



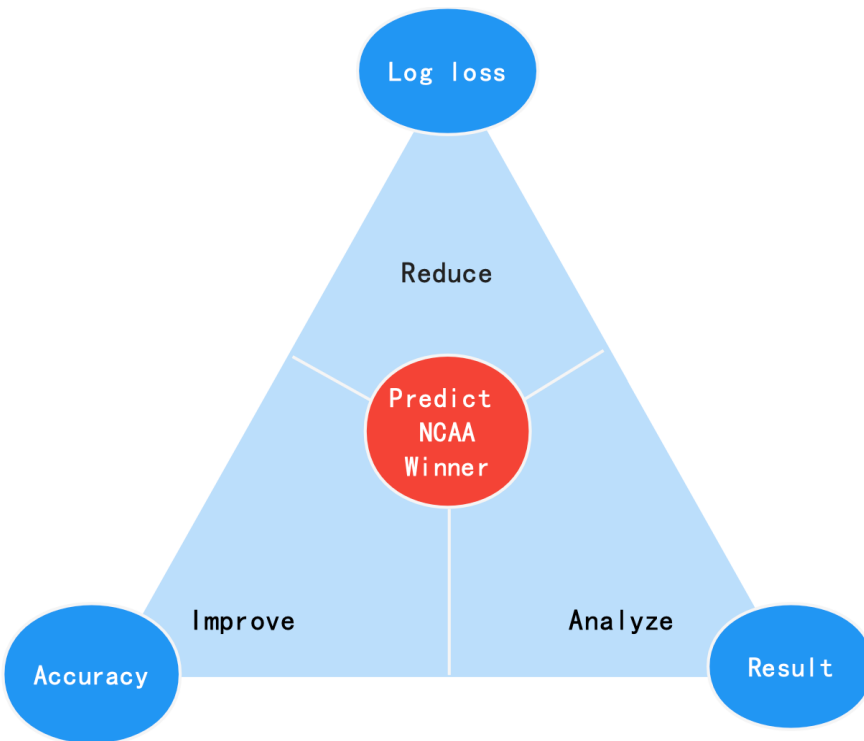
Presented by Goal Diggers

- Nancy Wang
- Stephanie Zhao
- Tony Wang
- Vandana Agarwal



2022 March Madness: NCAA Men's Tournament Bids Report

Objectives



We predict NCAA winner with following objectives:



1. Reduce log loss



2. Improve accuracy



3. Analyze result

Hypotheses



The team's seed number has a great impact in their probability of winning because a higher seeding means a weaker opponent and a home court advantage

Offensive Rebounds Off, Defense Rebounds Def, Tempo Factors play an equally important role in each game outcomes

We ignore the impact of athlete turnover due to graduation on the level of the team Predictors will remain sustainable in the foreseeable future.

