Group 14

# PREDICTING ALZHEIMER'S DISEASE

TEXAS
The University of Texas at Austin
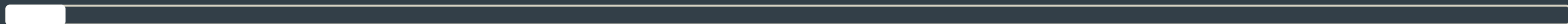
**Aileen Li, Grace Lin, Kristen Lowe, Liyan Deng, Sidharth Saha**
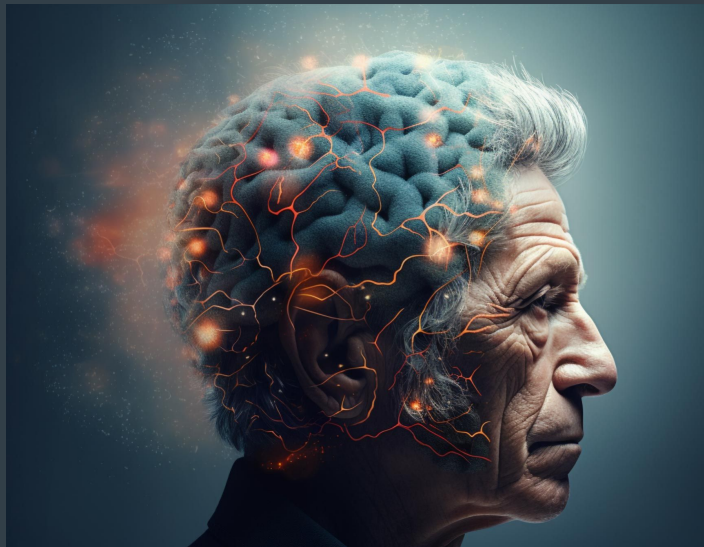MSBA, The University of Texas at Austin

# Agenda

- Introduction and Exploratory Data Analysis
- KNN
- Logistic regression
- Trees (Bagging and boosting)
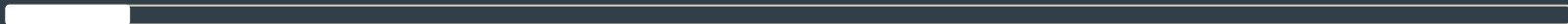- Conclusion and reflection

# Alzheimer's Disease Dataset

- From Kaggle
- 2,149 patient observations and 35 columns
- Target variable: Alzheimer's diagnosis (Yes/No)
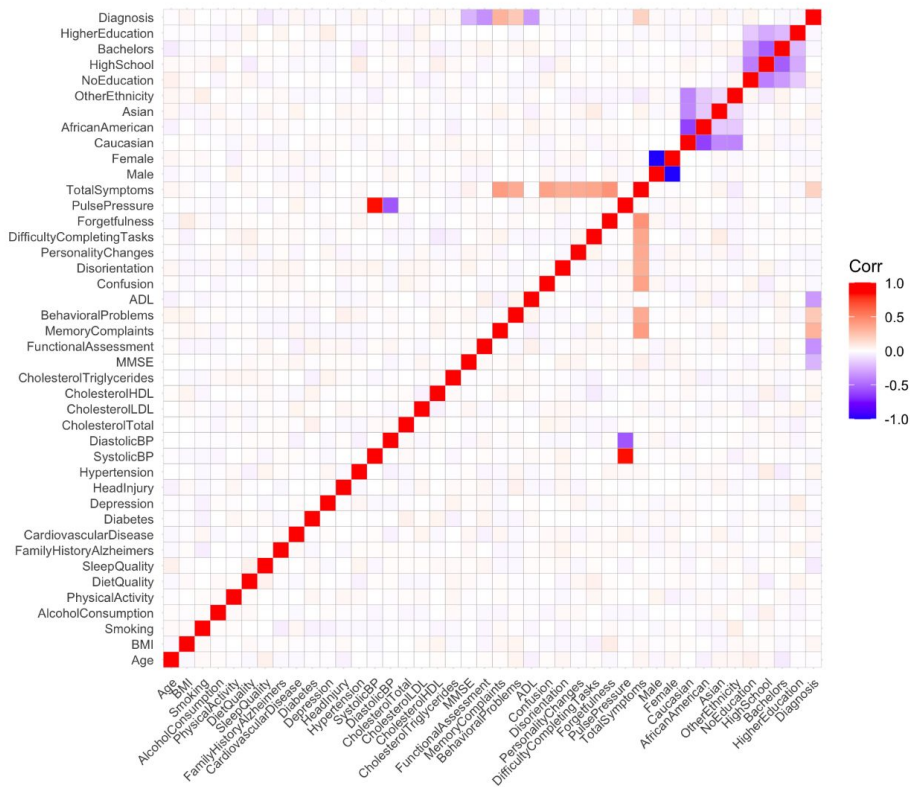- Data mildly imbalanced
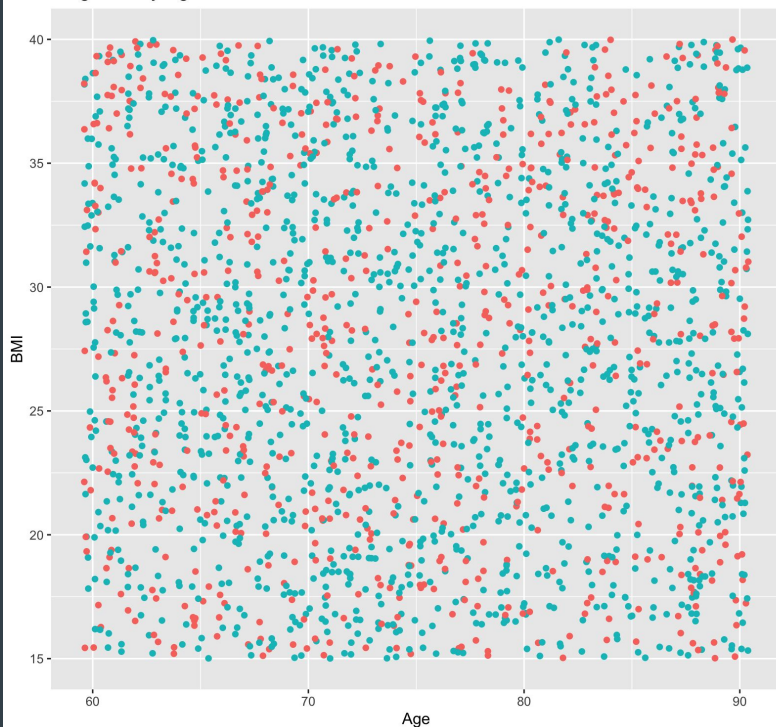
# Feature Engineering

