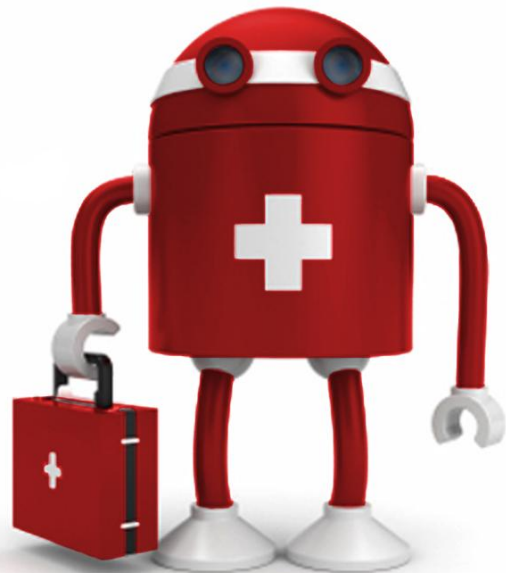


Predicting Diabetes Diagnosis in African Americans Using Ensemble Machine Learning

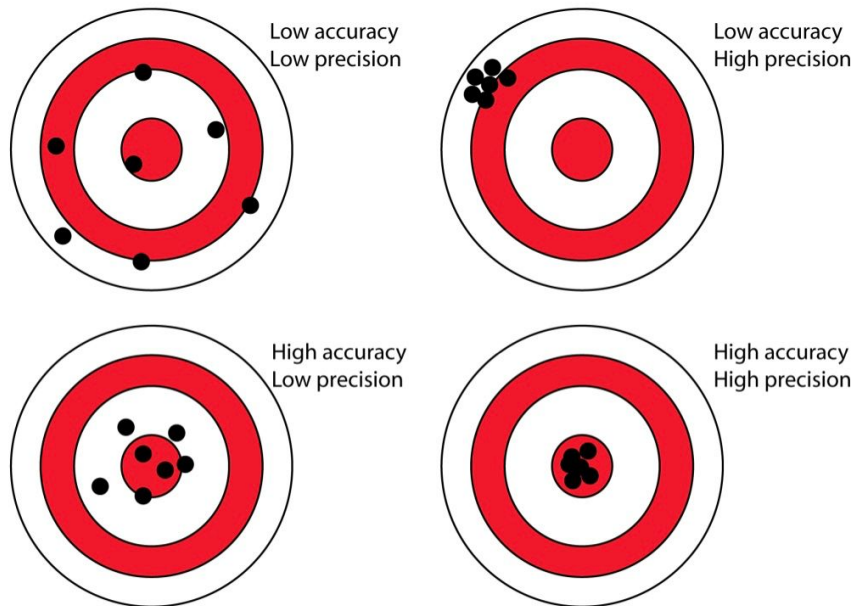
Kimberlyn Bailey, Jarvis Miller and Valerie Santiago-Gonzalez
Under the mentorship of Dr. Sherri Rose and
doctoral candidate Savannah Bergquist



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH

Objective of Our Work

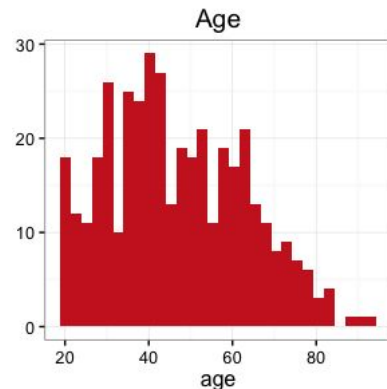
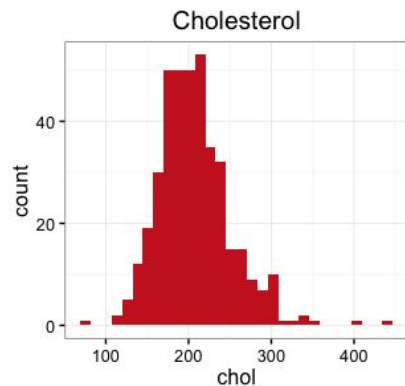
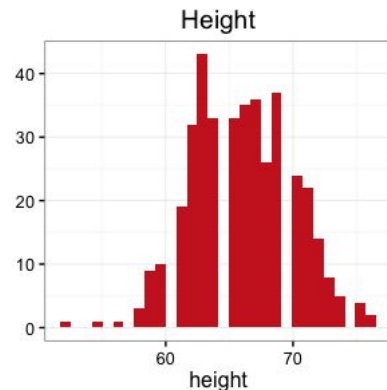
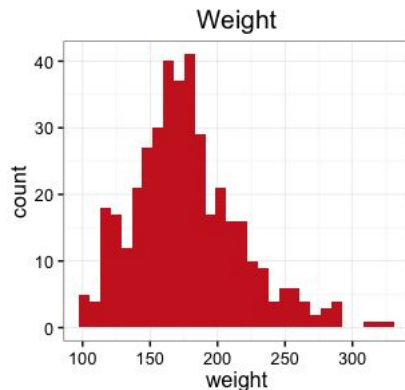
Build an ensemble machine learning algorithm that can accurately predict type II diabetes diagnoses from datasets of relevant patient characteristics.