Challenges

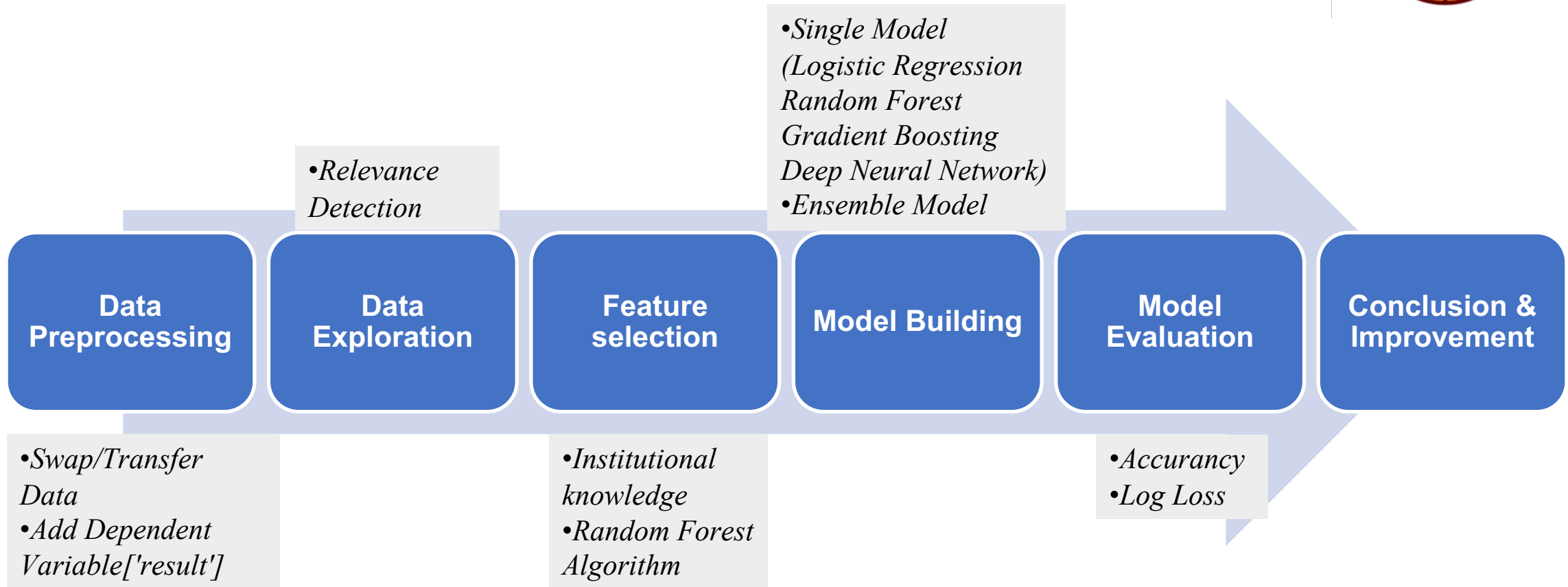


Our dataset has 104 features, and determining the impact of these features on the outcome of each match may be a complex process

All the observations share the same result (team1 wins) which will not satisfy the rule of model training part

A single model may lead to some errors in the results due to its limitations so that we need to build an ensemble model to improve the accuracy of our prediction

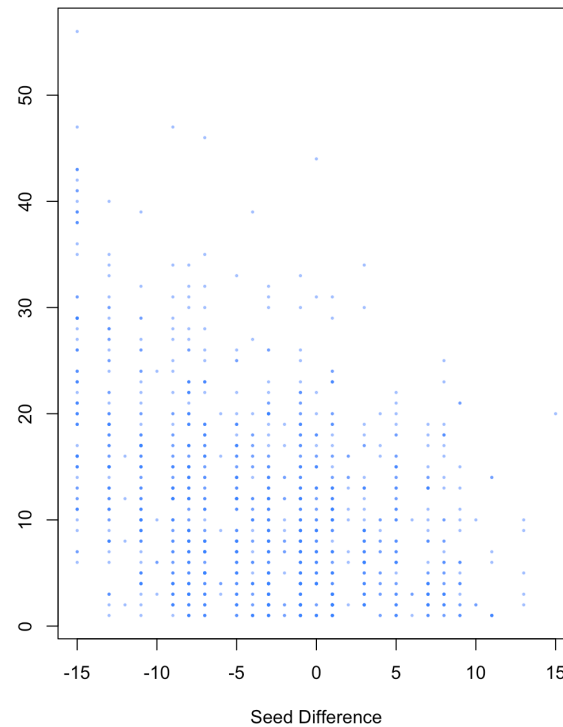
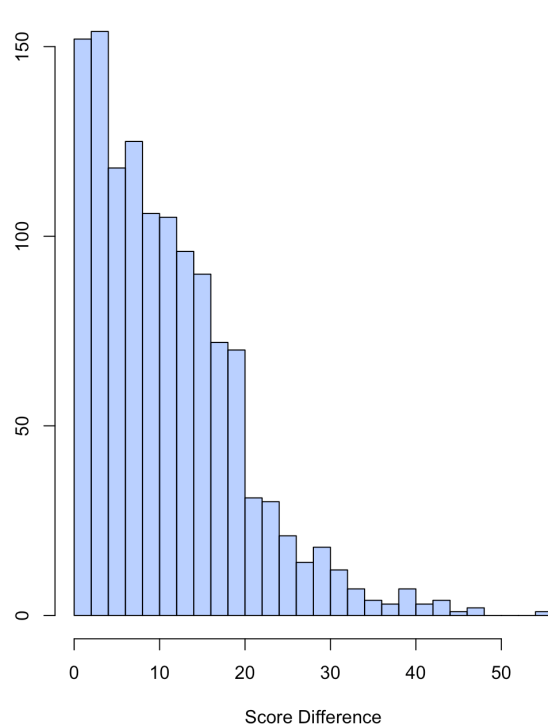
Methodology



