

PREDICTING NCAA MEN'S BASKETBALL TOURNAMENT RESULTS

Frederick McCollum



1

OVERVIEW



Background



Objectives



Methodology



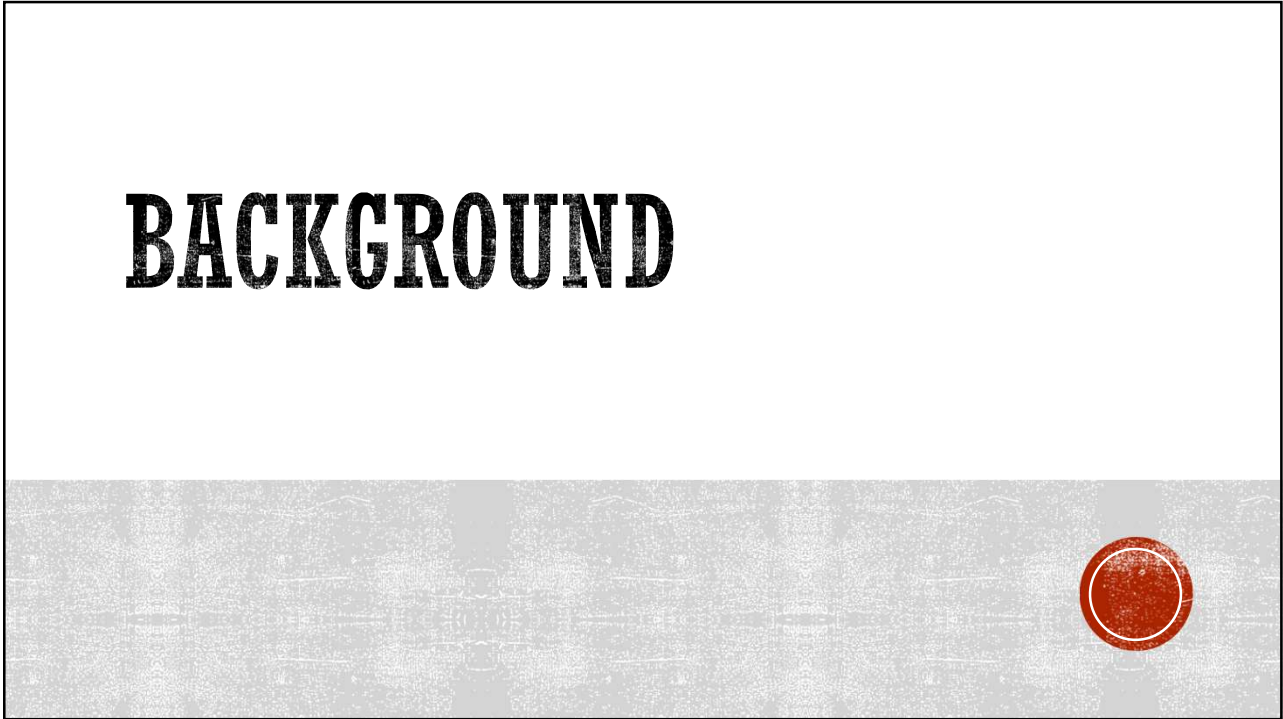
Results



Uses of this
Model



2



3



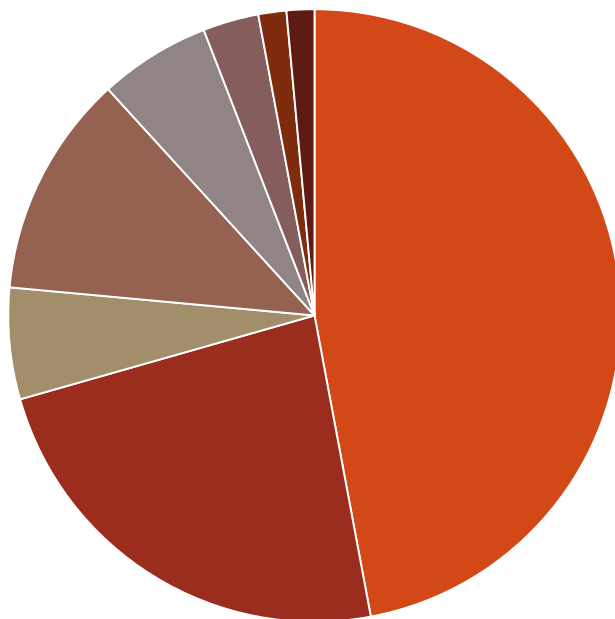
4

DATA OVERVIEW

- Data obtained from Kaggle (Andrew Sundberg)
- NCAA men's basketball team statistics
- 2015 – 2019 (5 years)
- Contains team's final result