- **Made two new predictor variables**
  - Pulse Pressure
  - Total Symptoms
- **One hot encoding**
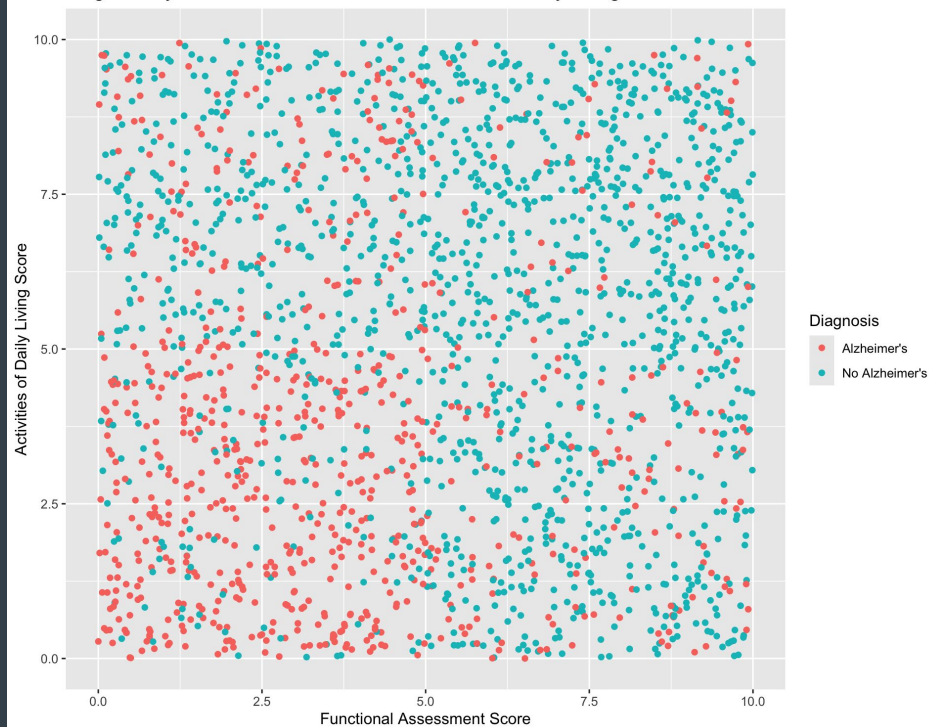  - Gender
  - Ethnicity
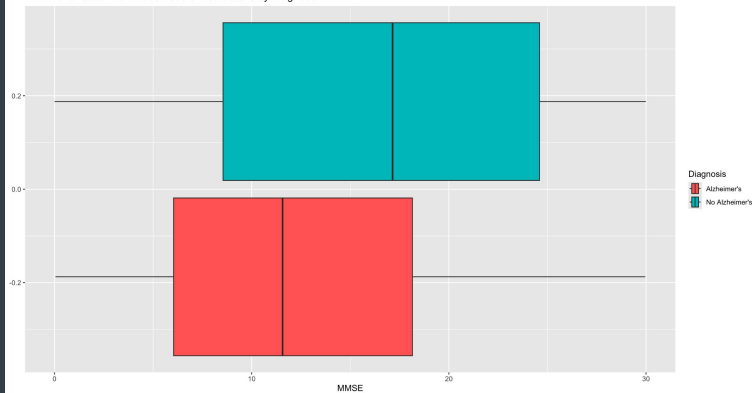  - Education Level

