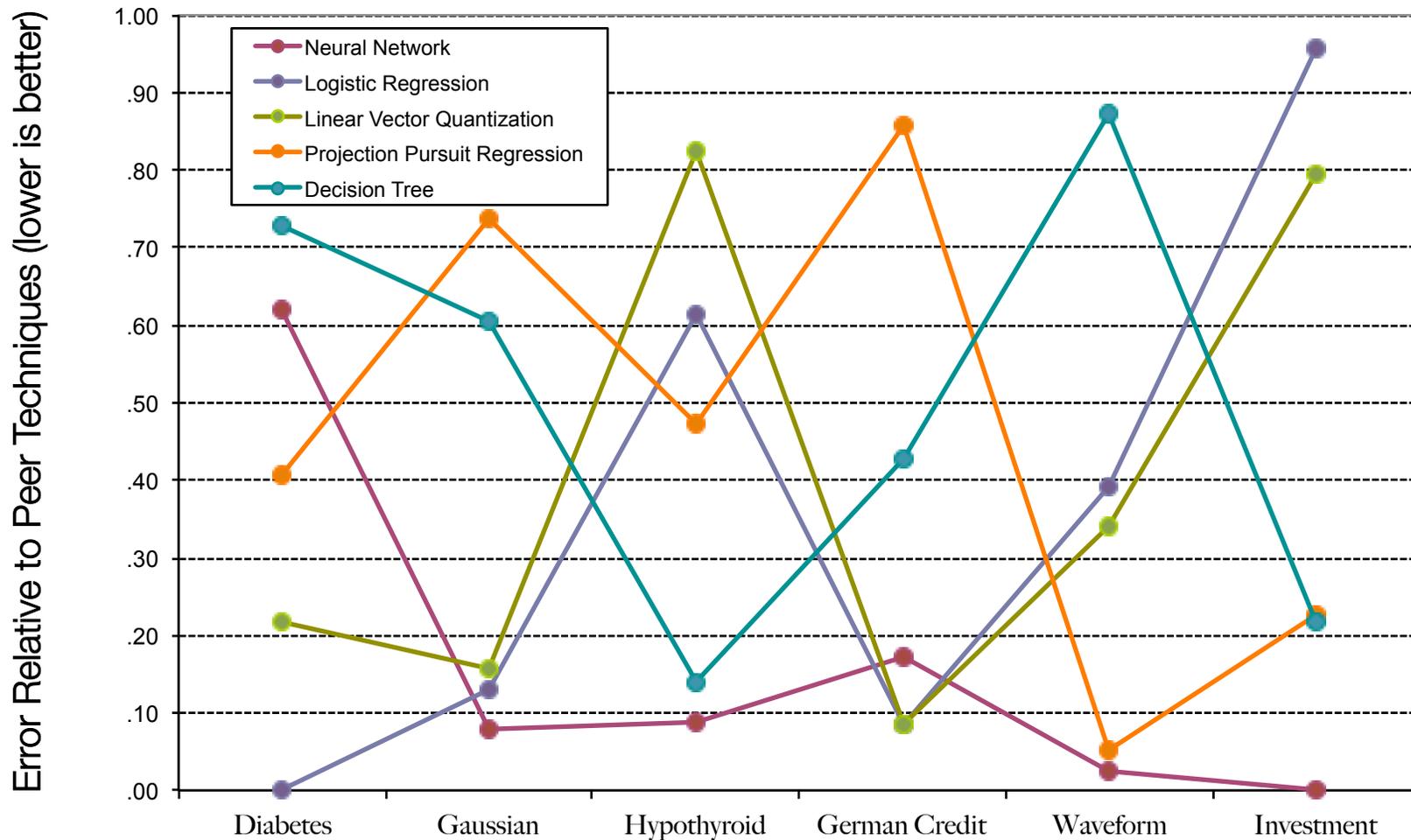
3. Robust Loss Functions

- Huber Loss

Regularization is “...any part of model building which takes into account – implicitly or explicitly – the finiteness and imperfection of the data and the limited information in it, which we can term ‘variance’ in an abstract sense.”

