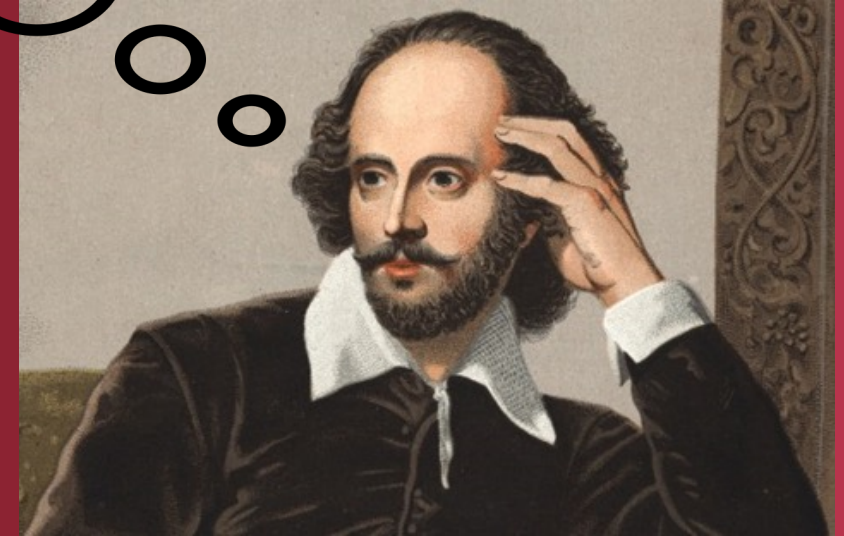


"To Boost,  
or not to Boost?"



# An Analysis of Gradient-Boosting Methods

Romith C., Troy J., Adi B.



UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL OF  
ENGINEERING & COMPUTER SCIENCE

# Research Question:

*Does an increase in sophistication of tree-based boosting algorithms impact the accuracy & efficiency of NFL play predictions?*

## Methodology:

- Split multiple, consecutive years of processed **play-by-play** data into *training* & *test* sets
- Run a basic **decision-tree** model for baseline classification performance
- Implement a standard **gradient-boosting** algorithm & an **extreme-gradient boosting** algorithm (XG-Boost)
- Evaluate performance based on:
  - **Accuracy** rate of predicting a multiclass categorical output (**RUN** vs. **PASS** vs. **FIELD-GOAL** vs. **PUNT**)
  - Relative **time** to run models



UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL OF  
ENGINEERING & COMPUTER SCIENCE

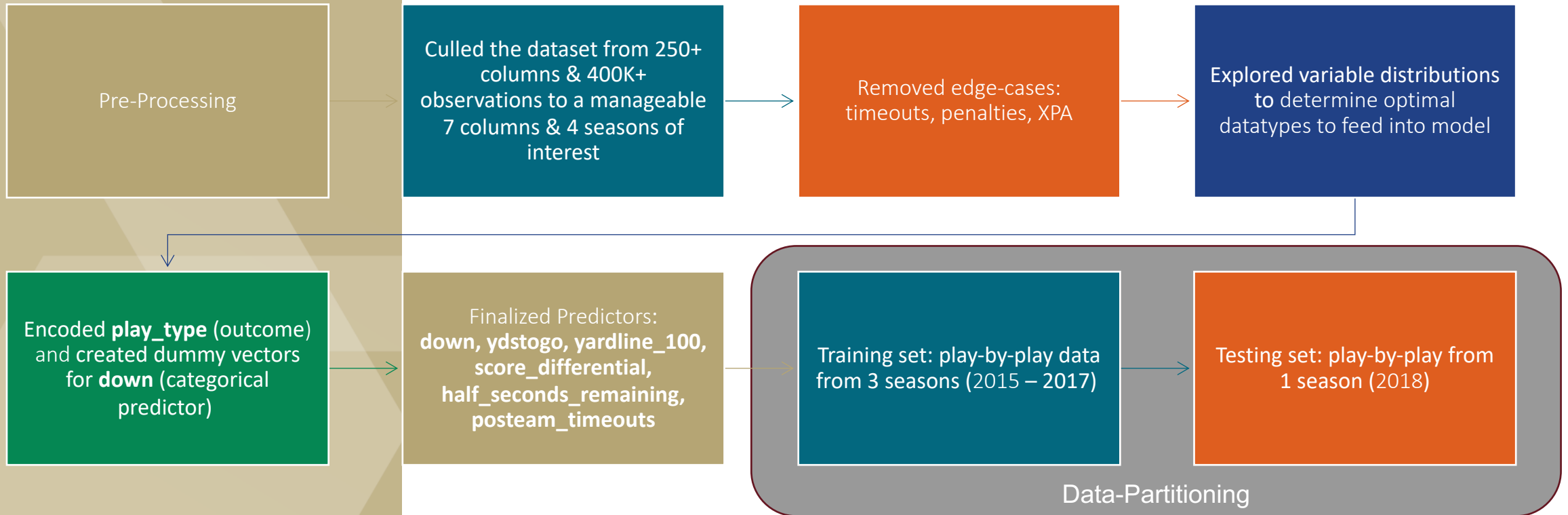
## Data Source and Definitions

Variable	Description	Data Type	Type
yardline_100	Yards until the goal line	Integer from [0, 100]	Independent
half_seconds_remaining	Seconds remaining in the half	Integer from [0, 1800]	Independent
down	The current play down	Integer from [1, 4]	Independent
ydstogo	Yards until first down	Integer from [0, 100]	Independent
posteam_timeouts	Offensive team timeouts remaining	Integer from [0, 3]	Independent
score_differential	The difference in score between teams	Integer from $(-\infty, \infty)$	Independent
play_type	The resulting type of play	Categorical: run, pass, field_goal, punt	Dependent

**Source:** <https://www.kaggle.com/maxhorowitz/nflplaybyplay2009to2016>

Data was scraped and uploaded onto **Kaggle** by the Carnegie Mellon Sports Analytics Club