Correlation Heatmap
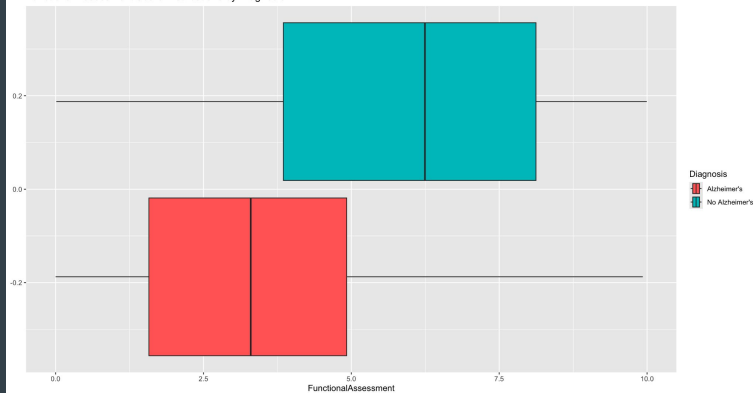
## Diagnosis by Age and BMI



## Diagnosis by Functional Assessment and Activities of Daily Living Scores
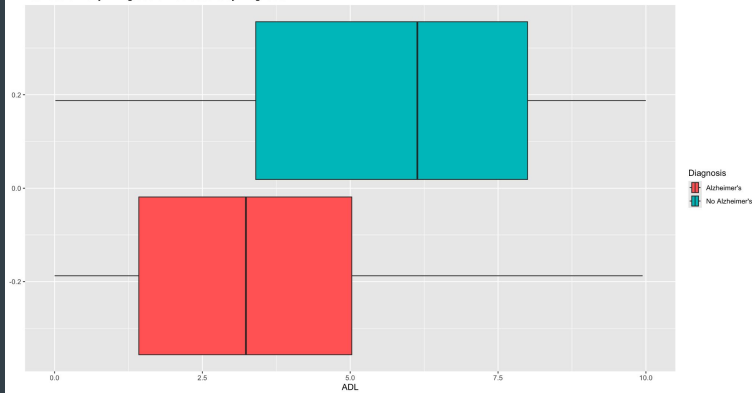


Diagnosis

- Alzheimer's
- No Alzheimer's

Mini-Mental State Examination Score Distributions by Diagnosis
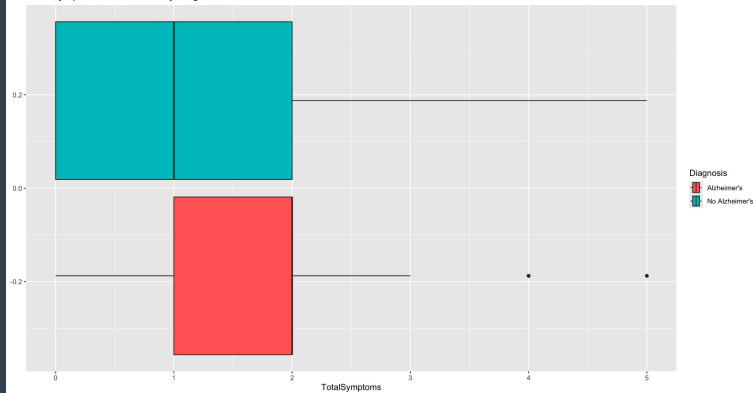


Functional Assessment Score Distributions by Diagnosis



Activities of Daily Living Score Distributions by Diagnosis



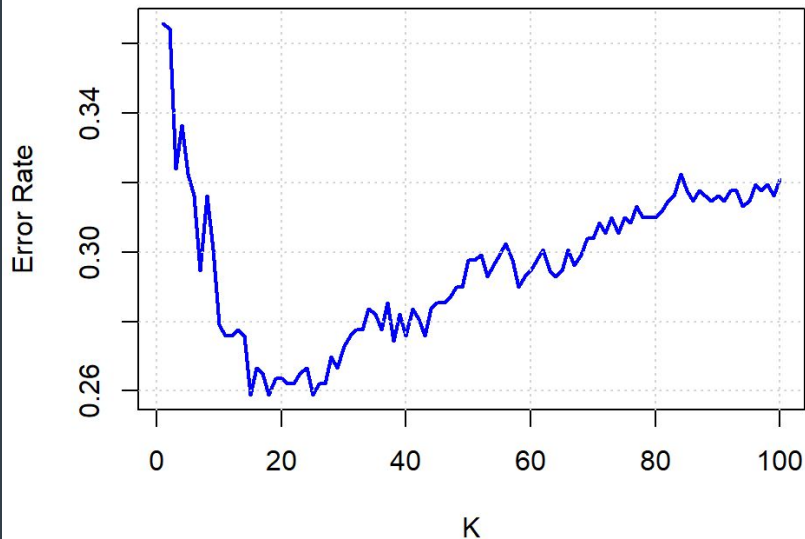Total Symptoms Distributions by Diagnosis