5



■ R64 ■ R32 ■ R68 ■ S16 ■ E8 ■ F4 ■ 2ND ■ 1ST

DATA DISTRIBUTION



6

OBJECTIVES



7



Predict NCAA Men's
Basketball Tournament results



Identify the best type of
predictive model



Determine variables with the
most significant impact on
predictions

OBJECTIVES



8

METHODOLOGY



9



1. Load data



2. Select features



3. Preprocess data



4 Train model(s)



5. Test model(s)



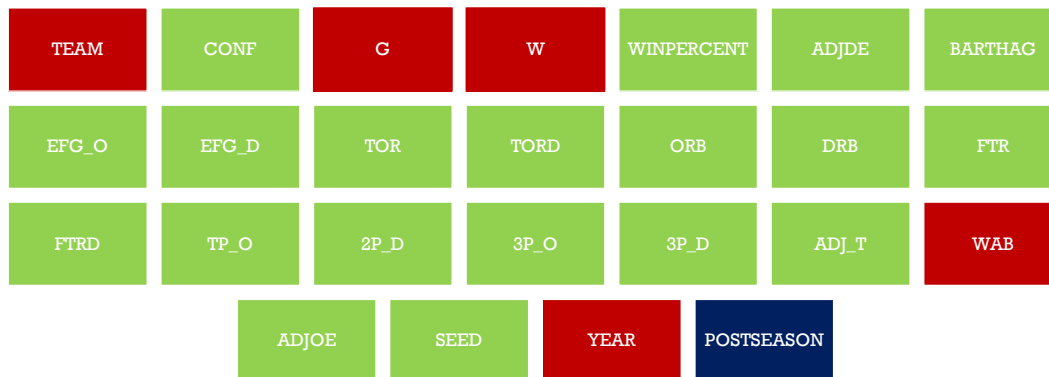
6. Compare results

METHODOLOGY OVERVIEW



10

FEATURE SELECTION



11

DATA WAS PREPROCESSED TO INCREASE PREDICTIVE EFFICACY OF THE DATA

- Identify and remove near-zero variance predictors (none removed)
- Remove highly correlated predictors (none removed)
- Center and scale numeric predictors
- Apply Yeo-Johnson transformation to numeric predictors
- Convert factor variables to dummy variables