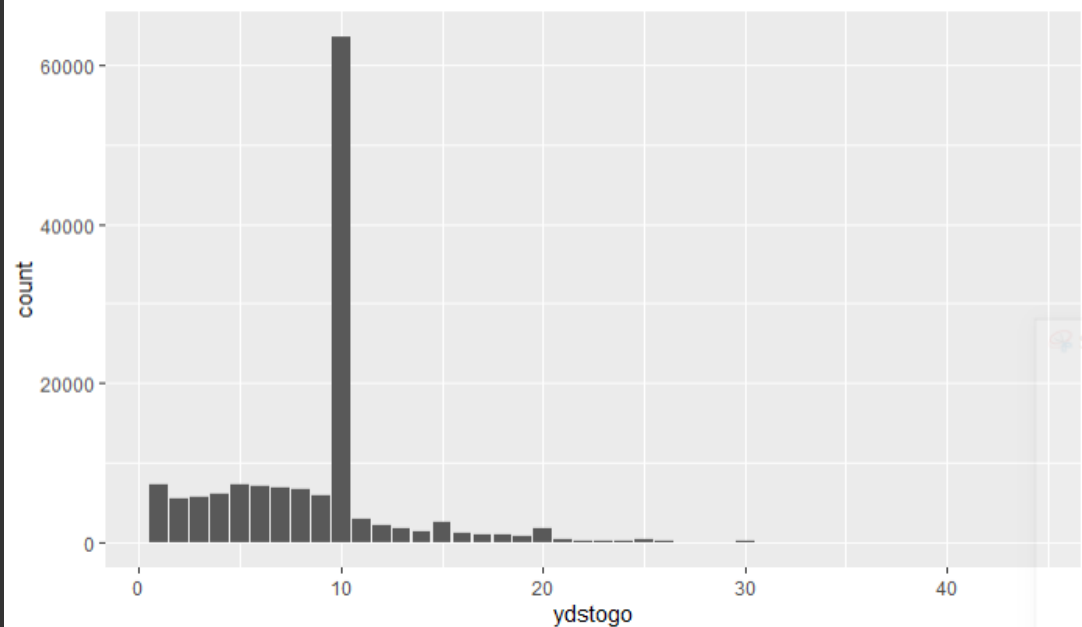
# EDA and Data Preparation



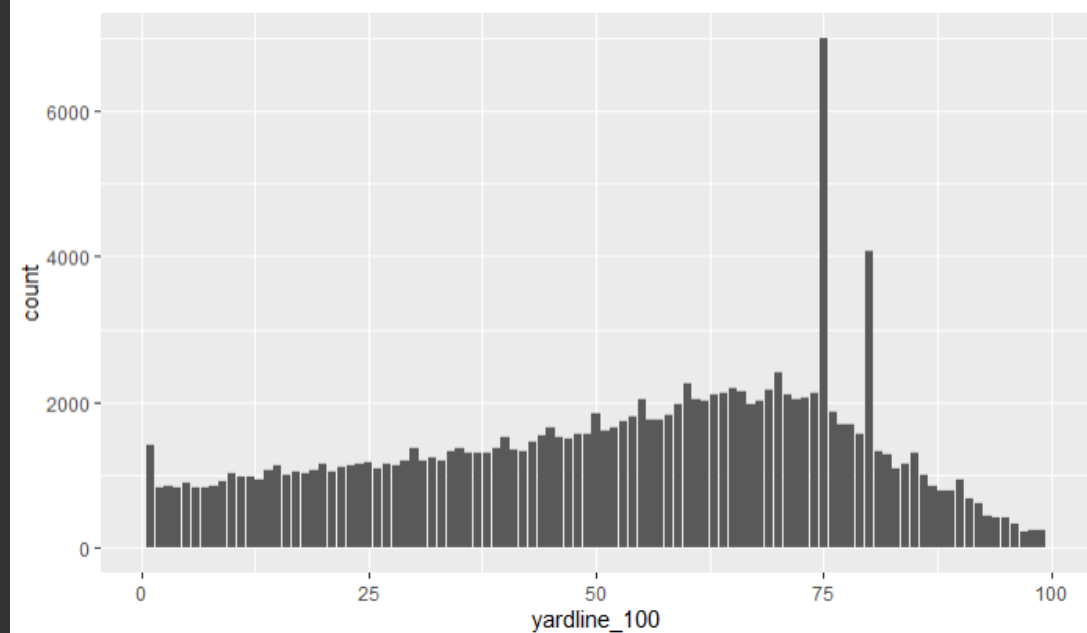
UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE