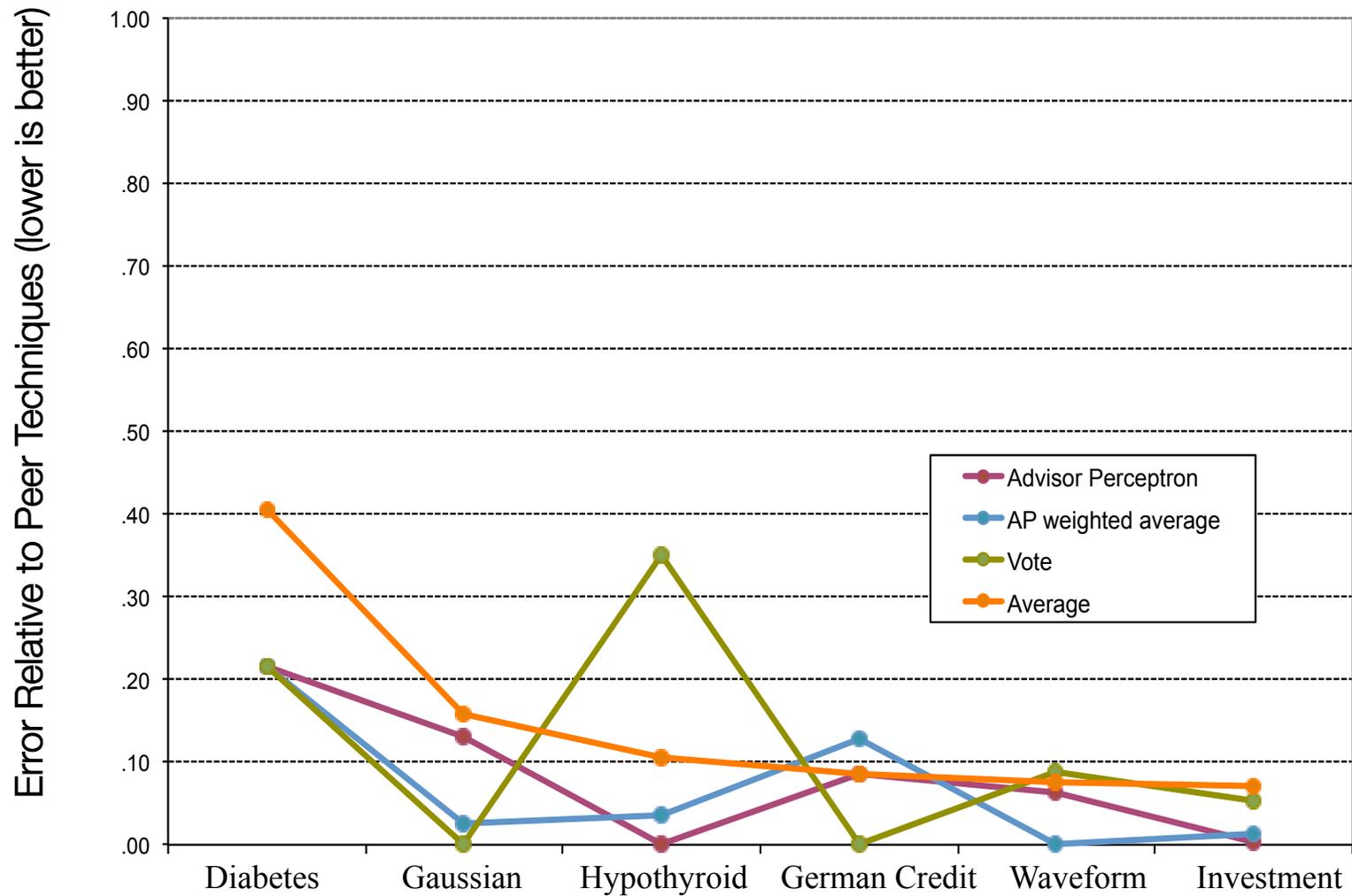
- Saharon Rosset (2003), *Ph.D. Thesis*

Single Model Performance



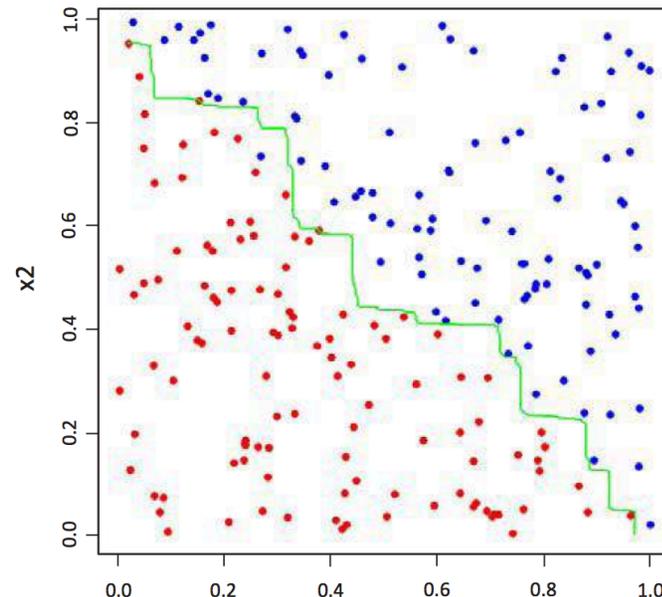
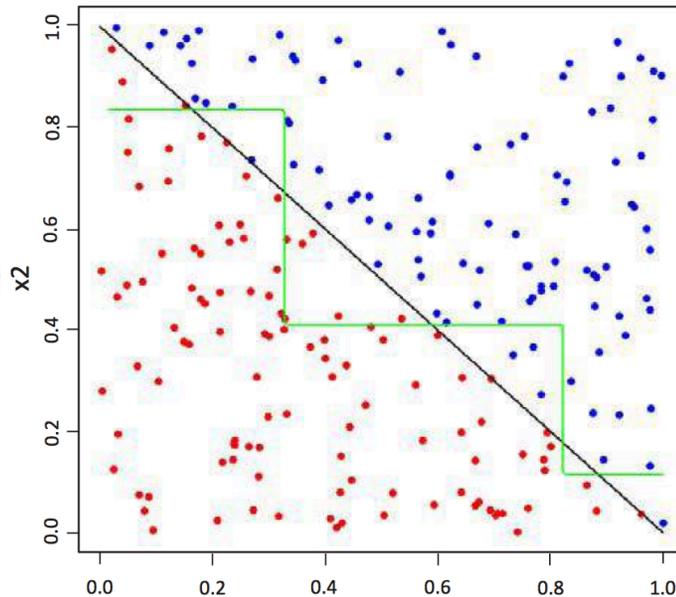
Comparison of 5 Algorithms on 6 Datasets (from Elder and Lee, 1997)

“Bundling” Models Improve Performance



Traditional Ensembles

- Combine multiple models together into a single result
- Usually leads to improved accuracy and generalization