# KNN

- Utilized all features
- Seed = 1
- k = 24

# KNN Performance

- Confusion Matrix

Accuracy: 73%

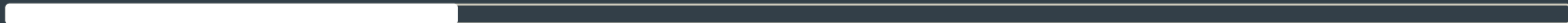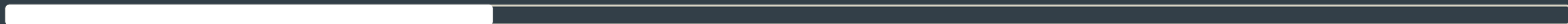|  |  | Test Data | |
| --- | --- | --- | --- |
|  |  | YES | No |
| Prediction | YES | 262 | 112 |
|  | No | 6 | 50 |

# Logistic Regression

- **Methods**
  - All predictors
  - All statistically significant predictors
  - Lasso
  - Ridge
  - Stepwise

# Logistic Regression

- **Use the following predictors:**
  - MMSE
  - Functional Assessment
  - Memory Complaints
  - Behavioral Problems
  - ADL
- Threshold = 0.5
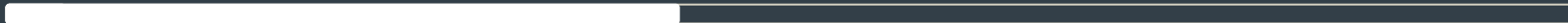
```
## Deviance Residuals:
##       Min        1Q    Median        3Q       Max
## -0.95553  -0.27230  -0.04281   0.26716   1.14373
##
## Coefficients:
##                          Estimate Std. Error t value Pr(>|t|)
## (Intercept)             0.9694298  0.0267783   36.20   <2e-16 ***
## MMSE                   -0.0142408  0.0009952  -14.31   <2e-16 ***
## FunctionalAssessment   -0.0536886  0.0029883  -17.97   <2e-16 ***
## MemoryComplaints        0.3343039  0.0214397   15.59   <2e-16 ***
## BehavioralProblems      0.3368824  0.0235439   14.31   <2e-16 ***
## ADL                    -0.0523070  0.0029245  -17.89   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```
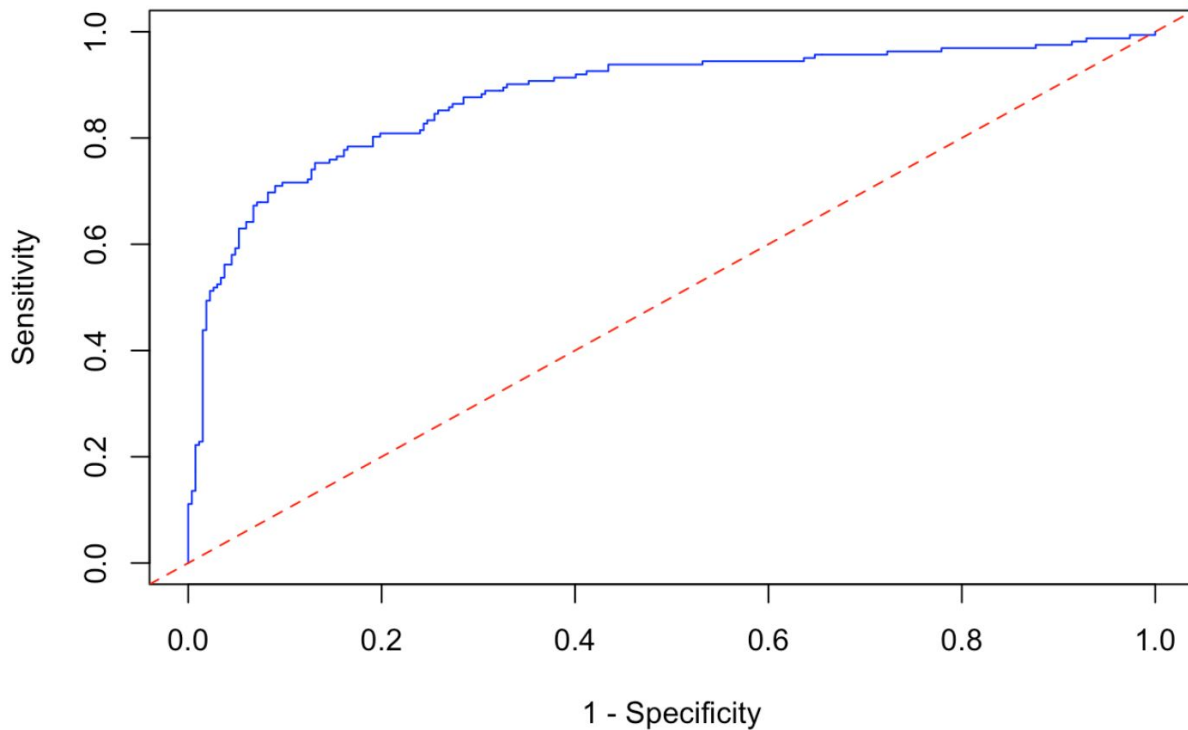
# Logistic Performance

Accuracy = ~83%