12

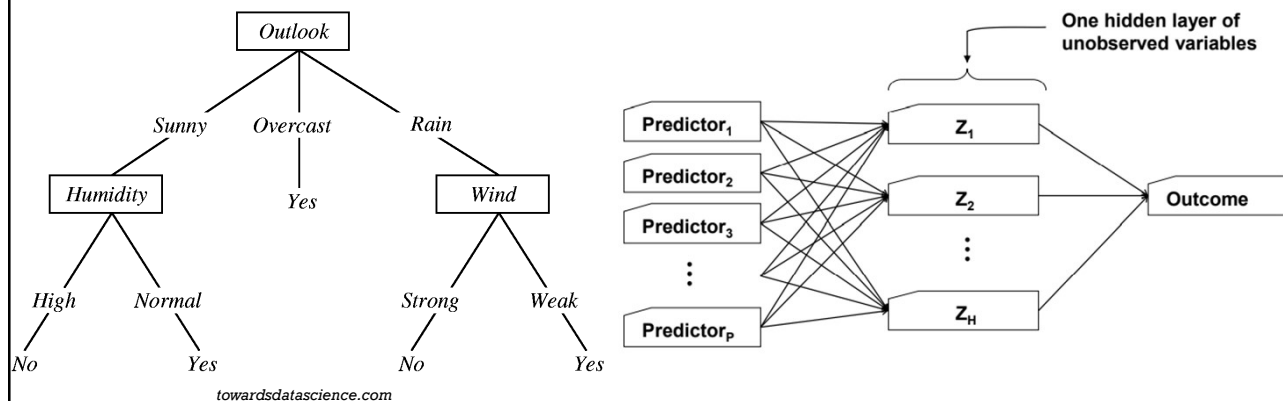
THREE DIFFERENT MODELS WERE CREATED

RANDOM FOREST	NEURAL NETWORK	STOCHASTIC GRADIENT BOOSTING
<ul style="list-style-type: none"> Decision tree based Ensemble model (aggregate) Random sample each time 	<ul style="list-style-type: none"> Constructed to resemble human brain Highly complex Large “network” of decision-making functions 	<ul style="list-style-type: none"> Decision tree based Ensemble model (stepwise) Each tree improves on the last



13

THREE DIFFERENT MODELS WERE CREATED



14

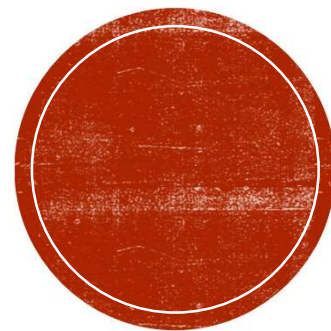
THREE DIFFERENT MODELS WERE CREATED

RANDOM FOREST	NEURAL NETWORK	STOCHASTIC GRADIENT BOOSTING
Pros <ul style="list-style-type: none"> • Less prone to overfitting • Less variance • Handles noisy data well Cons <ul style="list-style-type: none"> • May not perform as well for regression models • May not learn complexities as well as GBM 	Pros <ul style="list-style-type: none"> • Powerful learner • Flexible Cons <ul style="list-style-type: none"> • Complexity • Require more data to be effective • Processing power 	Pros <ul style="list-style-type: none"> • Variance similar to random forest • Perform well for unbalanced data Cons <ul style="list-style-type: none"> • More prone to overfitting than RF • More difficult to tune than RF