Data Exploratory



Score Difference : Distribution and correlation with Seed Difference



Findings:

Score difference shows exponential distribution

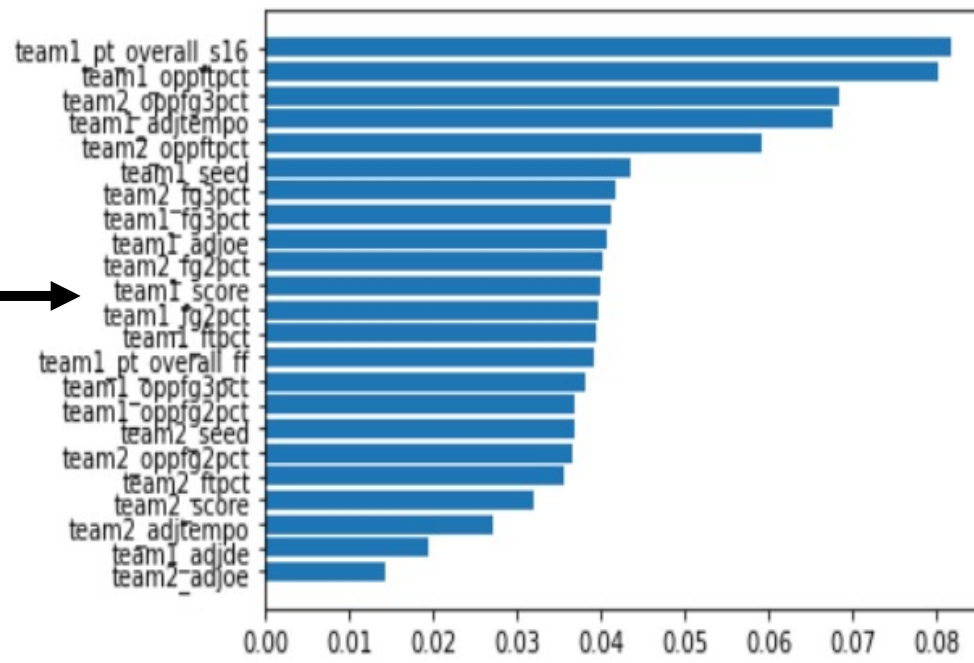
The score difference tends to be large when the seed difference is large negative