Diabetes Dataset

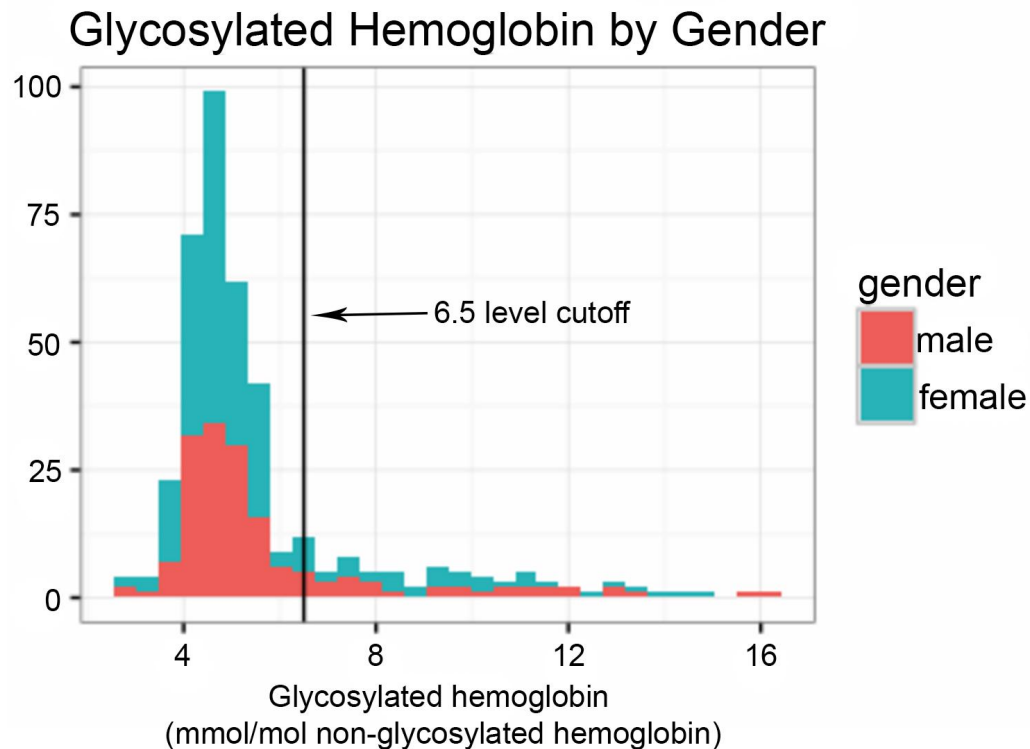
Background:

- Publically available dataset in R compiled by University of Virginia School of Medicine.
- Study designed to examine cardiovascular disease and diabetes trends of central Virginian African Americans.
- 403 subjects
- 14 variables

