Group
Erick  Calzadilla
Clayton Ragsdale

# What is March Madness

► March Madness- the "holy grail" of basketball tournaments

► 68 teams competing , single game elimination

► Fans and Analysts around the country decide to compete in their own way.

► Scientific formulas, dumb luck, favorite teams, and historical data

# 2015 March Madness Bracket

| First Round March 19-20 | Second Round March 21-22 | Sweet 16 March 26-27 | Elite 8 March 28-29 | Final Four April 4 | Championship April 6 | Final Four April 4 | Elite 8 March 28-29 | Sweet 16 March 26-27 | Second Round March 21-22 | First Round March 19-20 |
|---|---|---|---|---|---|---|---|---|---|---|

| Manhattan (19-13) | BYU (25-9) | March 17-18 | N. Florida (23-11) | Boise St. (25-8) |
|---|---|---|---|---|
| Hampton (16-17) | Ole Miss (20-12) | | R. Morris (19-14) | Dayton (25-8) |
| 16 Seed Midwest | 11 Seed West | | 16 Seed South | 11 Seed East |

## Midwest

- Kentucky (34-0) — 1 — Kentucky
- Hampton (16-17) — 16
- Cincinnati (22-10) — 8 — Cincinnati
- Purdue (21-12) — 9
- West Virginia (23-9) — 5 — West Virginia
- Buffalo (23-9) — 12
- Maryland (27-6) — 4 — Maryland
- Valparaiso (28-5) — 13
- Butler (22-10) — 6 — Butler
- Texas (20-13) — 11
- Notre Dame (29-5) — 3 — Notre Dame
- Northeastern (23-11) — 14
- Wichita St. (28-4) — 7 — Wichita St.
- Indiana (20-13) — 10
- Kansas (26-8) — 2 — Kansas
- New Mexico St. (23-10) — 15

Kentucky — Kentucky — Kentucky

West Virginia — Notre Dame — Notre Dame

Wichita St. — Wichita St.

Kentucky → Wisconsin

## West

- Wisconsin (31-3) — 1 — Wisconsin
- Coastal Carolina (24-9) — 16
- Oregon (25-9) — 8 — Oregon
- Oklahoma St. (18-13) — 9
- Arkansas (26-8) — 5 — Arkansas
- Wofford (28-6) — 12
- UNC (24-11) — 4 — UNC
- Harvard (22-7) — 13
- Xavier (21-13) — 6 — Xavier
- Ole Miss (20-12) — 11
- Baylor (24-9) — 3 — Georgia St.
- Georgia St. (24-9) — 14
- VCU (26-9) — 7 — Ohio St.
- Ohio St. (23-10) — 10
- Arizona (31-3) — 2 — Arizona
- Texas Southern (22-12) — 15

Wisconsin — Wisconsin — Wisconsin

Oregon — UNC

Arkansas — UNC

Xavier — Xavier — Arizona

Georgia St.

Ohio St. — Arizona

Wisconsin → Wisconsin

## Champions

Duke

## East

- Villanova (32-2) — 1 — Villanova
- Lafayette (20-12) — 16
- NC State (20-13) — 8 — NC State
- LSU (22-10) — 9
- N. Iowa (30-3) — 5 — N. Iowa
- Wyoming (25-9) — 12
- Louisville (24-8) — 4 — Louisville
- UC Irvine (21-12) — 13
- Providence (22-11) — 6 — Dayton
- Dayton (25-8) — 11
- Oklahoma (22-10) — 3 — Oklahoma
- Albany (24-8) — 14
- Michigan St. (23-11) — 7 — Michigan St.
- Georgia (21-11) — 10
- Virginia (29-3) — 2 — Virginia
- Belmont (22-10) — 15

NC State — Louisville — Louisville

Louisville

Dayton — Oklahoma — Michigan St.

Michigan St. — Michigan St.

Louisville — Michigan St.