- Techniques include:
 - Boosting
 - Bagging
 - Random Forests
 - Bayesian Model Averaging
 - Stacking



Importance Sampling Learning Ensembles

$$\{\hat{c}_m, \hat{p}_m\}_0^M = \min_{\{c_m, p_m\}_0^M} \sum_{i=1}^N L \left(y_i, c_0 + \sum_{m=1}^M c_m T_m(x; p_m) \right)$$

Minimize the loss function, L

Model Weights

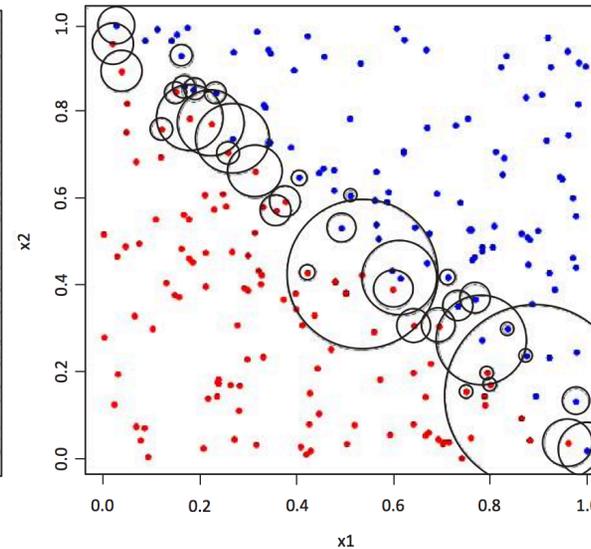
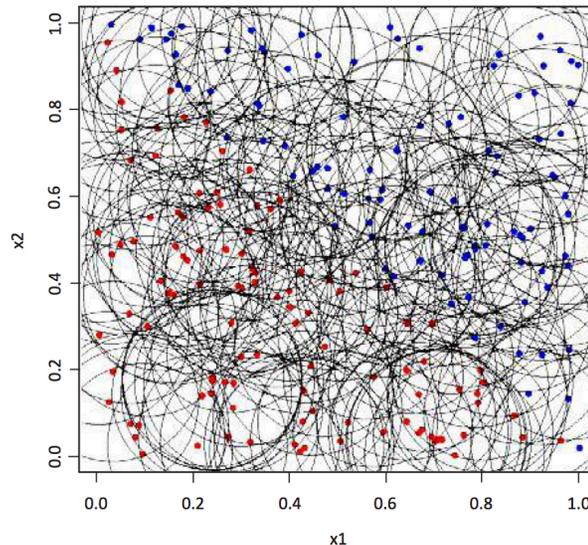
Collection of base learners (basis functions)