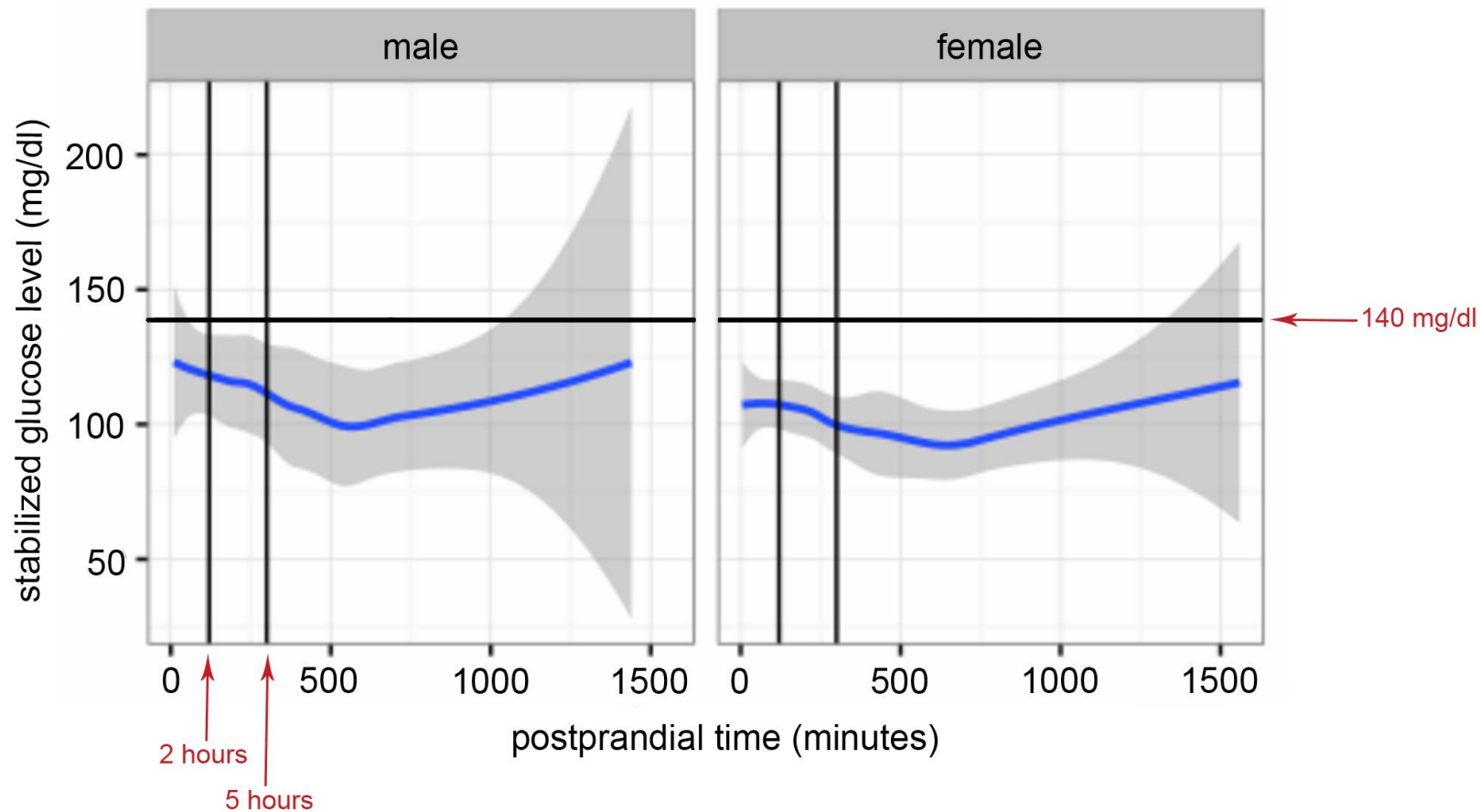


Data Transformation

- Missing data deleted → trimmed to 366 subjects
- Calculated *waist/hip* ratio
- Changed location variable to binary
- Created new variable:
bad cholesterol = total cholesterol - high density lipoprotein
- Changed glycosylated hemoglobin to binary with 6.5 cutoff