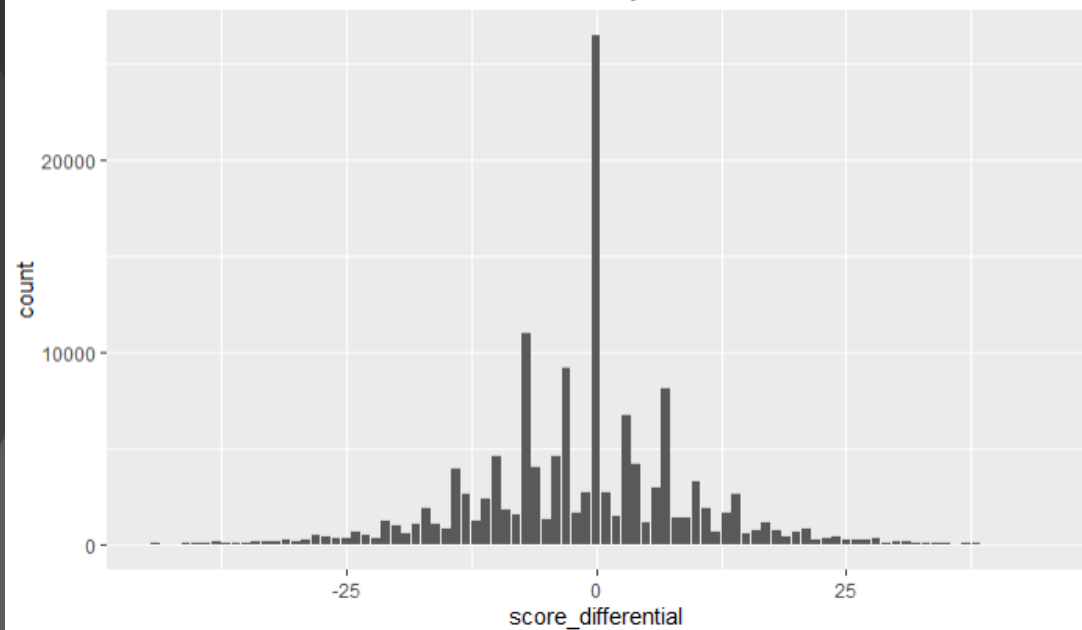
Distribution of Yards-Till-First-Down



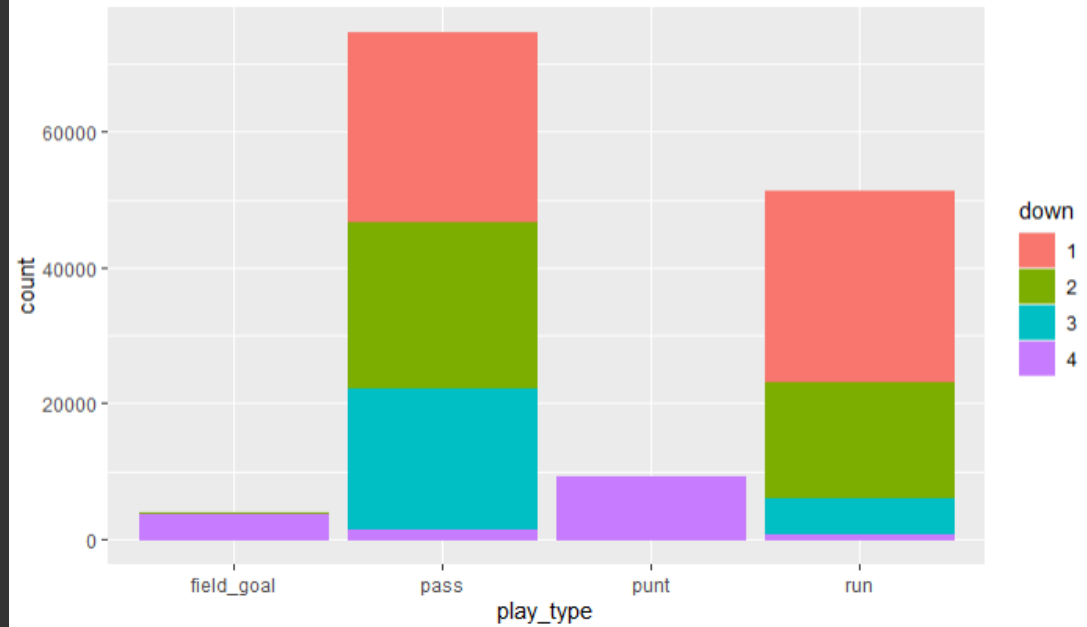
Distribution of Yards-Till-Endzone

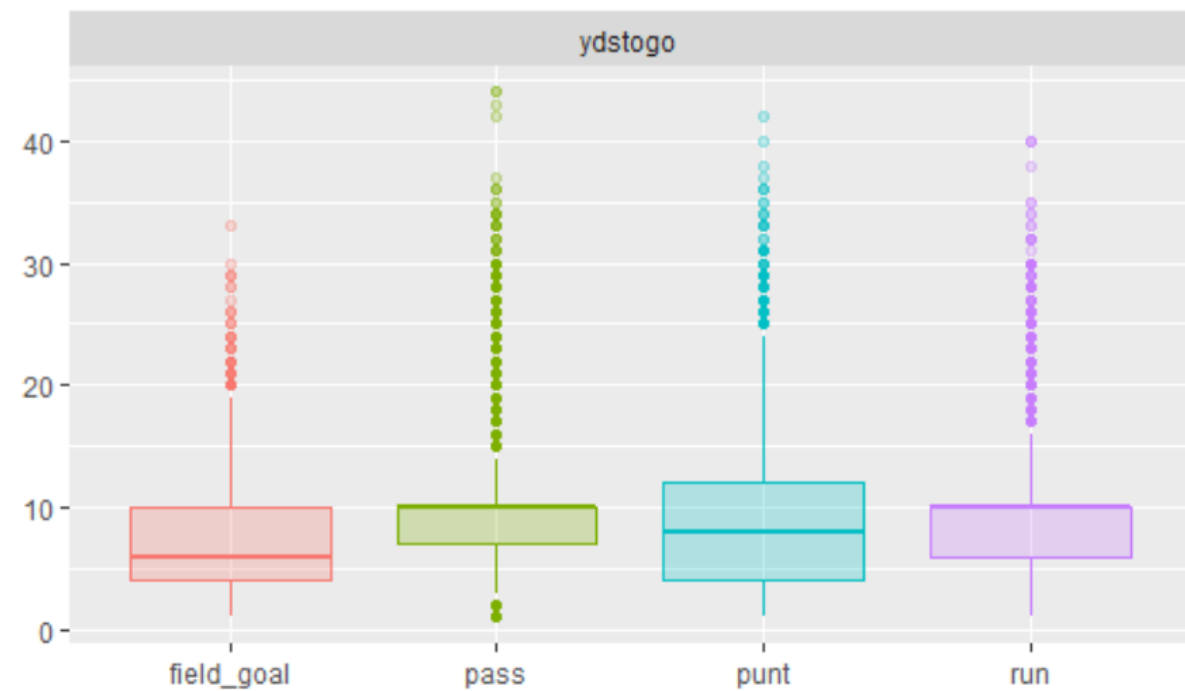
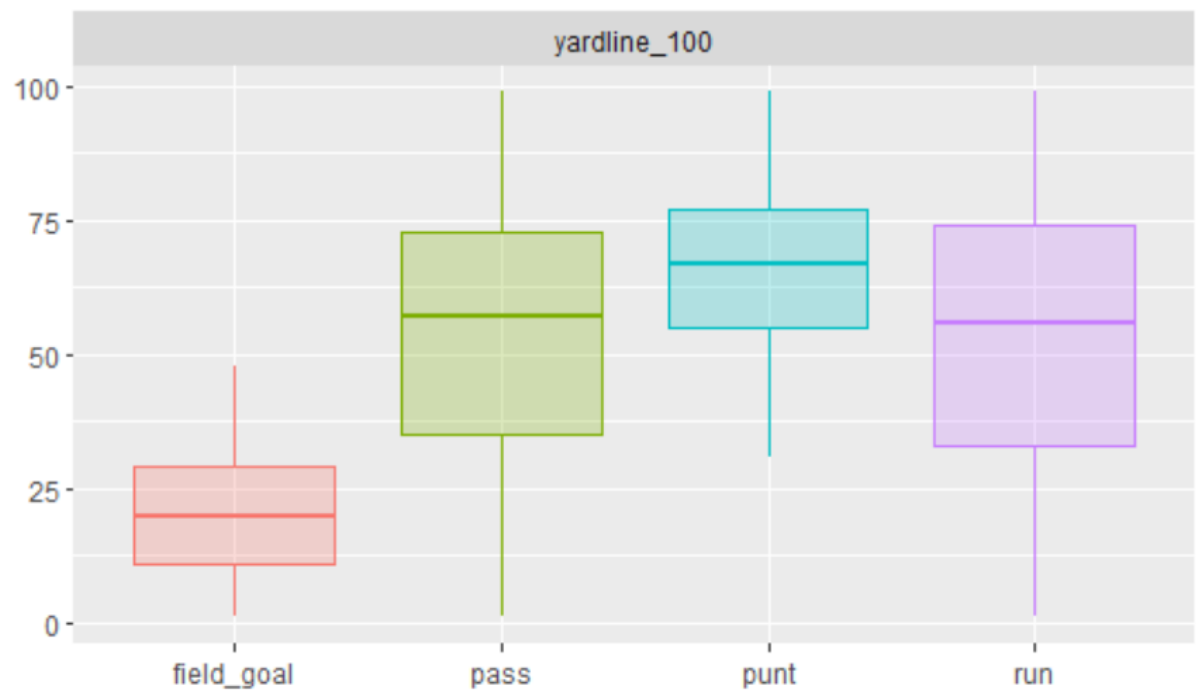
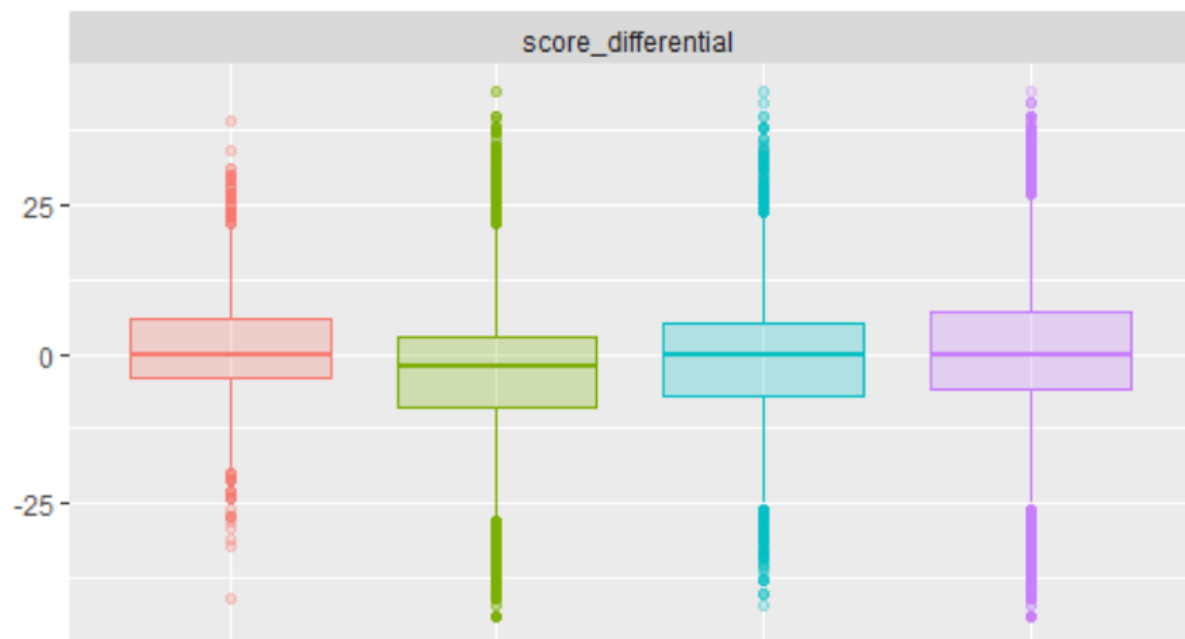
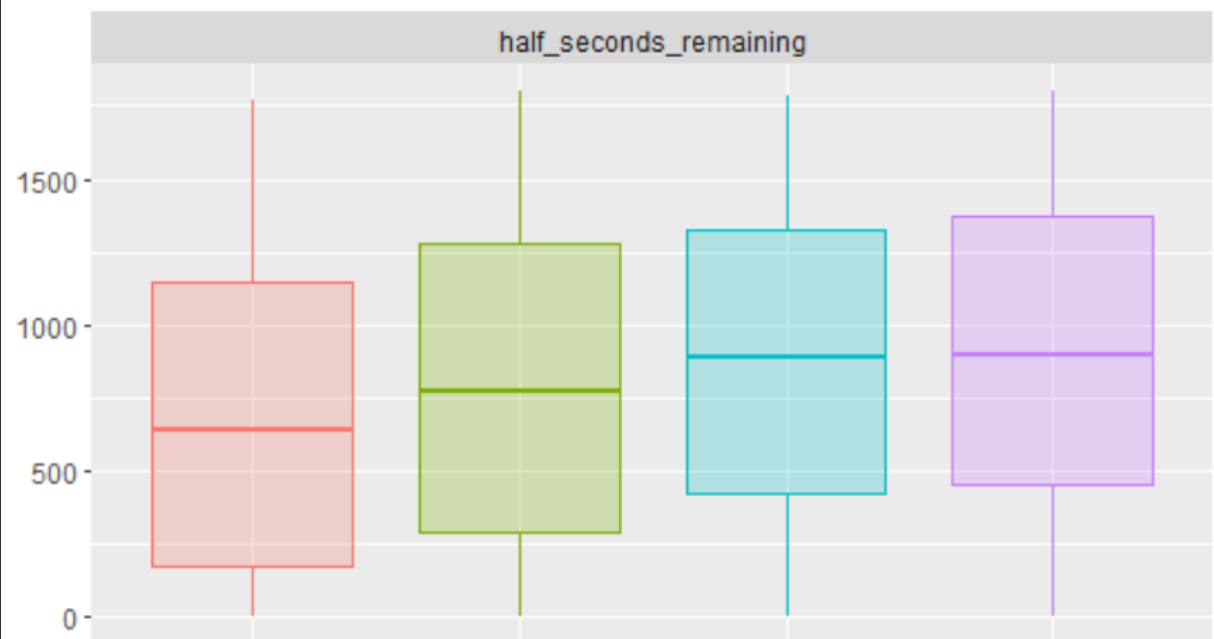