- “The Ensemble is a linear model in a (very) high dimensional space of derived variables” (Seni & Elder, pg 39)
- Must sample from space of all possible learners
 - Need to balance breadth and quality
- Difficult Optimization problem!
- Use Perturbation Sampling to find a solution
 - Perturb Data distribution
 - Perturb Loss Function
 - Perturb Search algorithm

Example: Gradient Boosting

- Boosting reweights the data inputs based on error (higher error means more weight)
- Generalization of AdaBoost to any differential loss function
- Gradient Boosting is doing shrinkage implicitly with a Lasso-type penalty (Seni & Elder, pg 61)

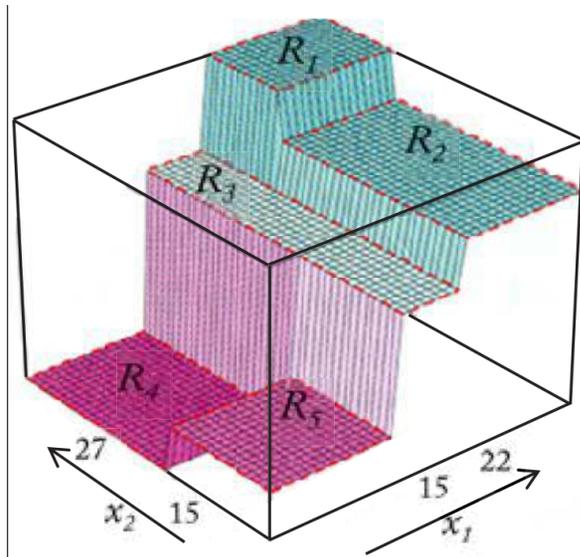
Before: All instances have the same weight



After: Only instances along the decision boundary have non-zero weight

Example: Rule Ensembles

- Ensembles lose the explainability of a single model
- Rule Ensembles provide improved explainability but with accuracy of the ensemble
- Convert trees to conjunctive rules
- Can be combined with other base learner functions
 - e.g., add linear terms