Stabilized Glucose Level When Labs Were Drawn

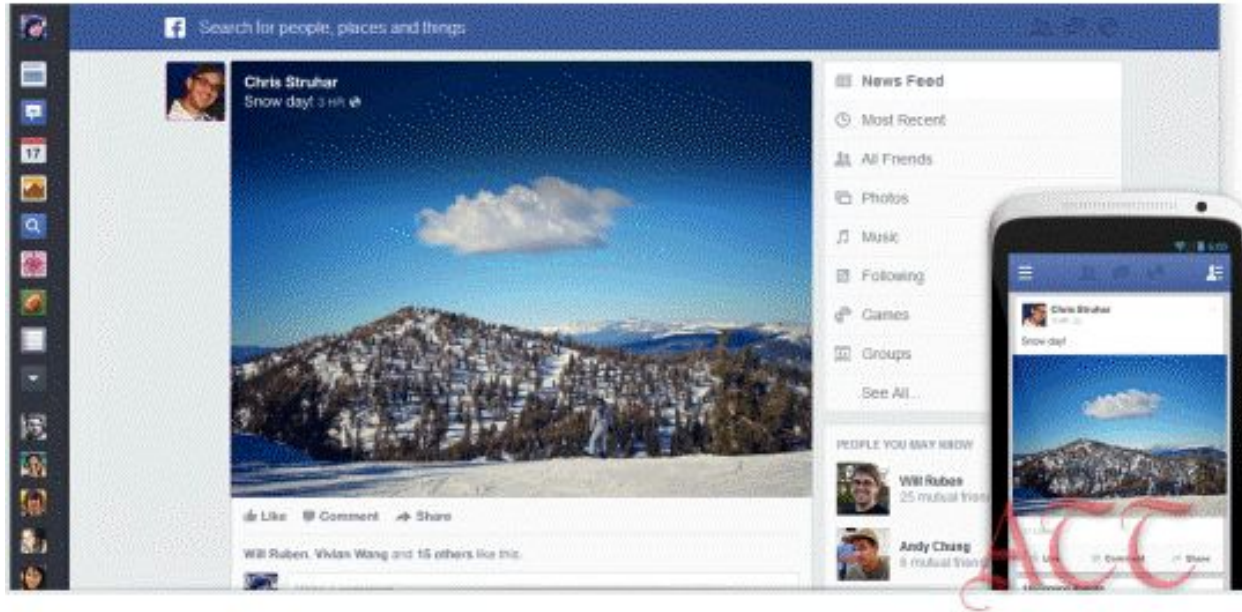




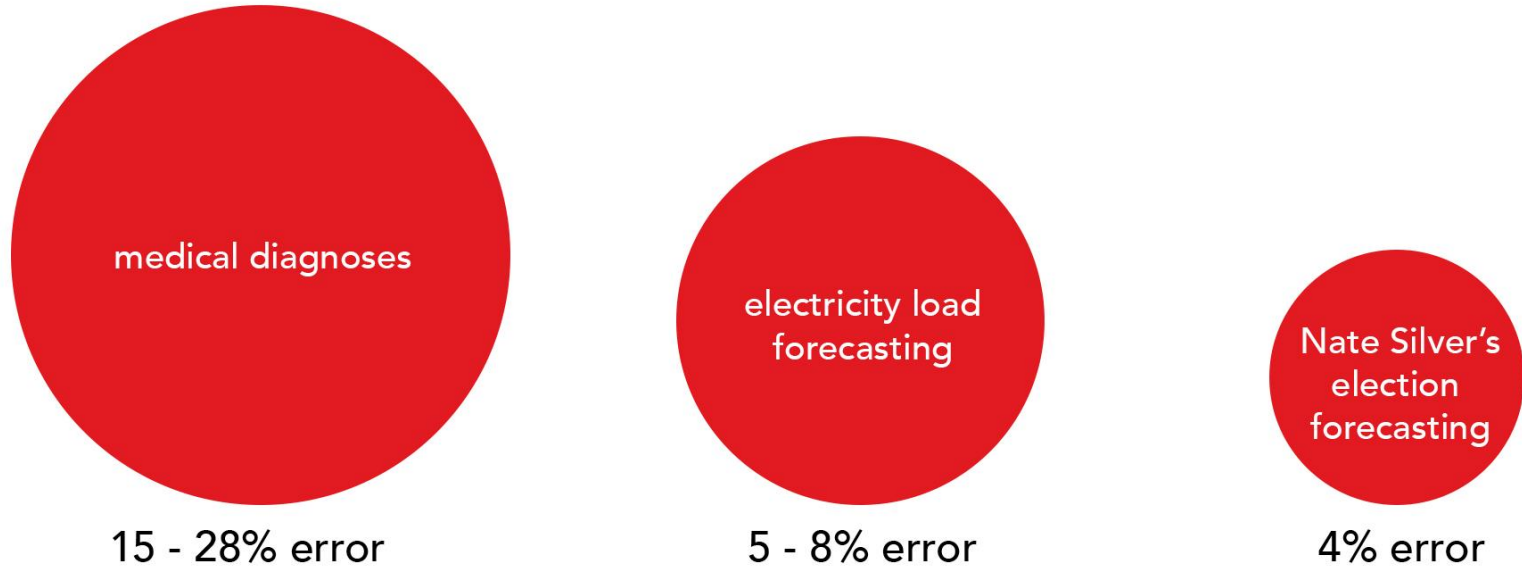
Machine Learning

Algorithms that can independently adapt or “learn” from new data.