** Definition:*

Score difference = team1 score – team2 score

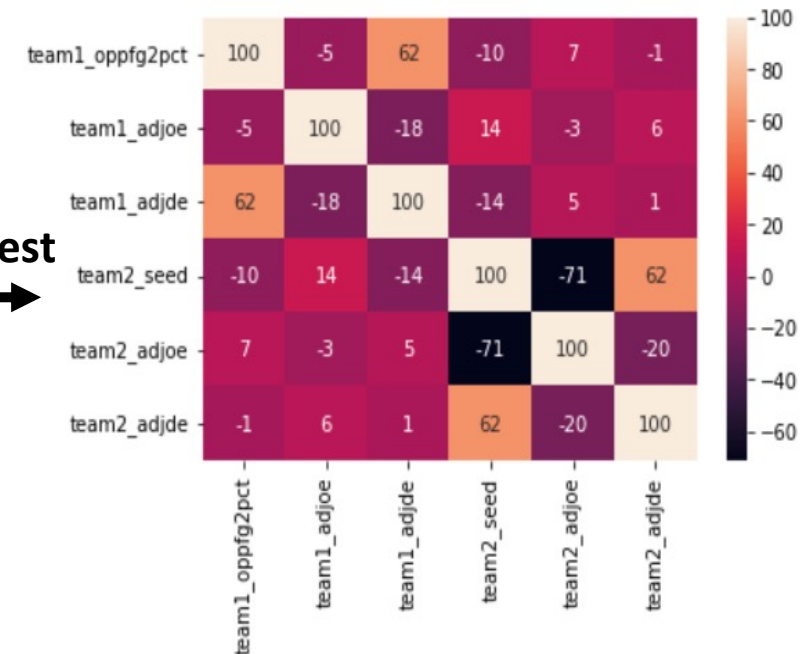
Seed difference = team1 seed – team2 seed

Feature Selection



Random Forest

cut_off = 0.05



Institutional Knowledge

104 Features

24 Features

6 Features

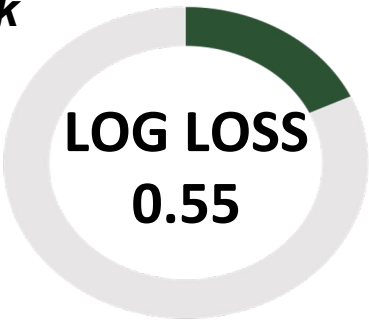
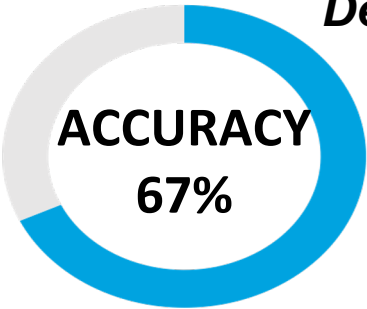
'team1_oppf2pct', 'team1_adjoe', 'team1_adjde', 'team2_seed', 'team2_adjoe', 'team2_adjde'