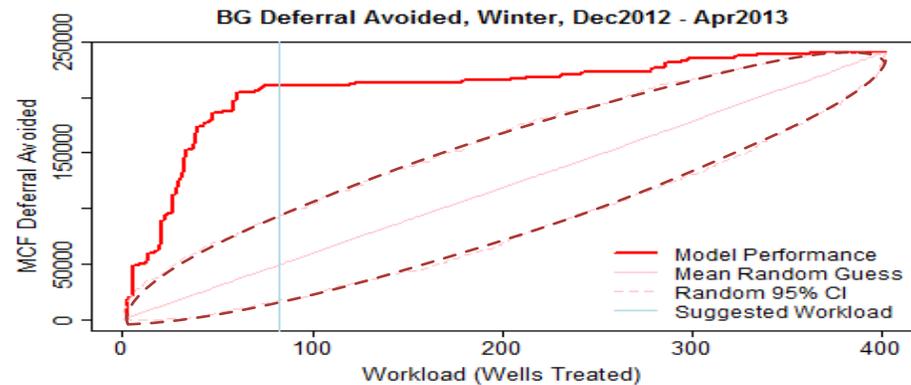
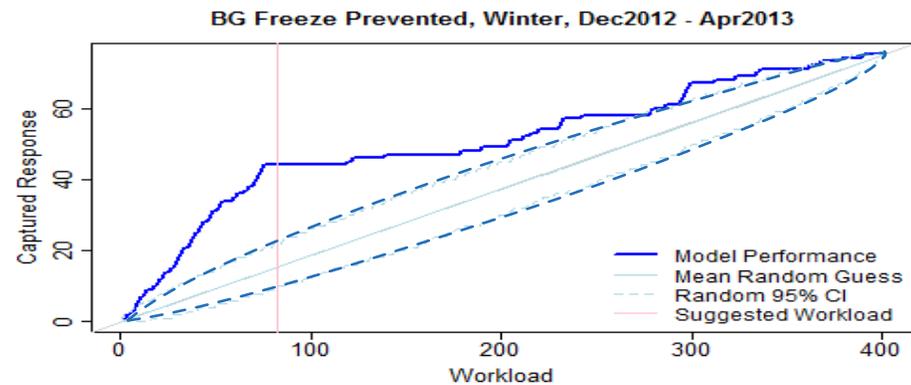
Region	Rule
R_1	$r_1(\mathbf{x}) = I(x_1 > 22) \cdot I(x_2 > 27)$
R_2	$r_2(\mathbf{x}) = I(x_1 > 22) \cdot I(0 \leq x_2 \leq 27)$
R_3	$r_3(\mathbf{x}) = I(15 < x_1 \leq 22) \cdot I(0 \leq x_2)$
R_4	$r_4(\mathbf{x}) = I(0 \leq x_1 \leq 15) \cdot I(x_2 > 15)$
R_5	$r_5(\mathbf{x}) = I(0 \leq x_1 \leq 15) \cdot I(0 \leq x_2 \leq 15)$

[Friedman and Popescu, 2005; Friedman and Bogdan, 2008]

Graphic adapted from Hastie et al., 2001

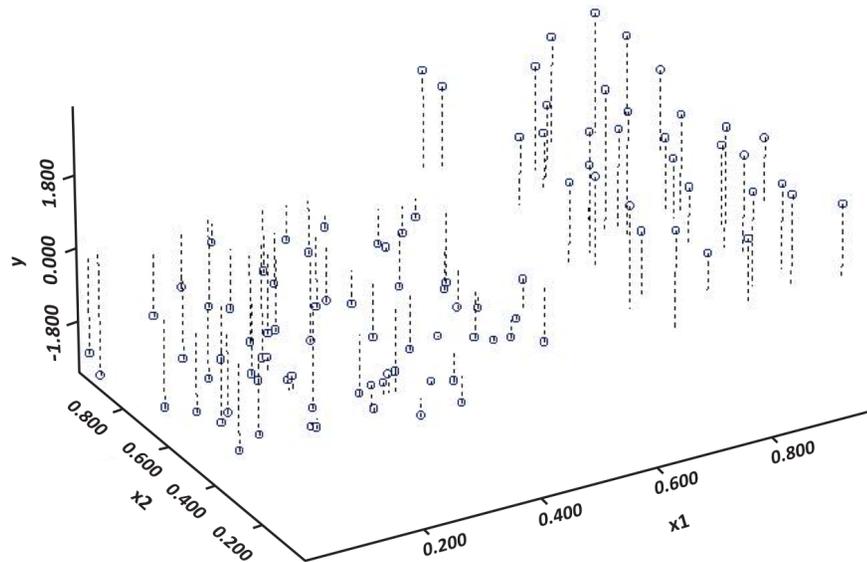
Well Maintenance Results

- Used Gradient Boosted Trees to predict which wells would require maintenance
- Dotted lines indicate the 95% percent confidence interval using permutation testing