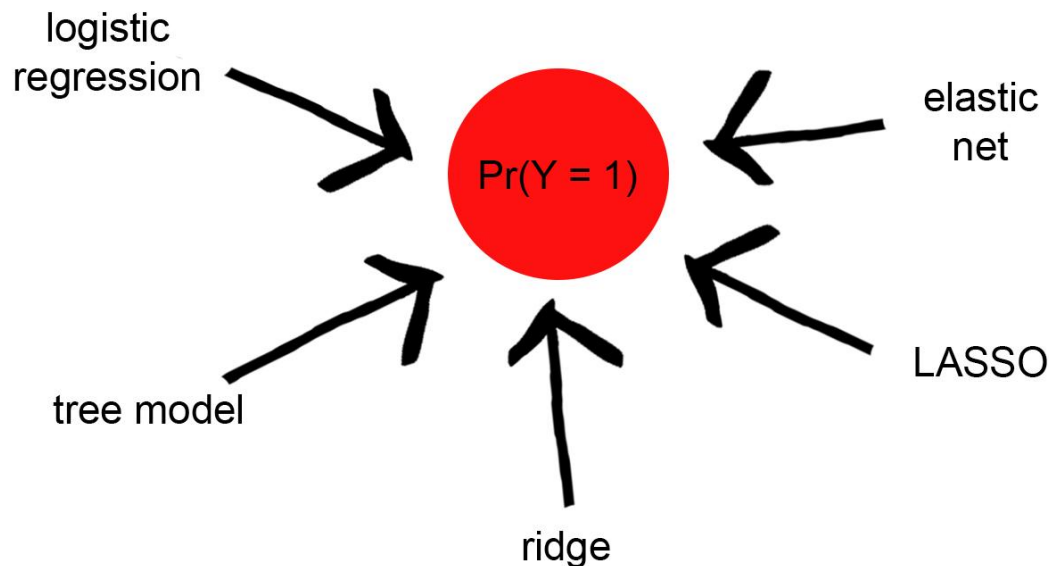
Facebook Newsfeed is Machine Learning



Room for Improvement in Diagnostic Accuracy

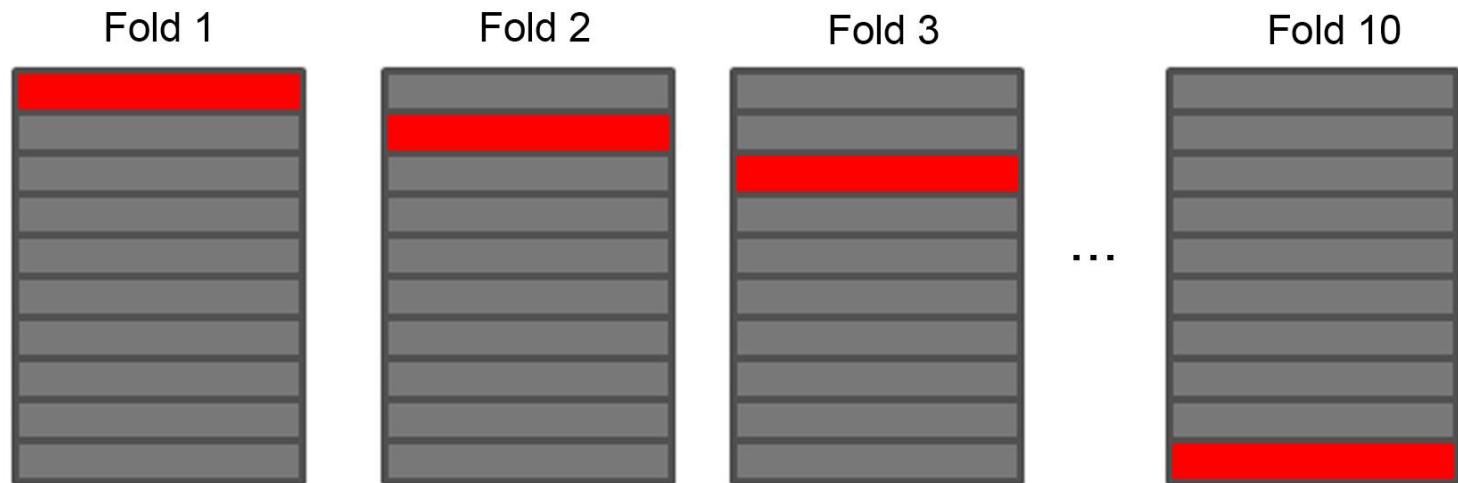
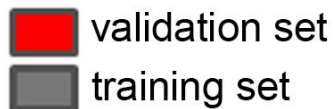


The Basics of Ensemble Machine Learning



Central question of ensemble learning:
How much should each algorithm contribute to the prediction?

Step 1: Determining Performance of the Algorithms

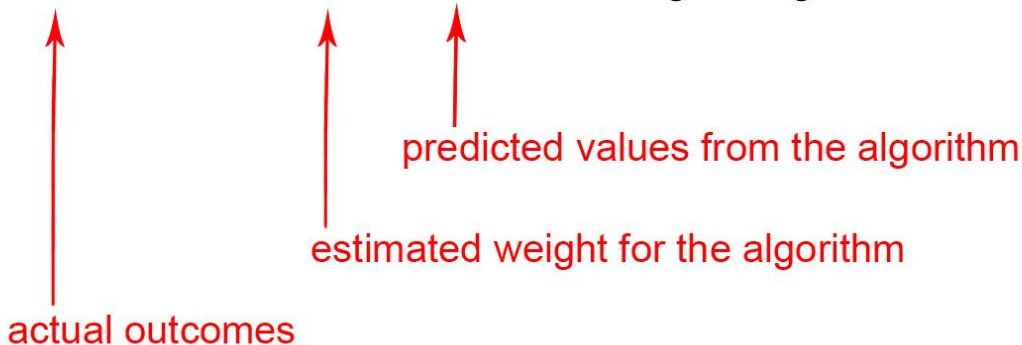


$$\text{CV MSE}_j = \frac{\sum_{i=1}^n (Y_i - D_{j,i})^2}{n}$$

Step 2: Optimizing the Aggregated Performance

Predicted values of each algorithm are used as the inputs in a regression to predict the actual outcomes, optimizing for minimum MSE. Coefficients become the weights on the individual algorithms.