Single Model Comparison

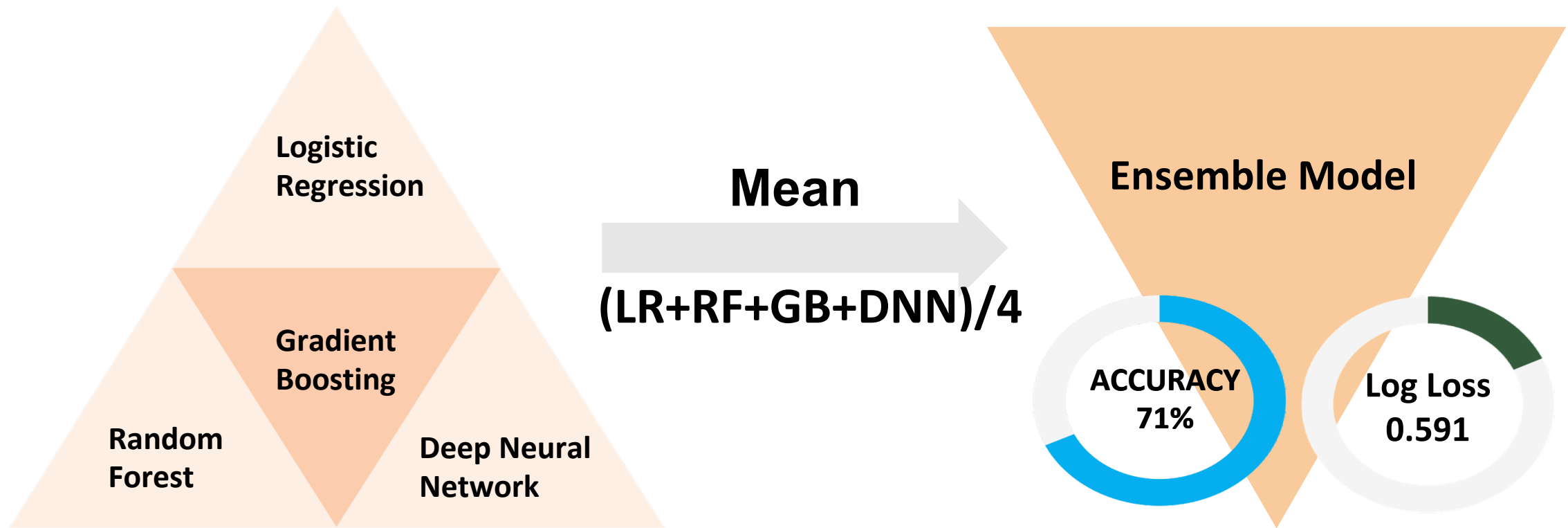


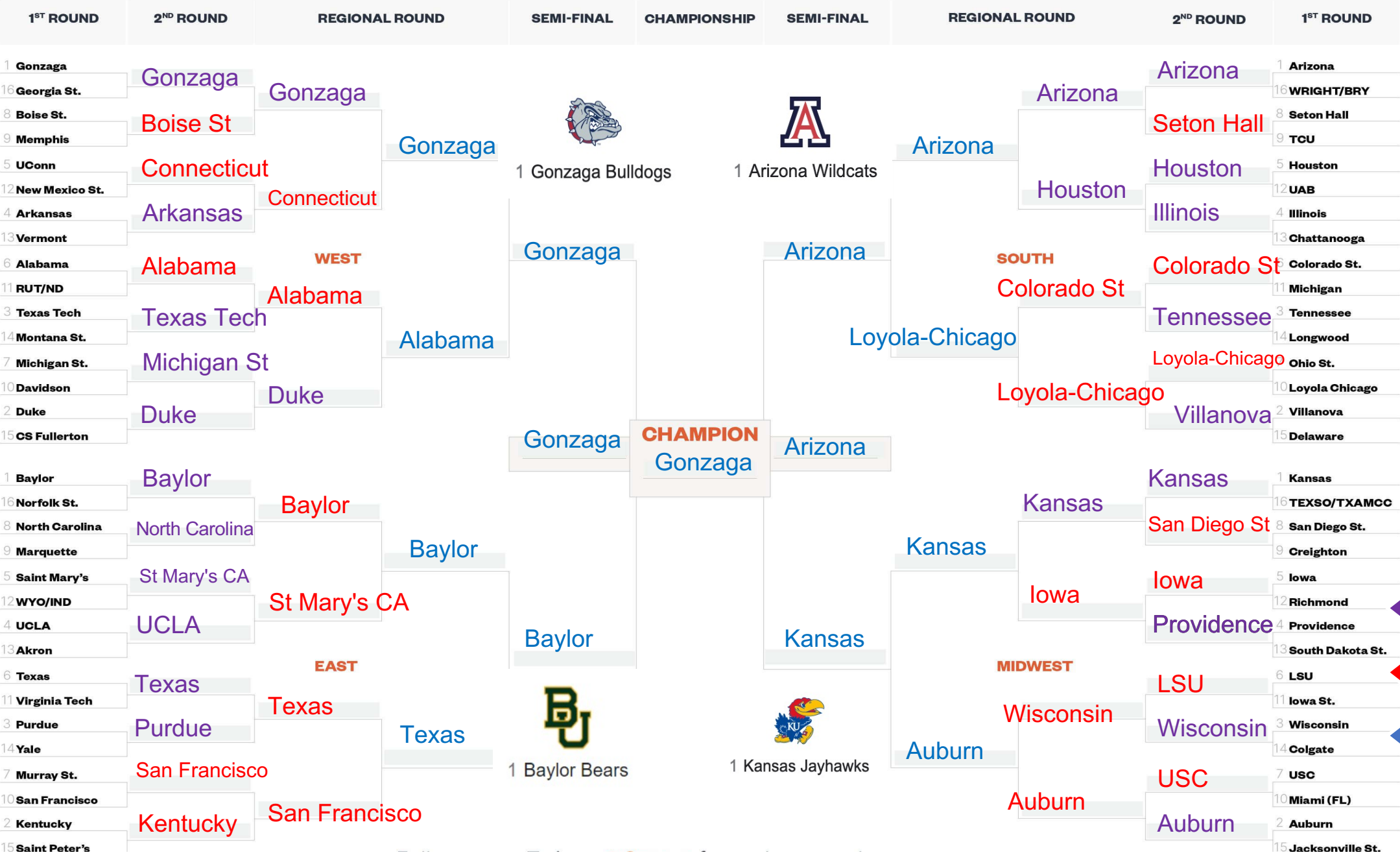
Model	Accuracy	Log Loss	Efficiency(high/modest/low)
<u>Logistic Regression</u>	0.7112	0.5641	HIGH
<u>Deep Neural Network</u>	0.6789	0.5527	HIGH
<u>Random Forest</u>	0.7005	0.6268	MODERATE
<u>Gradient Boosting</u>	0.6658	0.6618	MODERATE
<u>K-Nearest Neighbor</u>	0.7032	0.7316	LOW

