

# MARCH DATA CRUNCH MADNESS 2023

Predicting NCAA Tournament Champion



Team Jordan Year 

Yujoon Jang  
Jennifer McFadden  
Masaya Sugimoto

Sponsored by: **Deloitte.**

**FORDHAM**  
THE JESUIT UNIVERSITY OF NEW YORK

# Introduction

## Problem Statement

Can we predict the 2023 March Madness NCAA Men's Basketball tournament bracket winners?

## Objective

Use historical 2002-2022 tournament data including team and coach performance and find additional features to build a model that accurately predicts 2023 March Madness NCAA Men's Basketball tournament.

## Methodology

### Data Pre-Processing

- Data Validation
- Feature Engineering

### Feature Selection

- Principal Component Analysis
- Recursive Feature Elimination

### New Features

- Research
- Data Cleaning
- Data Merge
- Correlation

### Model Selection

- Logistic Regression
- Random Forest
- Decision Tree
- XGBoost
- LightGBM

### Model Evaluation

- Log Loss
- ROC AUC
- Precision
- Recall
- F-Measure

# Data Pre-Processing & Feature Engineering

New features derived from existing data to enhance model accuracy

**Pythagorean Win Percentage** (log5)

Team's expected win percentage based on adjusted offensive & defense efficiency

**Distance From Home**

Distance between home court and game location court

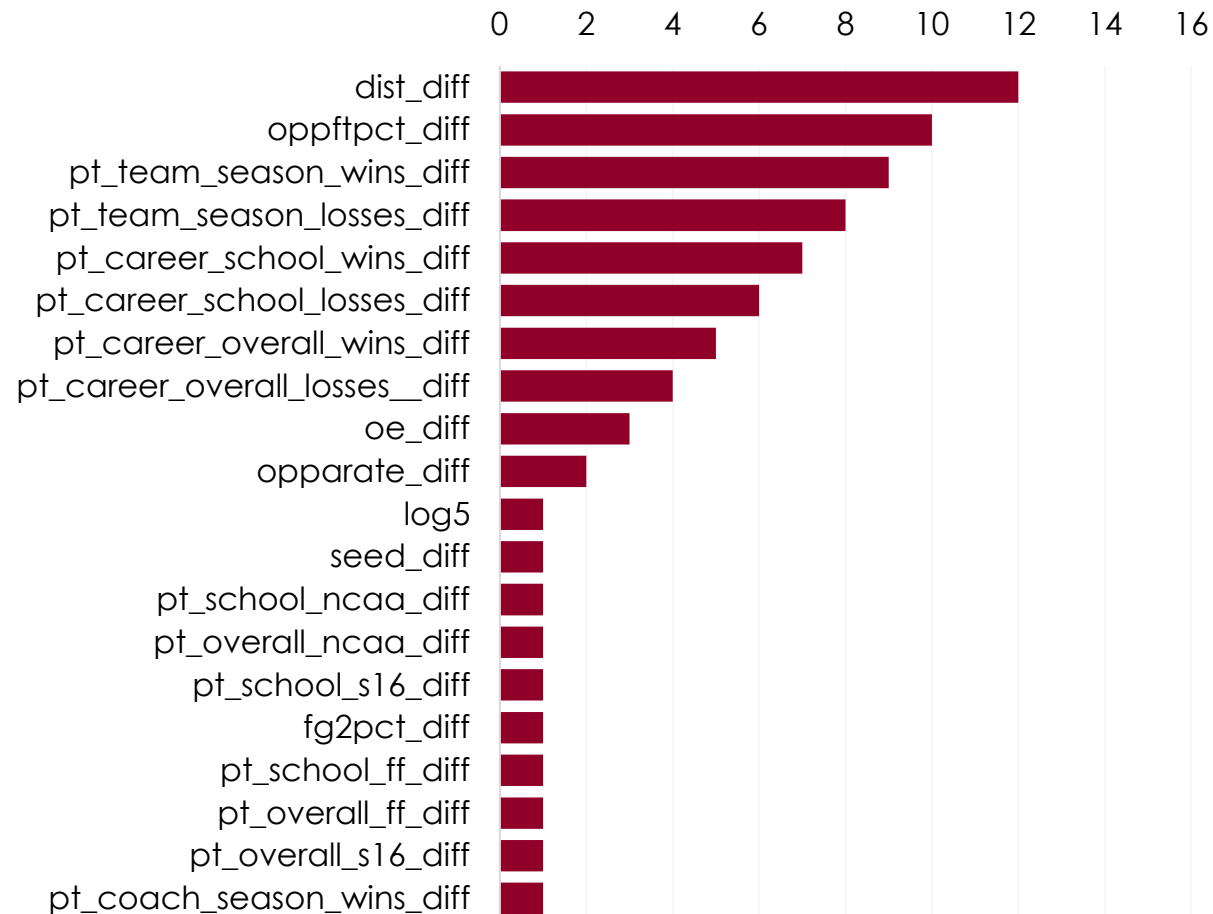
**Team 1 & Team 2 Difference**

Difference between Team 1 and Team 2 for all performance metrics

# Feature Selection

Recursive Feature Engineering (RFE) and Principal Component Analysis (PCA) selection techniques identified top 22 features with predicting potential from initial list of 104

## RFE Feature Ranking



## PCA

### Principal Component 1

Represent 87% of the dataset