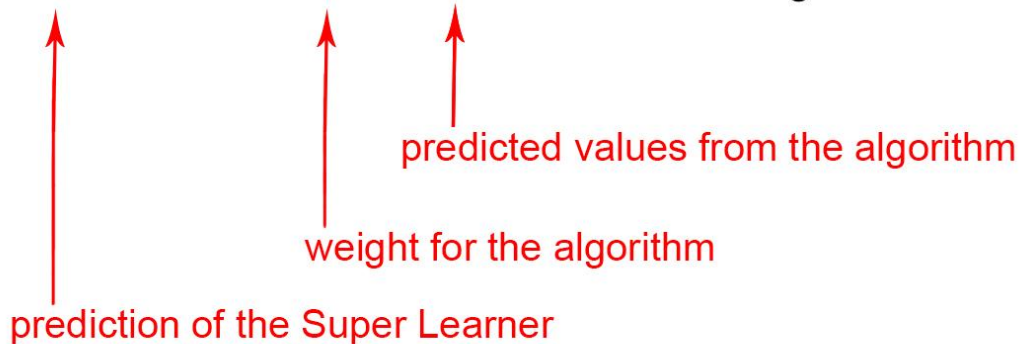
$$\Pr(Y=1 | D) = \alpha_{\text{lasso}} D_{\text{lasso}} + \alpha_{\text{ridge}} D_{\text{ridge}} + \dots + \alpha_{\text{logistic}} D_{\text{logistic}}$$



Step 2: Optimizing the Aggregated Performance

Resultant equation produces predictions of the Super Learner from the weighted aggregate of predictions from each individual algorithm.

$$\Pr(Y=1| D) = 0.46D_{\text{lasso}} + 0.13D_{\text{ridge}} + \dots + 0.27D_{\text{logistic}}$$



Regression Methods

Logistic Regression

- Similar to standard regression but dependent variable is binary
- Coefficients represent effect on log-odds of event occurrence

Ridge Regression

- Reduces variance of model
- Shrinks/Penalizes coefficients

LASSO (least absolute shrinkage and selection operator)

- Can shrink variables to zero → variable selection

Elastic Net

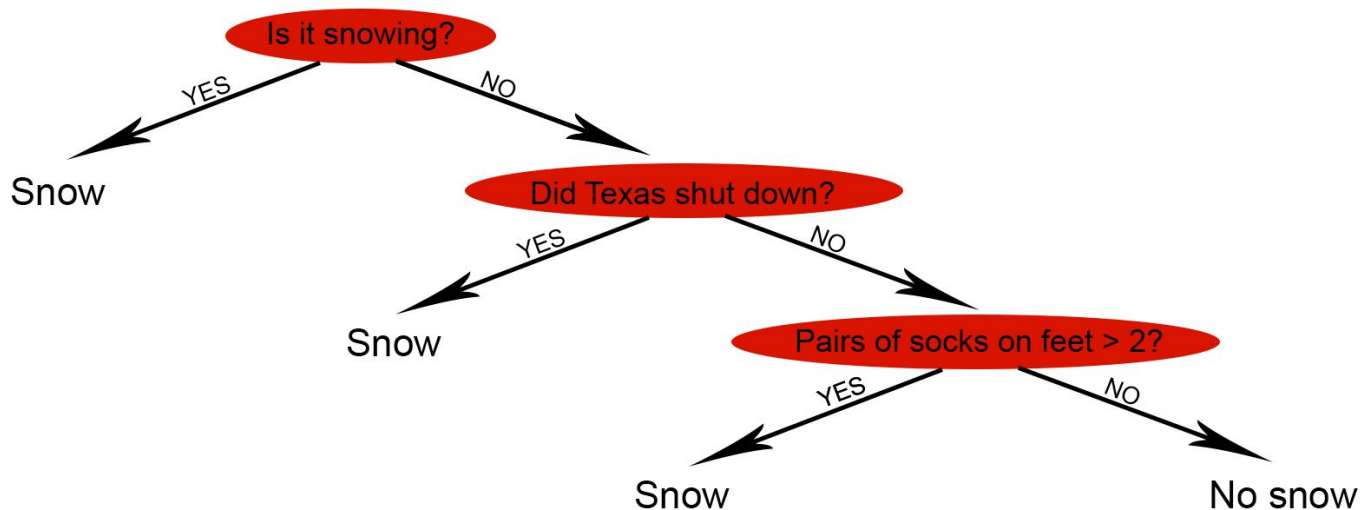
- Compromise between Ridge and LASSO
- Variable selection + correlated predictors