Distribution of Score-Differential for all Plays



Frequency of Play Types based on Down





play\_type

# How do tree-based boosting algorithms address our research question?



## Decision Tree

- A single model of decision-making branches to reach a predicted classification
- Serves as a "baseline" model for our study



## Gradient-Boosting

- Combine multiple models that sequentially learns from past decision trees
- Leads to higher accuracy but can be time-intensive

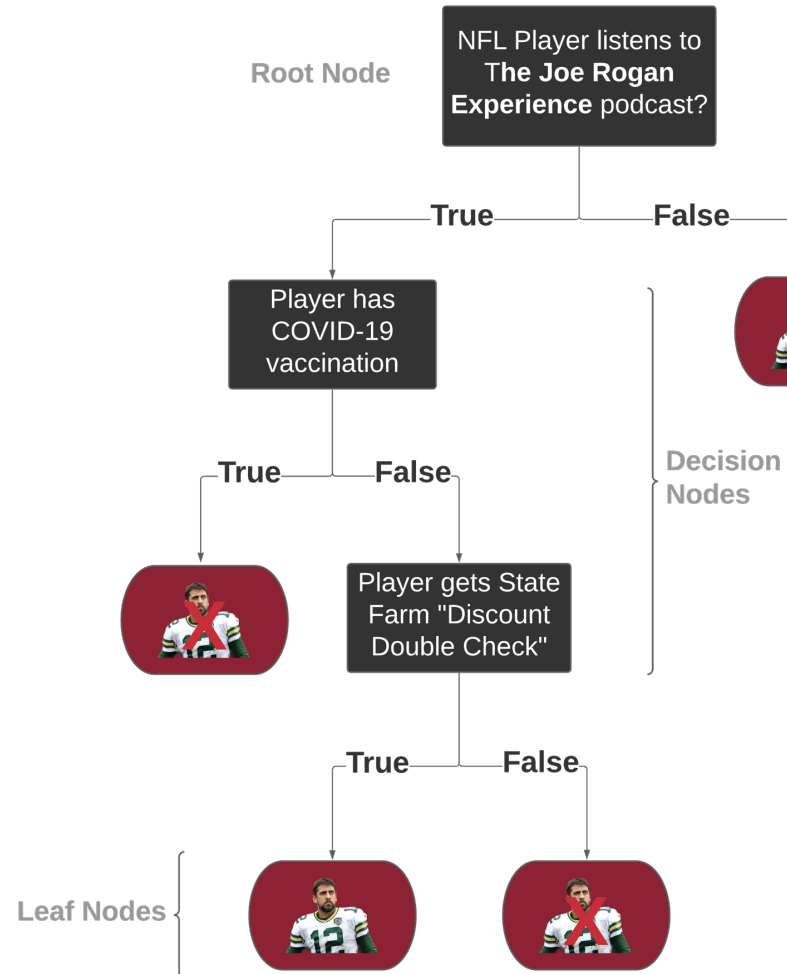
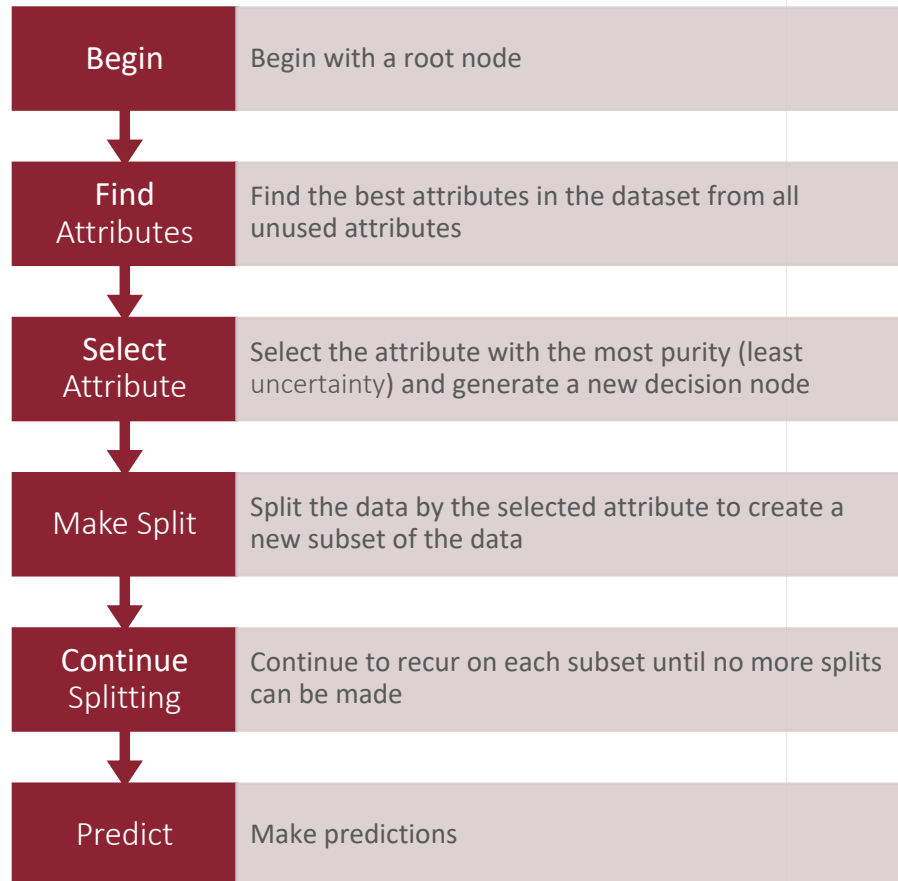


## Extreme Gradient-Boosting (XG-Boost)

- Optimizes for speed, while retaining predictive power of standard gradient-boosting
- Allows for more control in fine-tuning parameters

# Background — Decision Trees

## Principles of Decision Trees (for Classification)



Binary classification decision tree for "Is an NFL player Aaron Rodgers?"



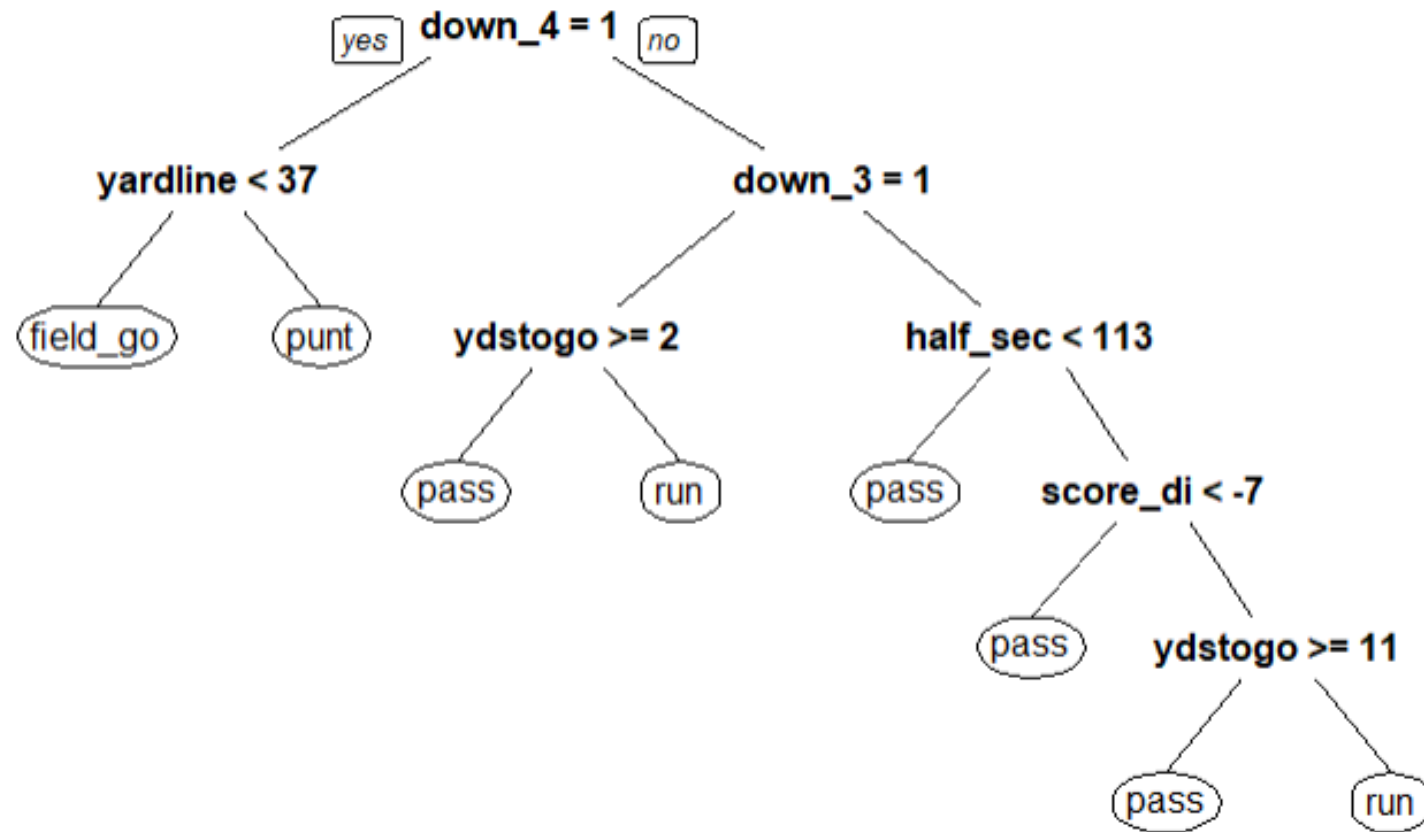
UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE