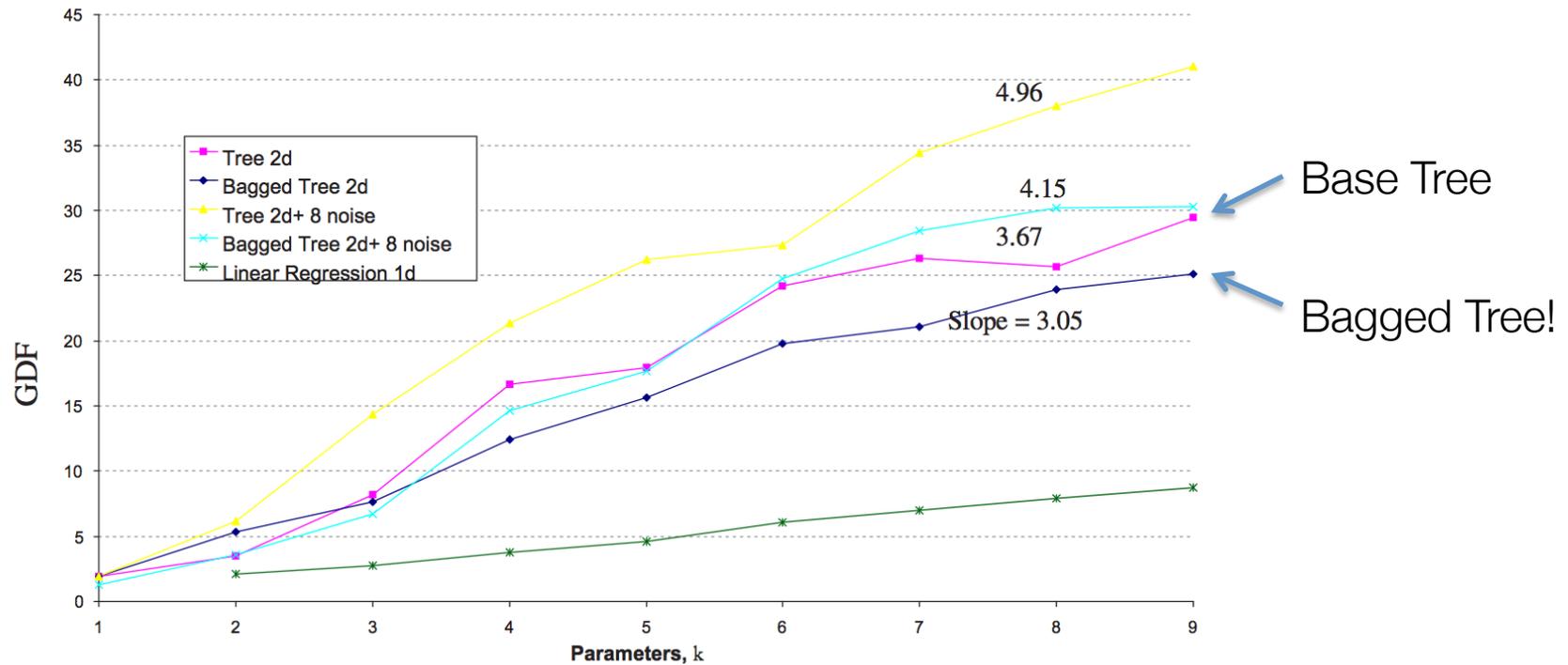


Ensembles and Complexity

- Traditionally measured by model size (number of model parameters).
 - Clearly this fails for Ensembles
- Generalized Degrees of Freedom (Ye 1998)
 - Empirical measure of the flexibility of the entire modeling process.
- True degrees of freedom can be greater than the number of parameters.



Generalized Degrees of Freedom



- Bagged Tree has fewer degrees of freedom than a single tree with the same number of splits

Summary

- Ensembles are a powerful tool for obtaining accurate results on real world problems.
- Despite multiple models, *flexibility* of the modeling process is actually lower in ensembles than a single model due to an implicit regularization process
- R packages and examples available in the book.
 - Increased adoption in COTS software, too

Contact Information

Andrew Fast, Ph.D.
Chief Scientist

fast@datamininglab.com
(434) 973-7673
www.datamininglab.com



Andrew Fast Chief Scientist, Elder Research, Inc.



DR. ANDREW FAST LEADS RESEARCH IN TEXT MINING AND SOCIAL NETWORK ANALYSIS AT ELDER RESEARCH, THE NATION'S LEADING DATA MINING CONSULTANCY. ERI WAS FOUNDED IN 1995 AND HAS OFFICES IN CHARLOTTESVILLE VA AND WASHINGTON DC, (WWW.DATAMININGLAB.COM). ERI FOCUSES ON FEDERAL, COMMERCIAL, INVESTMENT, AND SECURITY APPLICATIONS OF ADVANCED ANALYTICS, INCLUDING STOCK SELECTION, IMAGE RECOGNITION, BIOMETRICS, PROCESS OPTIMIZATION, CROSS-SELLING, DRUG EFFICACY, CREDIT SCORING, RISK MANAGEMENT, AND FRAUD DETECTION.