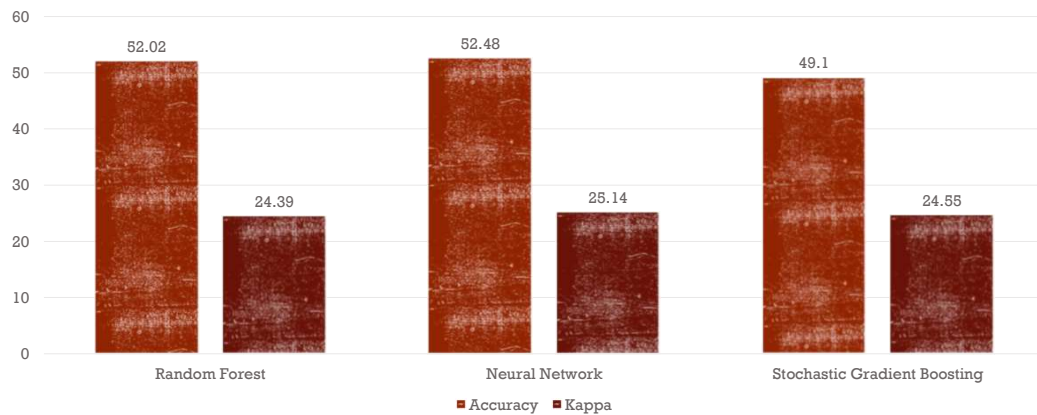
15

RESULTS



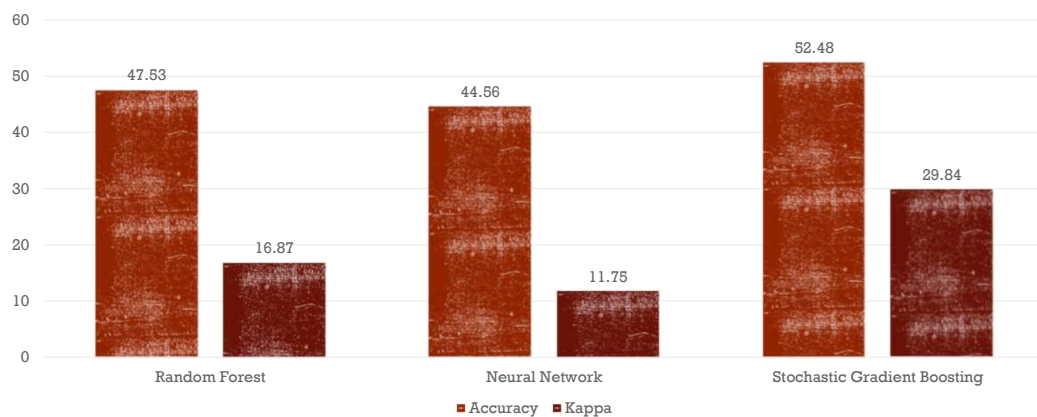
16

MODEL RESULTS — TRAINING DATA



17

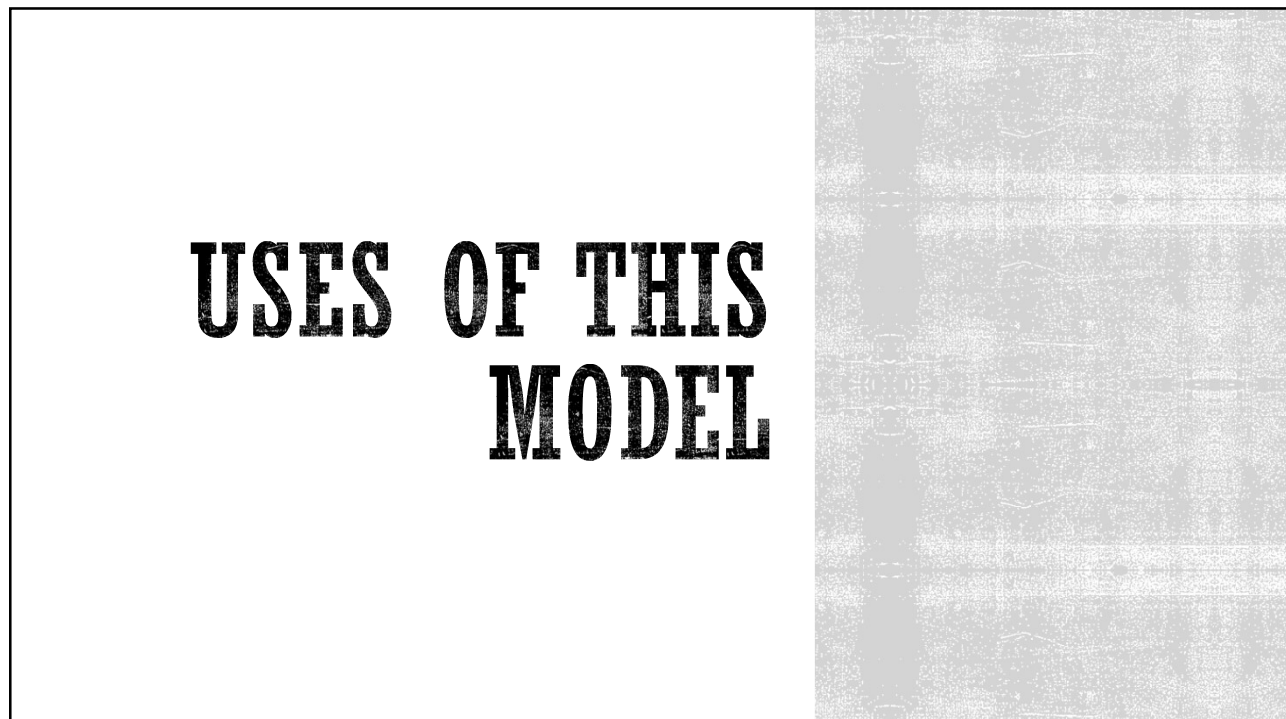
MODEL RESULTS — TEST DATA



18



19



20

POTENTIAL MODEL USES



BRACKET CHALLENGE
COMPETITIONS



TV NETWORK VIEWERSHIP
/ REVENUE FORECASTING



FUN



21



22



23