# Decision Tree

## Baseline Decision Tree

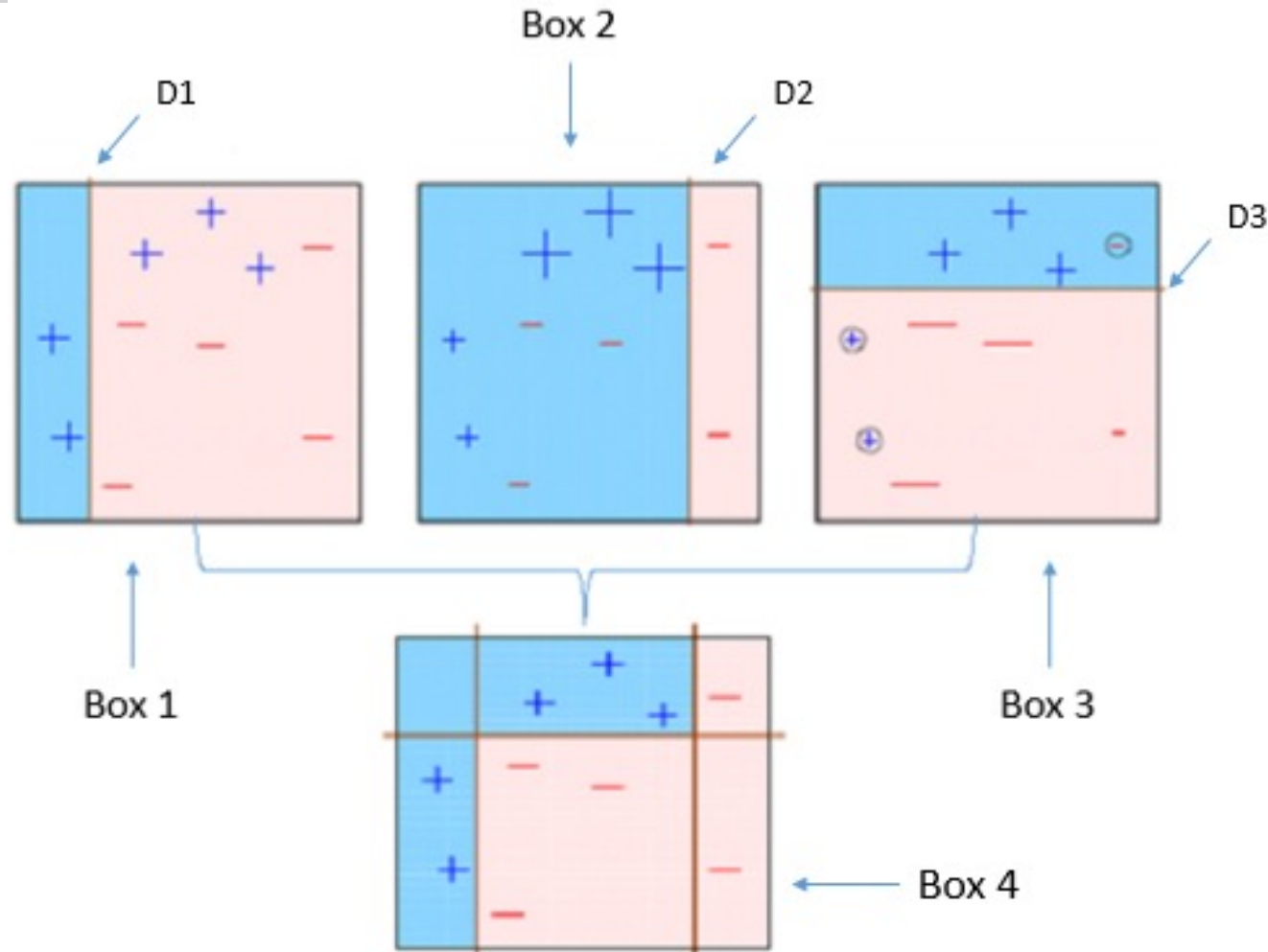


UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE

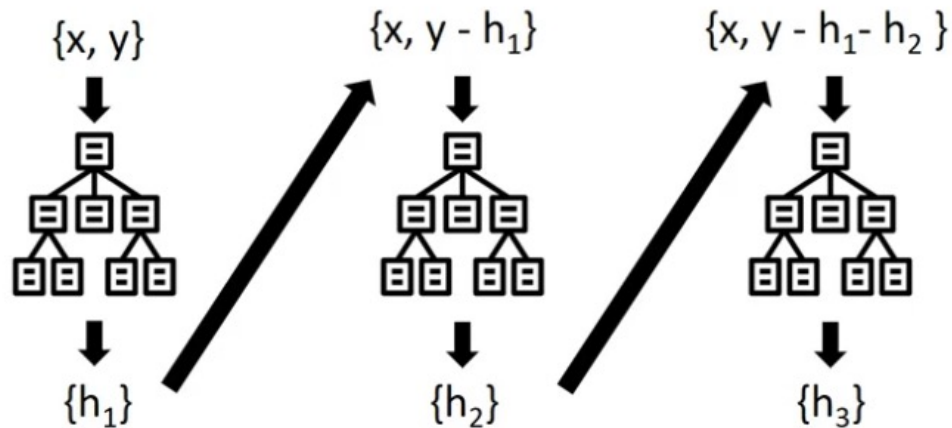
# Primary Method: Boosting

Boosting – An ensemble method using multiple Decision Trees



# Primary Method: Boosting

## Gradient-Boosting



Source: <https://sefiks.com>

- Base model
- Compute residual errors
  - Differentiable loss function (classification: **logloss**)
  - Gradient-Descent
- Parameters to avoid overfitting and output a final model with low variance and low bias