DR. FAST GRADUATED MAGNA CUM LAUDE FROM BETHEL UNIVERSITY AND EARNED MASTER'S AND PH.D. DEGREES IN COMPUTER SCIENCE FROM THE UNIVERSITY OF MASSACHUSETTS AMHERST. THERE, HIS RESEARCH FOCUSED ON CAUSAL DATA MINING AND MINING COMPLEX RELATIONAL DATA SUCH AS SOCIAL NETWORKS. AT ERI, ANDREW LEADS THE DEVELOPMENT OF NEW TOOLS AND ALGORITHMS FOR DATA AND TEXT MINING FOR APPLICATIONS OF CAPABILITIES ASSESSMENT, FRAUD DETECTION, AND NATIONAL SECURITY.

DR. FAST HAS PUBLISHED ON AN ARRAY OF APPLICATIONS INCLUDING DETECTING SECURITIES FRAUD USING THE SOCIAL NETWORK AMONG BROKERS, AND UNDERSTANDING THE STRUCTURE OF CRIMINAL AND VIOLENT GROUPS. OTHER PUBLICATIONS COVER MODELING PEER-TO-PEER MUSIC FILE SHARING NETWORKS, UNDERSTANDING HOW COLLECTIVE CLASSIFICATION WORKS, AND PREDICTING PLAYOFF SUCCESS OF NFL HEAD COACHES (WORK FEATURED ON ESPN.COM). WITH JOHN ELDER AND OTHER CO-AUTHORS, ANDREW HAS WRITTEN A BOOK ON PRACTICAL TEXT MINING, THAT WAS AWARDED THE PROSE AWARD FOR COMPUTING AND INFORMATION SCIENCE IN 2012.



John F. Elder IV Founder & CEO, Elder Research, Inc.

DR. JOHN ELDER HEADS THE USA'S LARGEST AND MOST EXPERIENCED DATA MINING CONSULTING TEAM. FOUNDED IN 1995, ELDER RESEARCH, INC. HAS OFFICES IN CHARLOTTESVILLE, VIRGINIA, WASHINGTON DC, AND BALTIMORE MARYLAND (WWW.DATAMININGLAB.COM). ERI FOCUSES ON FEDERAL, COMMERCIAL, AND INVESTMENT APPLICATIONS OF ADVANCED ANALYTICS, INCLUDING TEXT MINING, CREDIT SCORING, PROCESS OPTIMIZATION, CROSS-SELLING, DRUG EFFICACY, MARKET TIMING, AND FRAUD DETECTION.

JOHN EARNED ELECTRICAL ENGINEERING DEGREES FROM RICE UNIVERSITY, AND A PHD IN SYSTEMS ENGINEERING FROM THE UNIVERSITY OF VIRGINIA, WHERE HE'S AN ADJUNCT PROFESSOR TEACHING OPTIMIZATION OR DATA MINING. PRIOR TO 19 YEARS AT ERI, HE SPENT 5 YEARS IN AEROSPACE DEFENSE CONSULTING, 4 HEADING RESEARCH AT AN INVESTMENT MANAGEMENT FIRM, AND 2 IN RICE'S *COMPUTATIONAL & APPLIED MATHEMATICS* DEPARTMENT.

DR. ELDER HAS AUTHORED INNOVATIVE DATA MINING TOOLS, IS A FREQUENT KEYNOTE SPEAKER, AND HAS CHAIRED INTERNATIONAL ANALYTICS CONFERENCES. JOHN WAS HONORED TO SERVE FOR 5 YEARS ON A PANEL APPOINTED BY PRESIDENT BUSH TO GUIDE TECHNOLOGY FOR NATIONAL SECURITY. HIS BOOK WITH BOB NISBET AND GARY MINER, *HANDBOOK OF STATISTICAL ANALYSIS & DATA MINING APPLICATIONS*, WON THE PROSE AWARD FOR TOP BOOK IN MATHEMATICS FOR 2009. HIS BOOK WITH GIOVANNI SENI, *ENSEMBLE METHODS IN DATA MINING*, WAS PUBLISHED IN 2010, AND HIS BOOK WITH COLLEAGUE ANDREW FAST AND 4 OTHERS ON *PRACTICAL TEXT MINING* WON THE 2012 PROSE AWARD FOR COMPUTER SCIENCE.