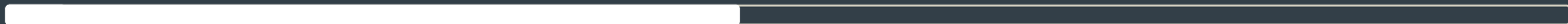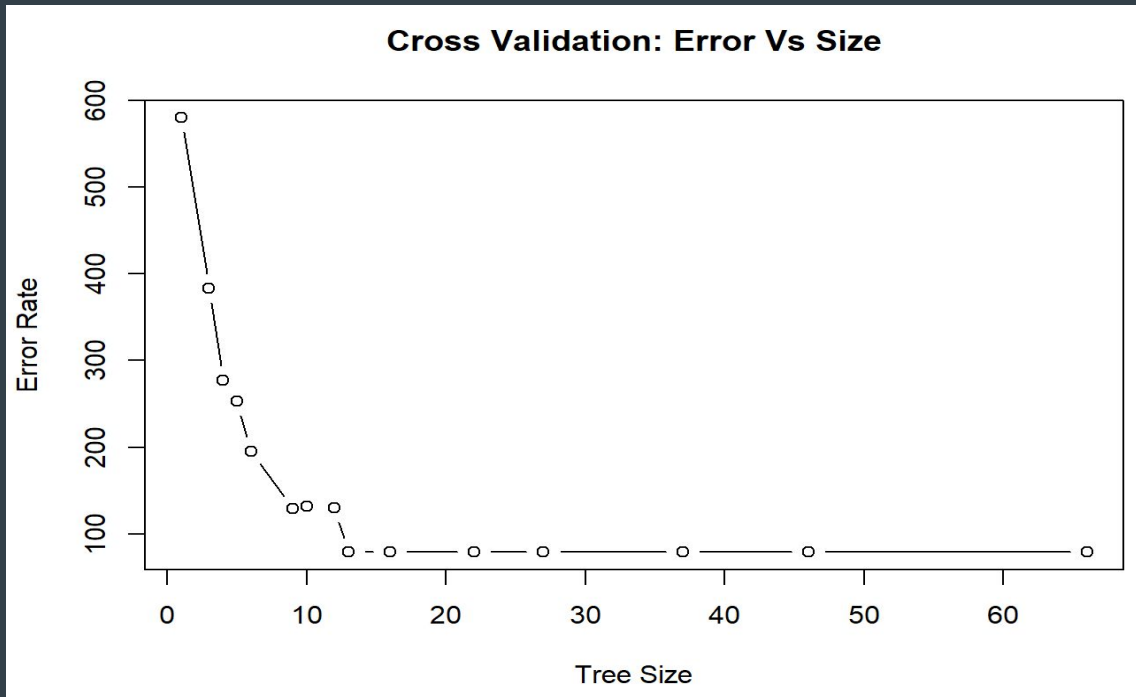|  |  | Test Data | |
|---|---|---|---|
|  |  | YES | No |
| Prediction | YES | 110 | 21 |
|  | No | 52 | 246 |

ROC Curve for Alzheimer's Diagnosis

AUC = 87.89%

# Trees

- **A Classification Tree**

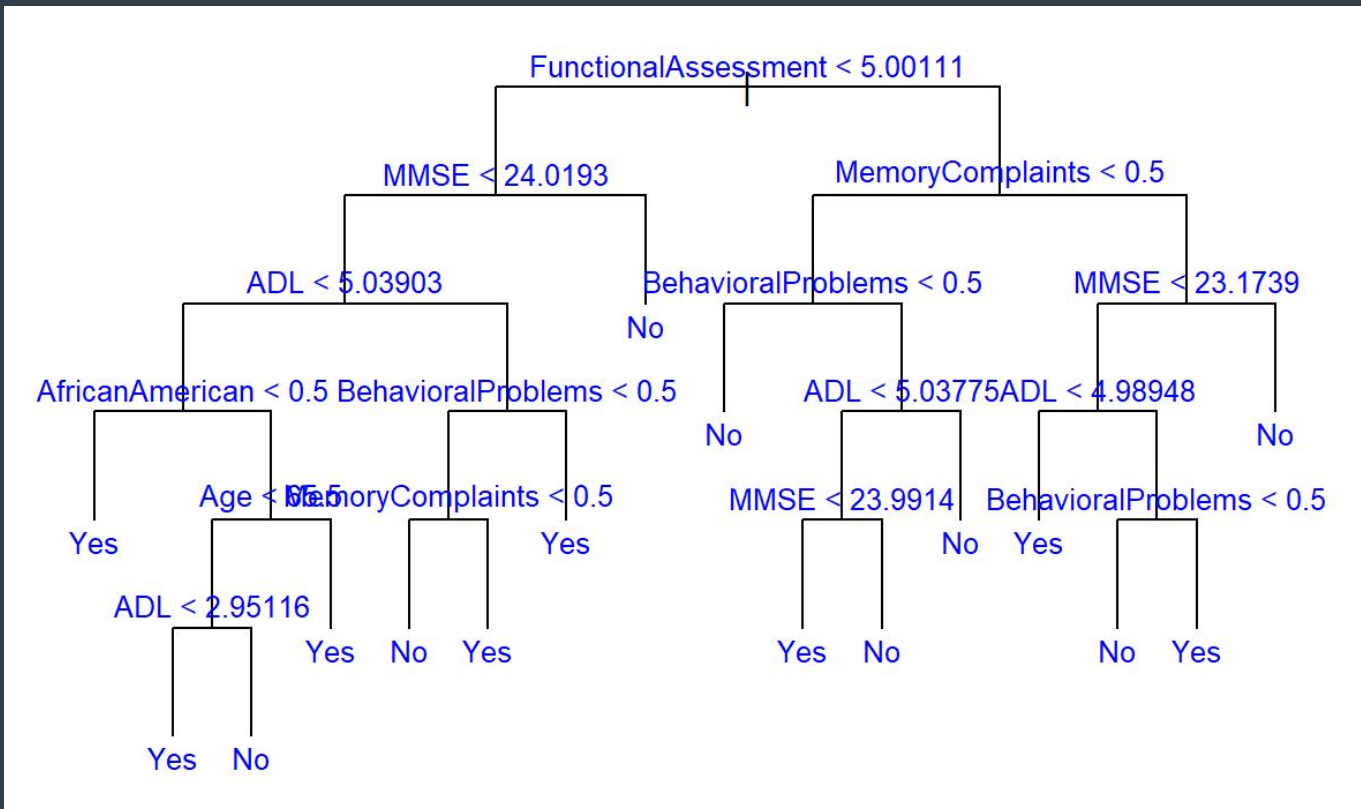- **Random Forest**

- **Bagging**

- **Boosting**
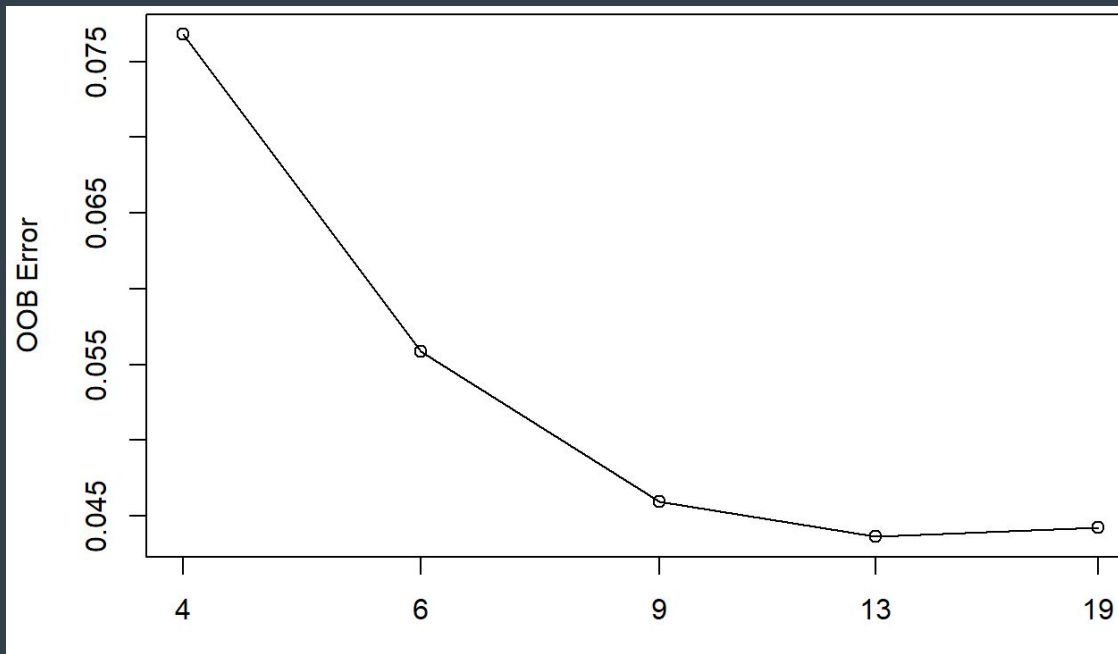
# Classification Tree

- Grow & Prune method



**Cross Validation: Error Vs Size**

# Tree Performance

Accuracy = ~94%

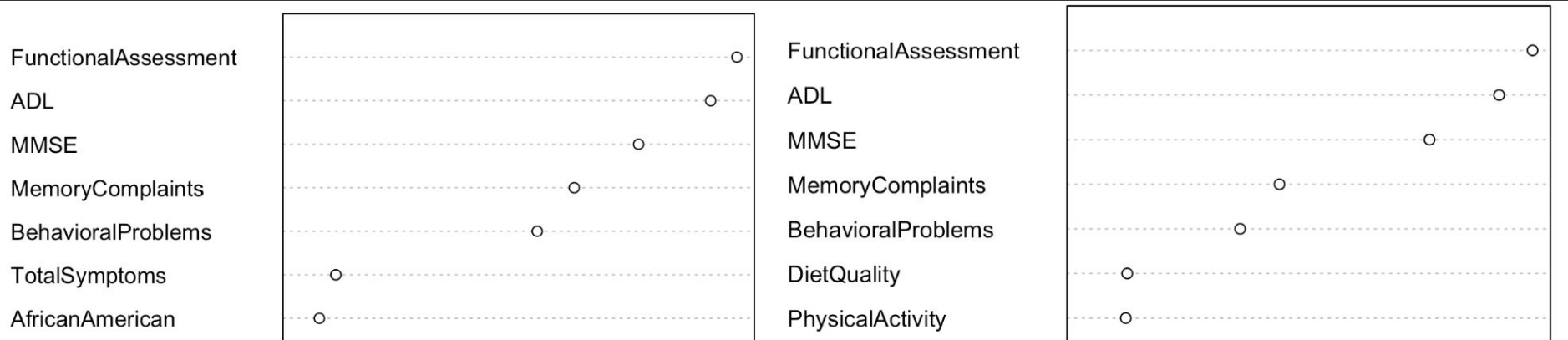| Pruned Tree | | Test Data | |
|---|---|---|---|
| | | YES | No |
| Prediction | YES | 145 | 9 |
| | No | 17 | 259 |

# Random Forest

- RFM1 : mtry = 6

- mtry tuning