#1

•

•

•

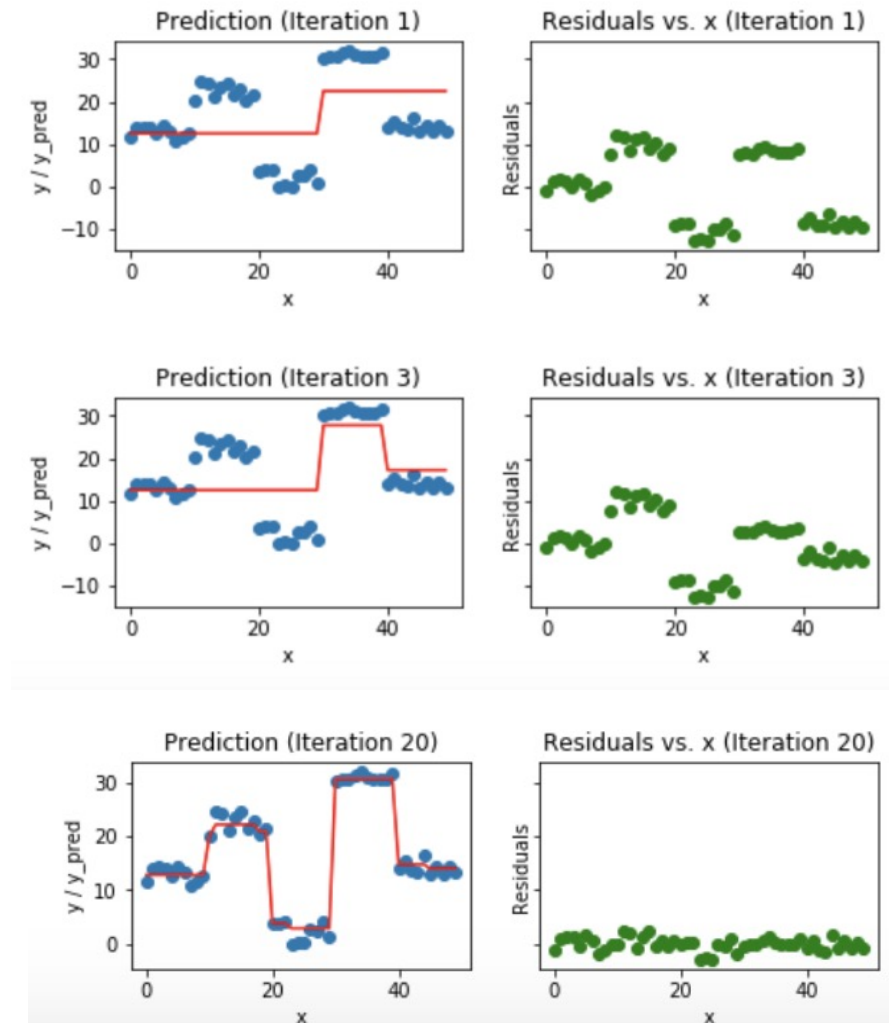
#3

•

•

•

#20

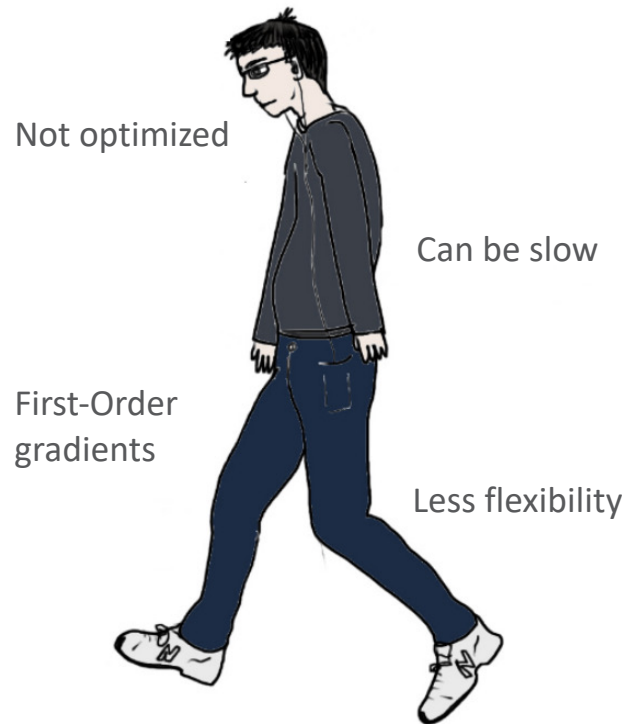


Source: <https://blog.mlreview.com/gradient-boosting-from-scratch-1e317ae4587d>

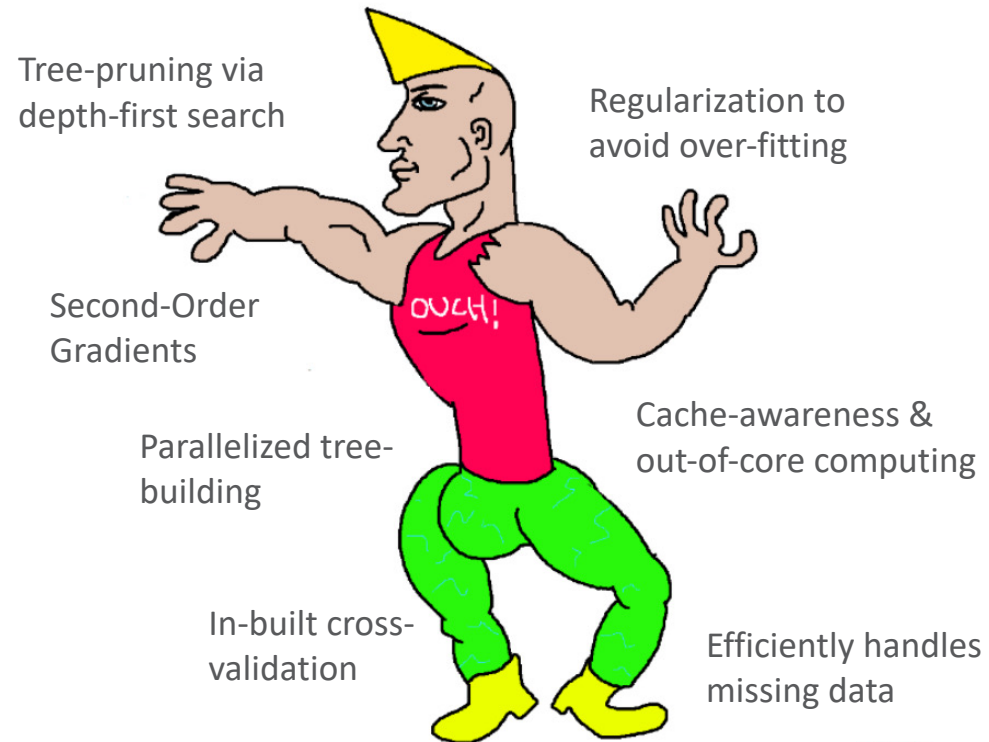
# Primary Method: Boosting

## XG-Boost (Extreme Gradient-Boosting)

### GRADIENT-BOOSTING



### THE "CHAD" XG-BOOST



UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE

# Primary Method: Boosting

## XG-Boost Implementation

Data Requirement	Approach to Satisfy Requirement
Factor variable for classification problems	<ul style="list-style-type: none"><li>- Predictor Factors: dummy-vectorized down</li><li>- Remaining variables kept in numeric form</li></ul>
Only numeric vectors ( <u>must</u> include 0 for converted factors)	Model matrix form of predictors and outcome converted to specialized matrix in xgb package (xgb.DMatrix)
Input data structure for model: sparse matrix (cells containing 0 not stored, so enforcing memory-efficiency)	Encoded Outcome Variable: fg → 0   pass → 1   punt → 2   run → 3

XG-Boost handles the following data characteristics: **correlated features** and **null values**



UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE

# Primary Method: Boosting

## XG-Boost Implementation

Parameter	Option Used	Description
<b>objective</b>	multi:softprob	Outputs predicted probability of observations belonging to each class
<b>eval_metric</b>	mlogloss	Type of metric for validation for each tree; negative log-likelihood for classifications
<b>eta</b>	0.1	Learning rate [0, 1]; lower eta avoids overfitting
<b>max_depth</b>	5	Maximum depth of each tree; lower max_depth avoids overfitting
<b>lambda</b>	1	Regularization term; higher lambda leads to reduce overfitting
<b>gamma</b>	1	Minimum loss reduction to make partition in tree; larger gamma avoids overfitting
<b>nrounds</b>	50	Number of iterations