Highest Model Performance
Deep Neural Network



Model Building - Ensemble





☐ Date by
March 23



Correct
Prediction



Incorrect
Prediction



Predicted
Top 8 Teams



Final 4 Conclusion



Champion Probability Prediction

89.71%



1 Gonzaga Bulldogs

54.10%



1 Baylor Bears

89.89%



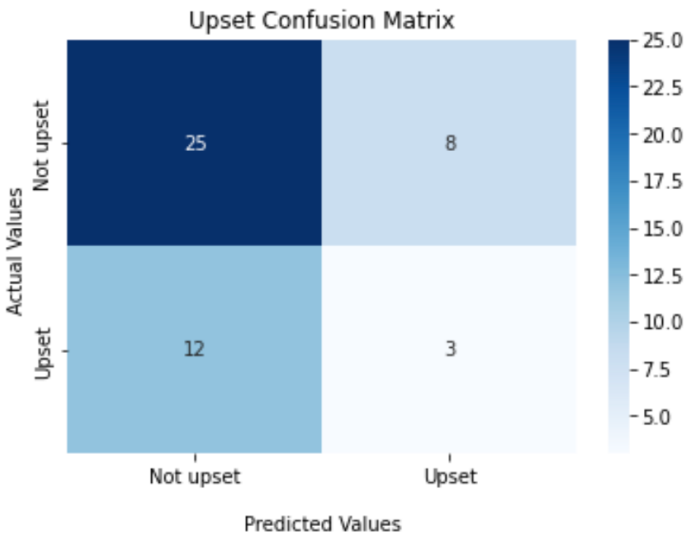
1 Arizona Wildcats

36.06%




1 Kansas Jayhawks


Upset Prediction



PREDICTED FIRST ROUND UPSETS



San Francisco
51.63%
10 vs 7

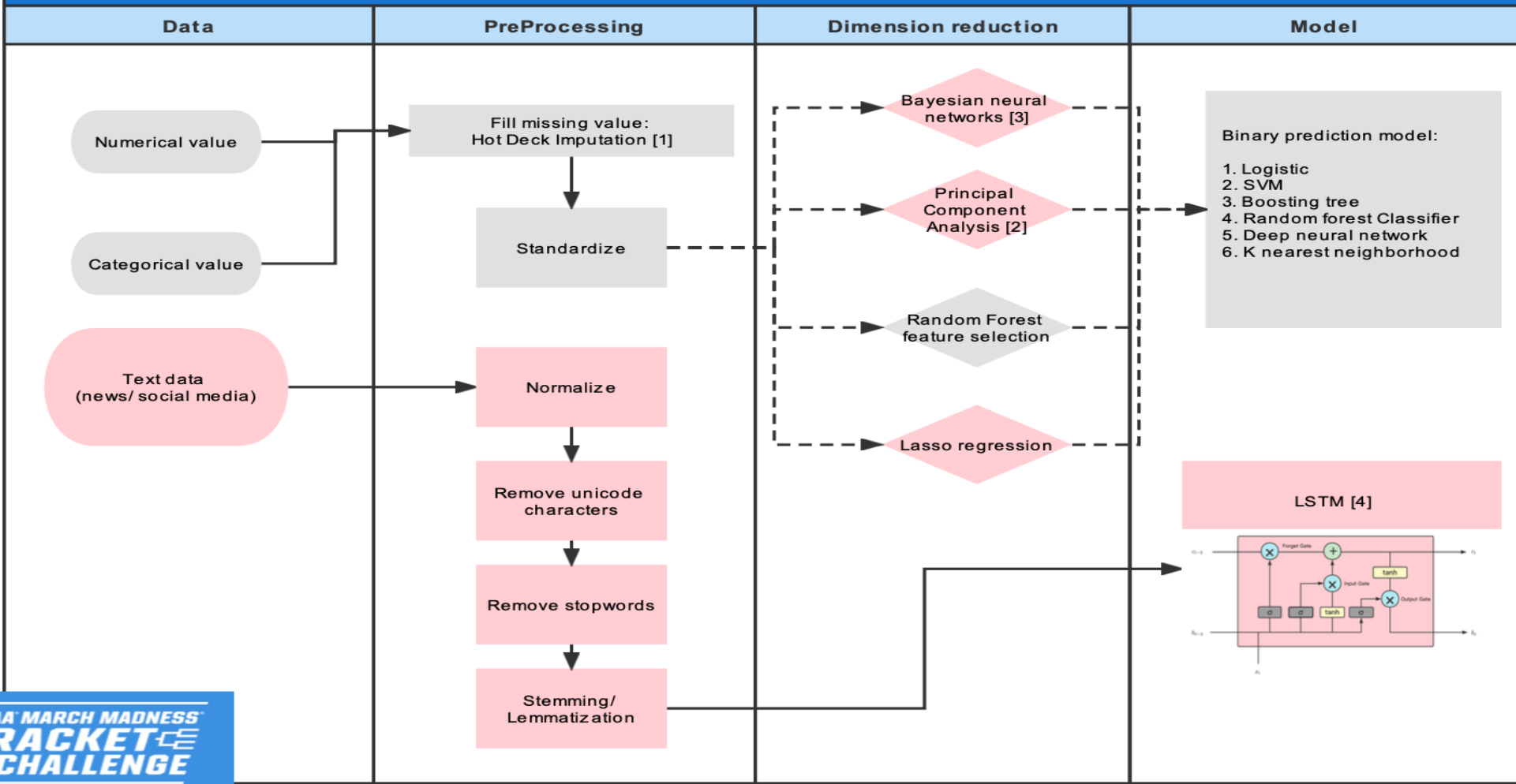


Loyola
53.57%
10 vs 7



Improvement

IMPROVEMENT FLOWCHART



What we have done

What can be improved

Reference



- [1] Andridge, R. R., & Little, R. J. (2010). A review of hot deck imputation for survey non-response. *International statistical review*, 78(1), 40-64.
- [2] Abdi, H., & Williams, L. J. (2010). Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4), 433-459.
- [3] Carreira-Perpinán, M. A. (1997). A review of dimension reduction techniques. Department of Computer Science. University of Sheffield. Tech. Rep. CS-96-09, 9, 1-69.
- [4] Yao, L., & Guan, Y. (2018, December). An improved LSTM structure for natural language processing. In 2018 IEEE International Conference of Safety Produce Informatization (IICSPI) (pp. 565-569). IEEE.
- [5] https://en.wikipedia.org/wiki/NCAA_Division_I_Men%27s_Basketball_Tournament