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2,$$
$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t,$$

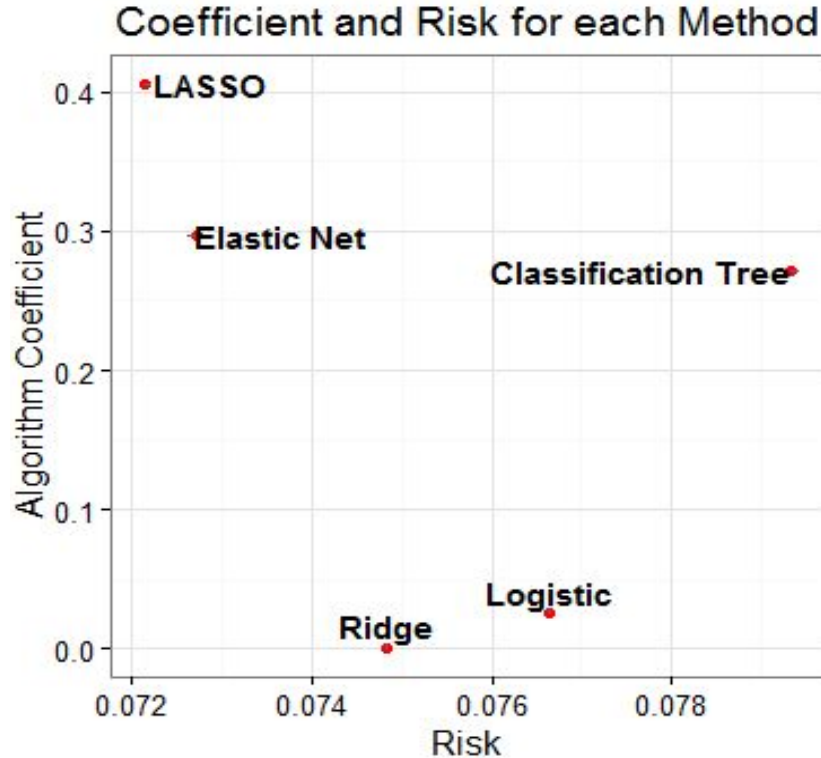
Tree-based Models

Classification Tree

- Each 'node' is a regressor
- Each 'branch' is a Y/N response



Super Learner Results vs. Included Algorithms



Algorithm	MSE	RE
Logistic	0.0766	1.141
Ridge	0.0748	1.114
LASSO	0.0721	1.074
Elastic Net	0.0727	1.083
Tree Model	0.0793	1.181
Super Learner	0.0672	1