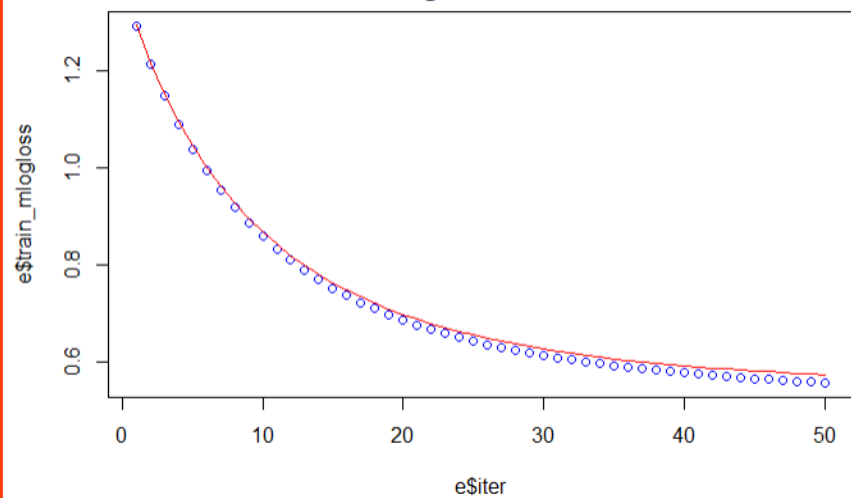
- Other notable parameters: **min\_child\_weight**, **colsample\_bytree**, **subsample**
- **Grid Search** for hyper-parameter tuning



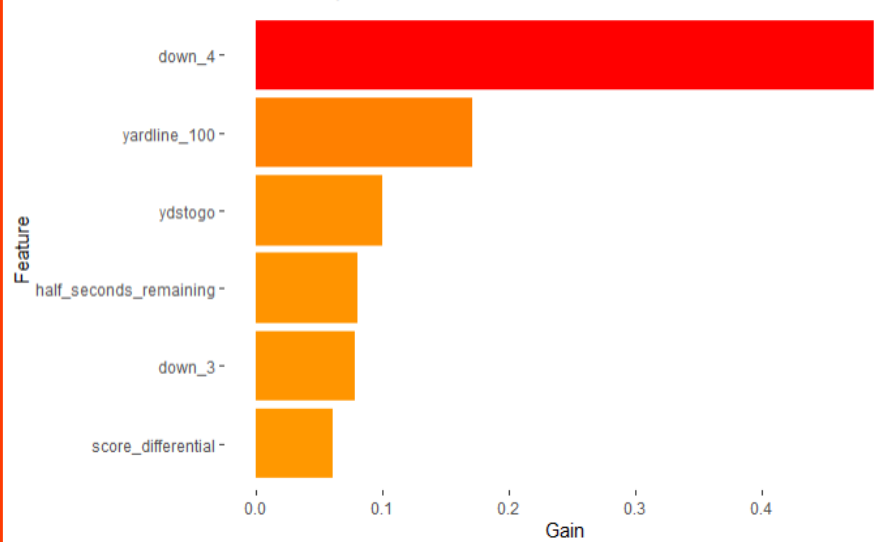
UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE

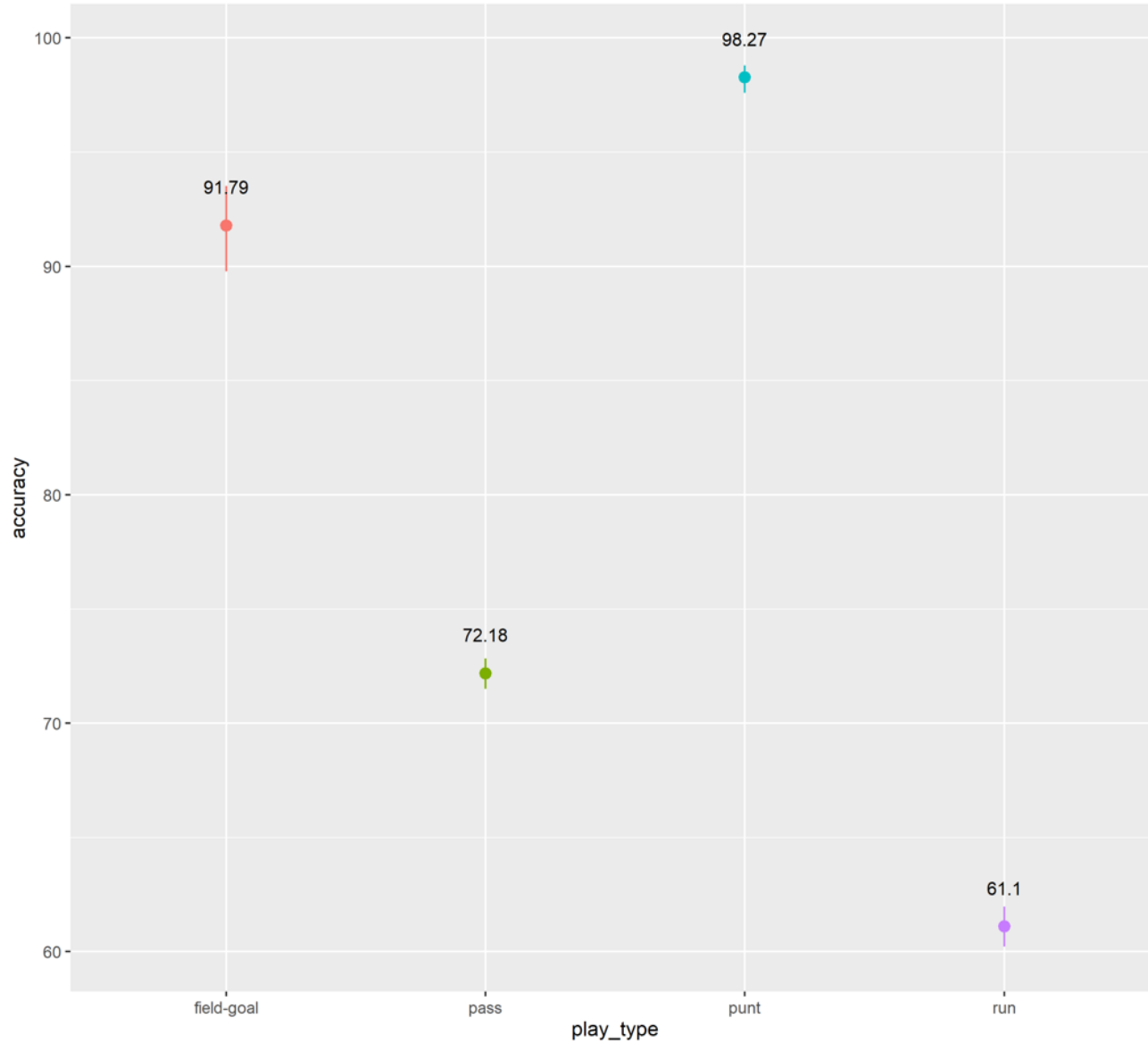
Train vs Test LogLoss Over Each Iteration