Michigan St. → Duke

## South

- Duke (29-4) — 1 — Duke
- R. Morris (19-14) — 16
- San Diego St. (26-8) — 8 — San Diego St.
- St. Johns (21-11) — 9
- Utah (24-8) — 5 — Utah
- SF Austin (29-4) — 12
- Georgetown (21-10) — 4 — Georgetown
- E. Washington (26-8) — 13
- SMU (27-6) — 6 — UCLA
- UCLA (20-13) — 11
- Iowa St. (25-8) — 3 — UAB
- UAB (19-15) — 14
- Iowa (21-11) — 7 — Iowa
- Davidson (24-7) — 10
- Gonzaga (32-2) — 2 — Gonzaga
- N. Dakota St. (23-9) — 15

Duke — Duke — Duke

San Diego St. — Utah

Utah — Georgetown

UCLA — UCLA — Gonzaga

UAB

Iowa — Gonzaga

Duke → Duke

# Objective

- Out of 67 games try to create a model using historical data

- From that model we will try to predict as many games possible correct

- To be able to predict the 2015 champion Duke University Blue Devils

- In order to establish a highly successful model using Random Forrest algorithm for classification
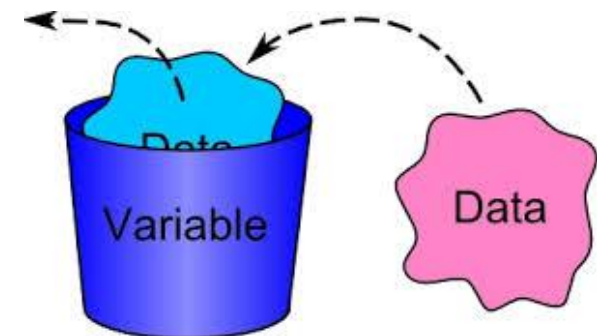
- 1 = Win, 0= Loss

# DATA

- Gathered data from Kaggle
- Regular season data from 2003-2016
- Created Model from half of data, and tested
- -on other half
- Post season data from 2015 tested on
- 71242 observations
- Started with 33 variables
- Ended with 23 variables

# Dependent Variable & Independent Variables

- Dependent Variables

- Win


- Independent Variables

- ftm

- fta

- or

- dr

- ast

- to

- stl

# Random Forrest

- It develops lots of decision tree based on randomly selecting random variables from a random selection of data

- Two principles

- - most of the trees are predicting correctly

- -Trees are making mistakes at different nodes

- Majority Rules

# Error Rate



hihi1000

# VIP



hihi1000

# Training and Testing

- Split the random sample half and half

- The error rate was 10.2% on the training data

- The error rate on the testing data was 8.53%

- Next goal to test on 2015 tournament teams based on their averages

# Explaining Averages

► From the tournament we now use season averages of 2015, instead of post game data to predict outcome

► By doing so we expect a higher error rate since we are now feeding the model average season data instead of post game data

► We took the Average of the 2015 Season

$$\overline{X} = \frac{\Sigma X}{n}$$

# By Round Model

- Each round was run separately, and does not include wins from previous rounds.

- Even if game in previous round is incorrect, the next round is reset with all correct winners.

- This allows model to perform better than it would on the traditional bracket

# Bracket by Round

65.67% Correct

# Results by Round

- ► Winner Duke Blue Devils
- ► ¼ - play in round
- ► 22/32– first round
- ► 11/16- second round
- ► 6/8- Sweet Sixteen
- ► ¾- Elite Eight
- ► ½-Final Four
- ► Winner predicted correctly

# Traditional Tournament

- Single game elimination

- Fill out bracket entirely even if previous winners are incorrect

- Assumption: error rate will be greater than by Round bracket since games are conditional upon previous rounds

# Traditional Bracket



2015 NCAA Division I MEN'S BASKETBALL CHAMPIONSHIP BRACKET

# Results by Traditional

- Winner Duke Blue Devils

- ¼ - play in round

- 21/32– first round

- 10/16- second round

- 6/8- Sweet Sixteen

- 2/4- Elite Eight

- ½-Final Four

- Winner predicted correctly

# How to improve model ?

▶ For loop can be ran in order to calculate averages that are updated every game per team.

▶ Model would now resemble more of our test data format

▶ Create a model in which in can do team by team comparisons.

▶ - or create model where it compares different conferences

▶ -assumption: conferences has specific style

▶ Each team has their own model would then predict winner off of predicted independent variables .

# What to take away

▶ Winners heavily dependent upon shooting percentage and there opponents overall shooting percentage.

▶ Though some predictions were incorrect since teams with a lower strength of schedule had higher shooting percentages for playing weaker teams

▶ Ex-New Mexico State beats Kansas

▶ Overall Random Forest algorithm provides exceptional results considering lack of data and descriptive independent variables.

▶ Shows us the power of many decision trees

# Questions?