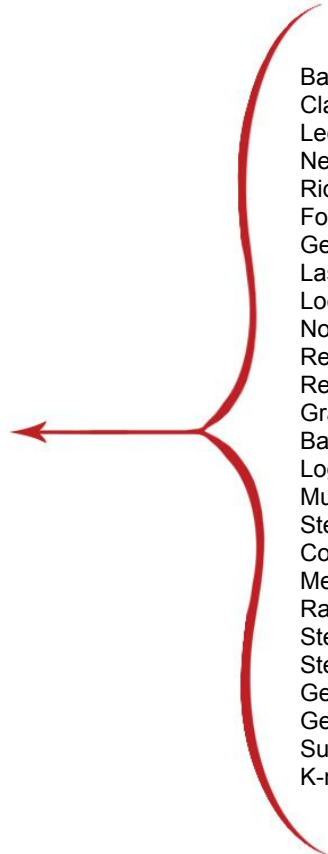
Limitations

- Full EHR dataset

- More predictors
- Actual diagnosis
- Representative

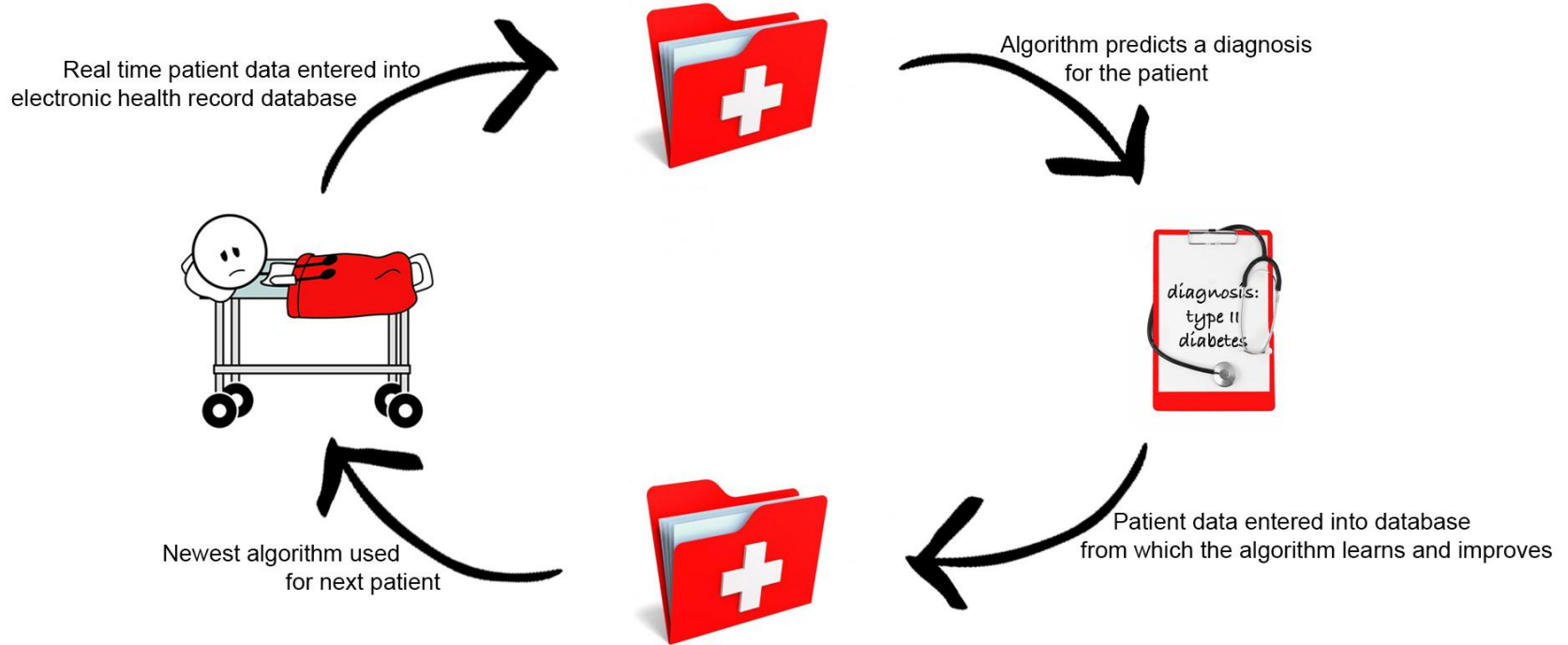
- Include additional algorithms

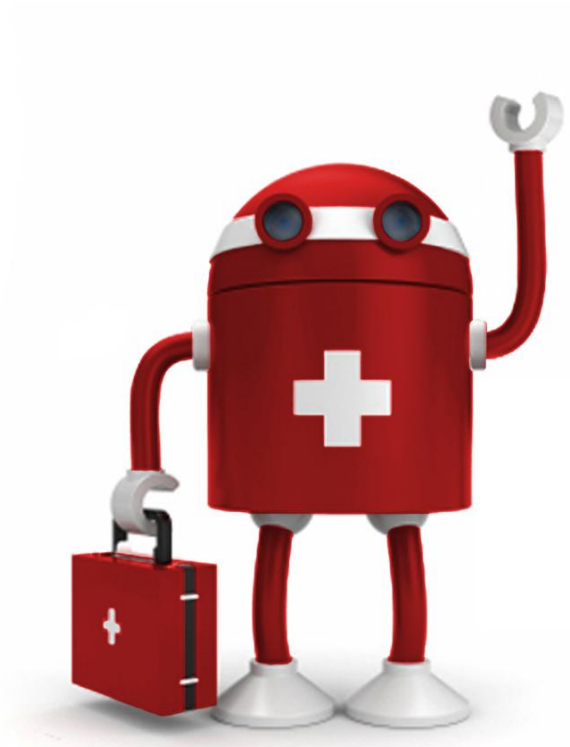
- Separate validation set



Bayesian generalized linear model
Classification and Regression Training
Leekasso
Neural nets
Ridge regression
Forward stepwise regression
Generalized additive models
Lasso and elastic net generalized linear models
Local regression
Non-Negative Least Squares (NNLS)
Recursive partitioning and regression trees
Regressive partitioning and regression trees with pruning
Gradient boosting method
Bagging classification trees
Logistic regression
Multivariate adaptive polyspline regression
Stepwise regression by Akaike information criterion (AIC)
Control forest
Mean algorithm
Random forest
Stepwise regression
Stepwise regression with interactions
Generalized linear model
Generalized linear model with interactions
Support vector machine
K-nearest neighbors algorithm

Ultimate Goal: Machine Learning Diagnostics in Real Time





Thank you for listening!
Any questions?