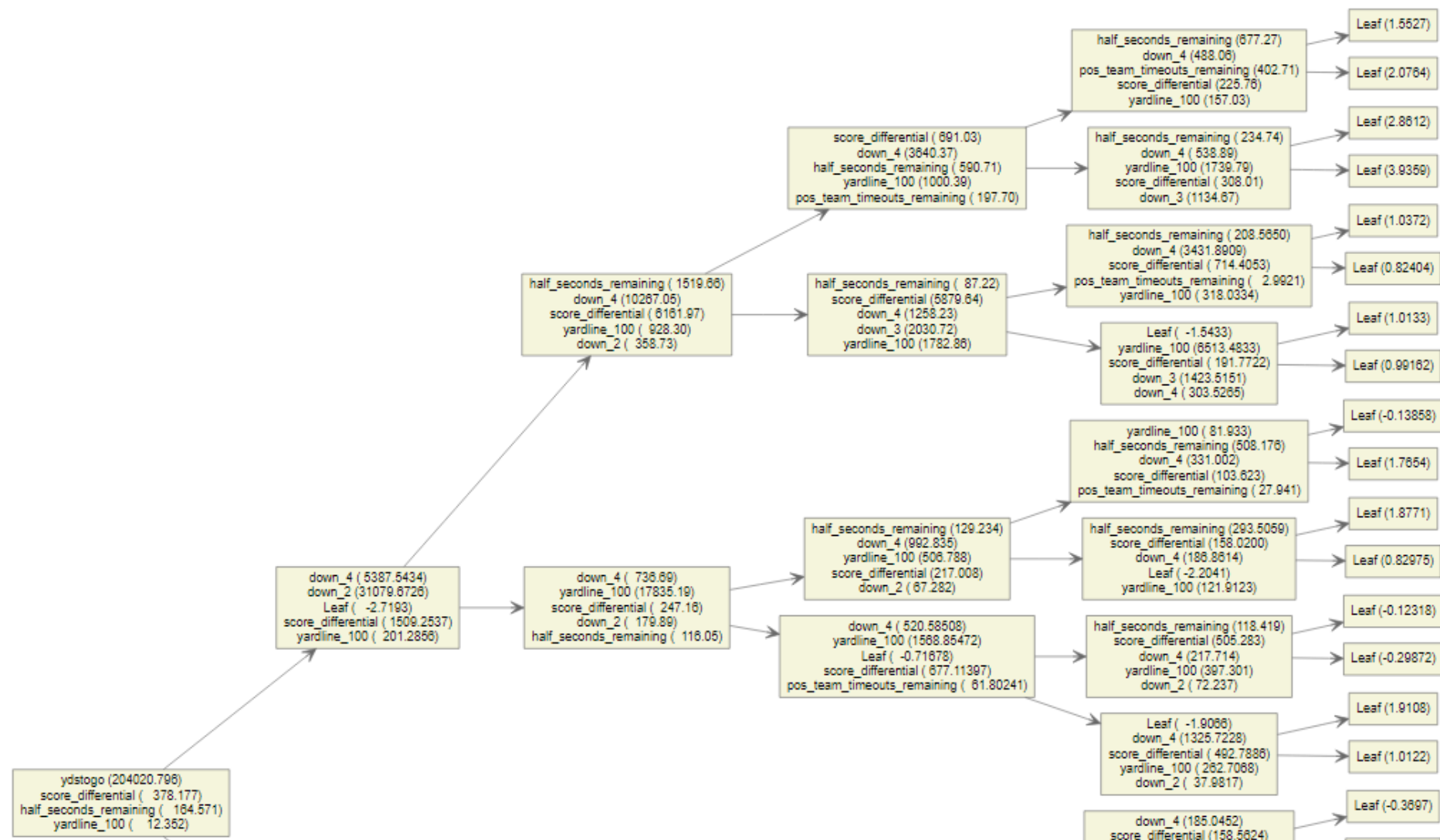


Feature Importance from XG-Boost



Accuracy (%) Breakdown by Play







# Does an increase in sophistication of tree-based boosting algorithms impact the accuracy & efficiency of NFL play predictions?

Method	Classification Accuracy	Runtime
Decision Tree (Baseline)	62.64 %	1.1 sec
Gradient-Boosting	69.01 %	36.8 sec (50 iterations)
XGBoost	70.25 %	2.9 sec (50 iterations)

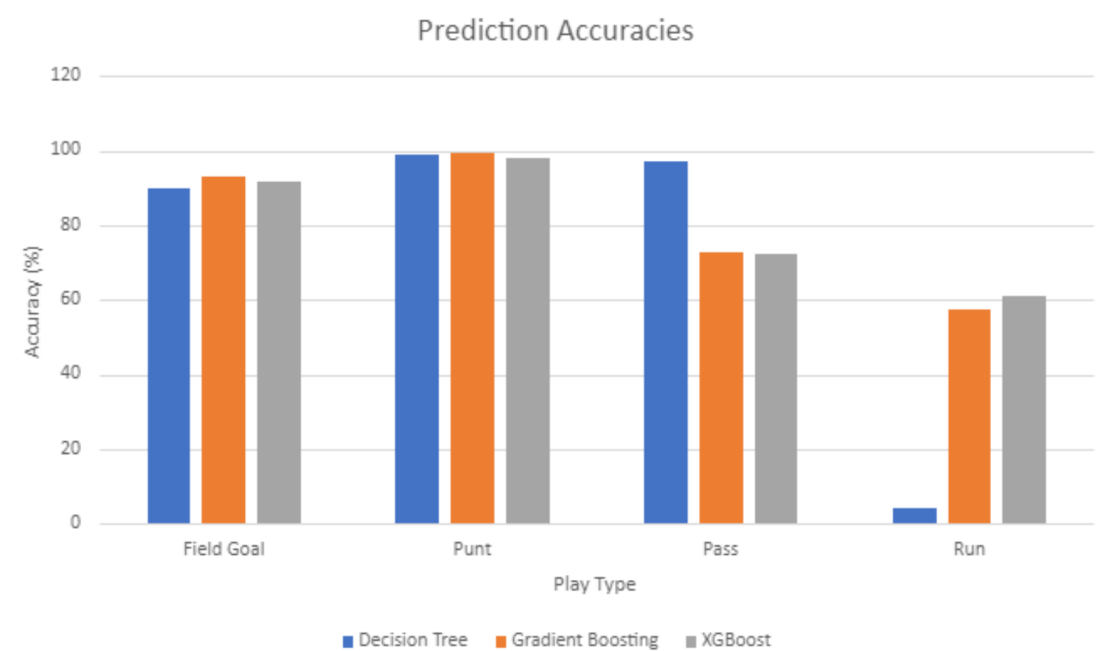
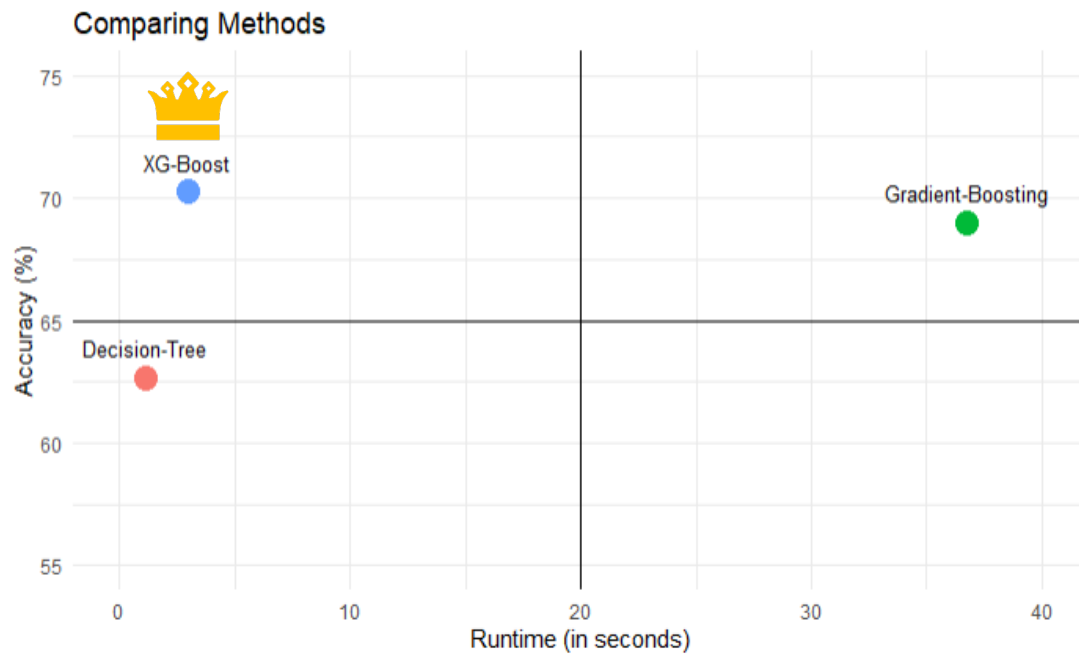


UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE

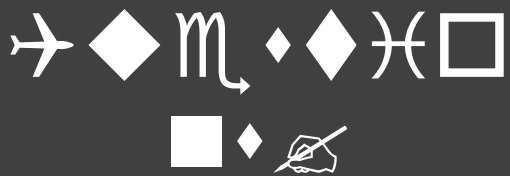
# Conclusions

## Summary of Model Analysis



UNIVERSITY of  
DENVER

DANIEL FELIX RITCHIE SCHOOL  
OF ENGINEERING & COMPUTER SCIENCE



"Questions?"



John Edwards  
@John\_B\_Edwards

was inspired by this to share my own, personal cheat sheet i rely on for my projects in case it's helpful

JOHN EDWARDS'  
PATENTED MACHINE  
LEARNING ALGORITHM  
CHEAT SHEET