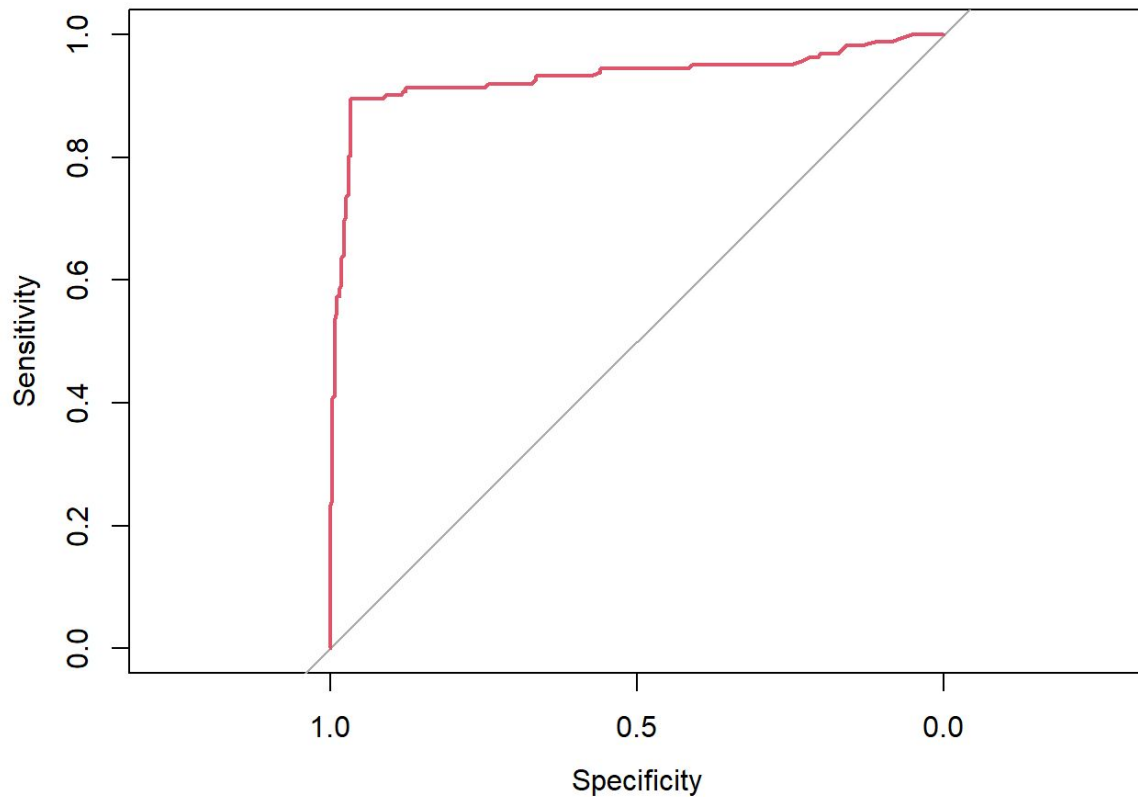- ntrees = 500

# Feature Importance



- Mean Decrease Accuracy
- Mean Decrease Gini

# RF Performance Accuracy = 95.6%*

| RF | | Test Data | |
|---|---|---|---|
| | | YES | No |
| Prediction | YES | 549 | 49 |
| | No | 27 | 1094 |

ROC Curve for Random Forest Model2

# Bagging

- Cross-validation
- folds = 5

# Bagging Feature Importance
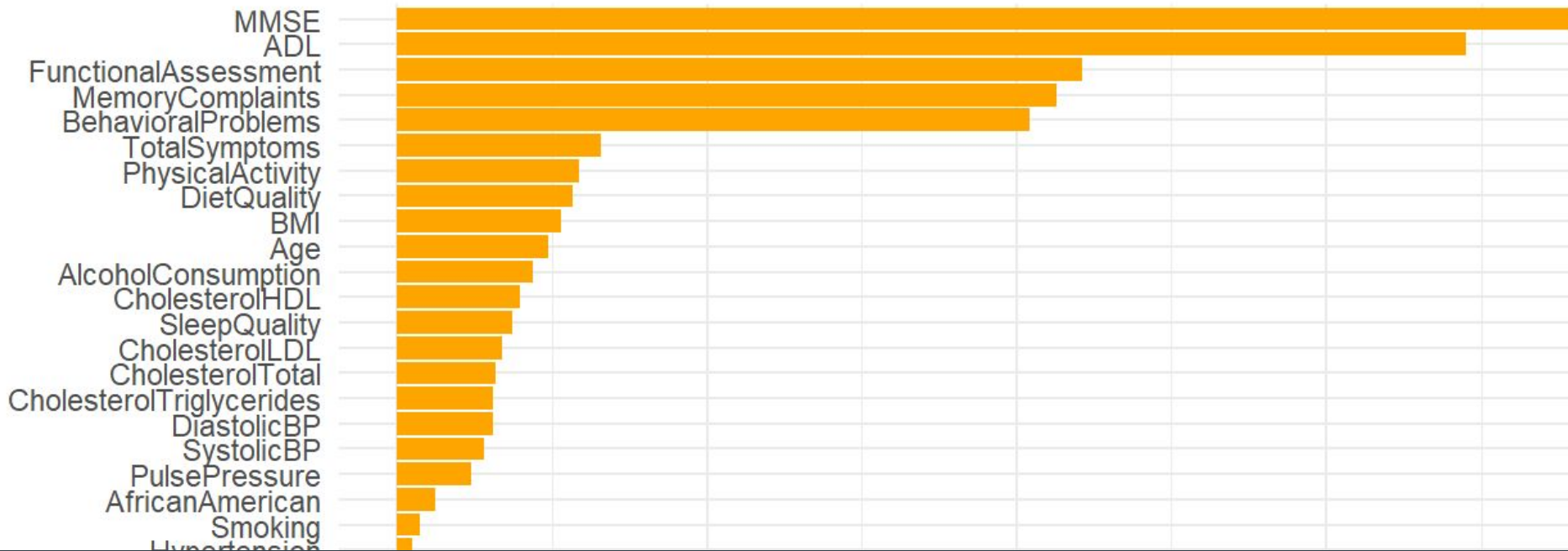


- Mean Decrease Accuracy
- Mean Decrease Gini

# Bagging



Variable Importance Plot

# Bagging Performance

Accuracy = 92.6%

| Bagging with CV | | Test Data | |
| --- | --- | --- | --- |
| | | YES | No |
| Prediction | YES | 141 | 11 |
| | No | 21 | 257 |

ROC Curve for Bagging Model

# Boosting

- cv =10
- n.trees = 500

# Boosting

Accuracy = 91.2%

| Boosting with CV | | Test Data | |
|---|---|---|---|
| | | YES | No |
| Prediction | YES | 133 | 29 |
| | No | 9 | 259 |

ROC Curve for Gradient Boosting Model

# Conclusion

- Best model: **Random Forest**
- High accuracy (95.6%)
- Robustness
  - Sensitivity: 95.31%
  - Specificity: 95.71%
  - Good for imbalanced data

# Reflection

- **Identifying causes/outcome relationships**
  - Feature engineering
  - Correlation analysis
- **Evaluation metrics**
  - Opportunity cost for 'False positive' vs 'False negative'

| Descriptive Analysis | Diagnostics Analysis |
|---|---|
| Predictive Analysis | Prescriptive Analysis |

# Questions?