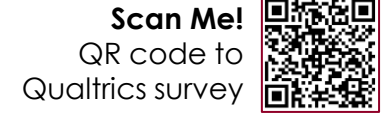
- dist\_diff
- pt\_career\_overall\_wins\_diff
- pt\_career\_school\_wins\_diff
- pt\_coach\_season\_wins\_diff
- pt\_team\_season\_losses\_diff
- seed\_diff
- adjoe\_diff
- adjde\_diff
- oe\_diff

### Principal Component 2

Combined with PC1 represent 95% of the dataset

- pt\_career\_overall\_losses\_diff
- pt\_career\_school\_losses\_diff
- pt\_overall\_ncaa\_diff
- pt\_school\_ncaa\_diff
- pt\_overall\_s16\_diff
- pt\_school\_s16\_diff

# New Feature Research



Conducted preliminary research to understand most impactful indicators for tournament success

## Research

- 1<sup>st</sup> party Qualtrics survey conducted at NCAA regular season game (15 total responses)
- NCAA super-fans family members and friends
- Articles and insights from top analysts

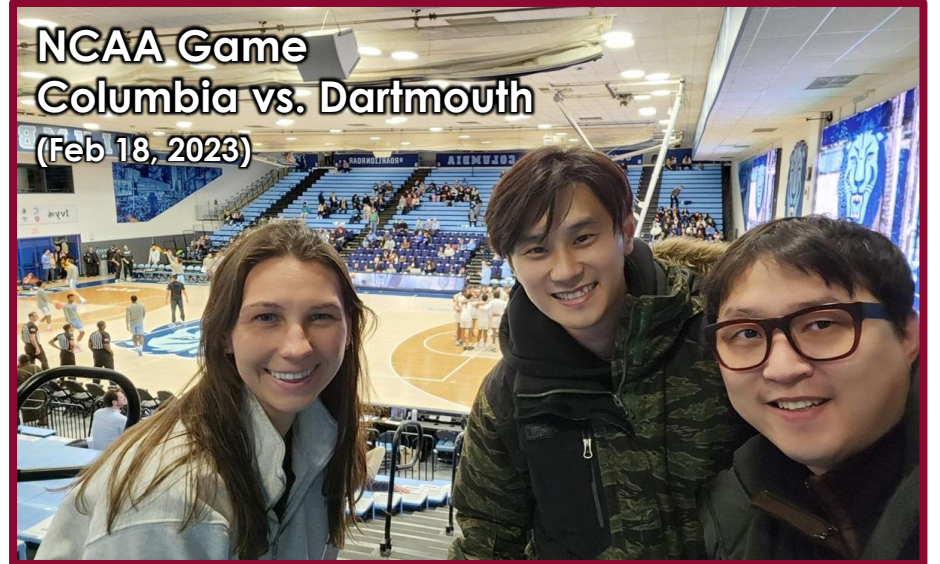


## Top-cited Indicators

- Money
- Team and coach historical performance
- Scoring ability



NCAA Game  
Columbia vs. Dartmouth  
(Feb 18, 2023)



## Areas of opportunity based on data source availability

- 1) Team budget
- 2) Top scorers per team
- 3) Injured players

# New Features



## Expense & Revenue

- Men's team expenses
- Revenue generated by Men's Team

Source: US Department of Education Equity in Athletics Data Analysis



## Star Player Rating

- Top scoring player's avg. points-per-game (ppg) proportion of team avg. points-per-game

Source: NCAA Statistics



## Player Injury Adjustment

- 10% reduction on adjusted offensive & defensive efficiency for season-ending key player injuries

Source: Rotowire & News Coverage

\* Applied to 2023 data only

---

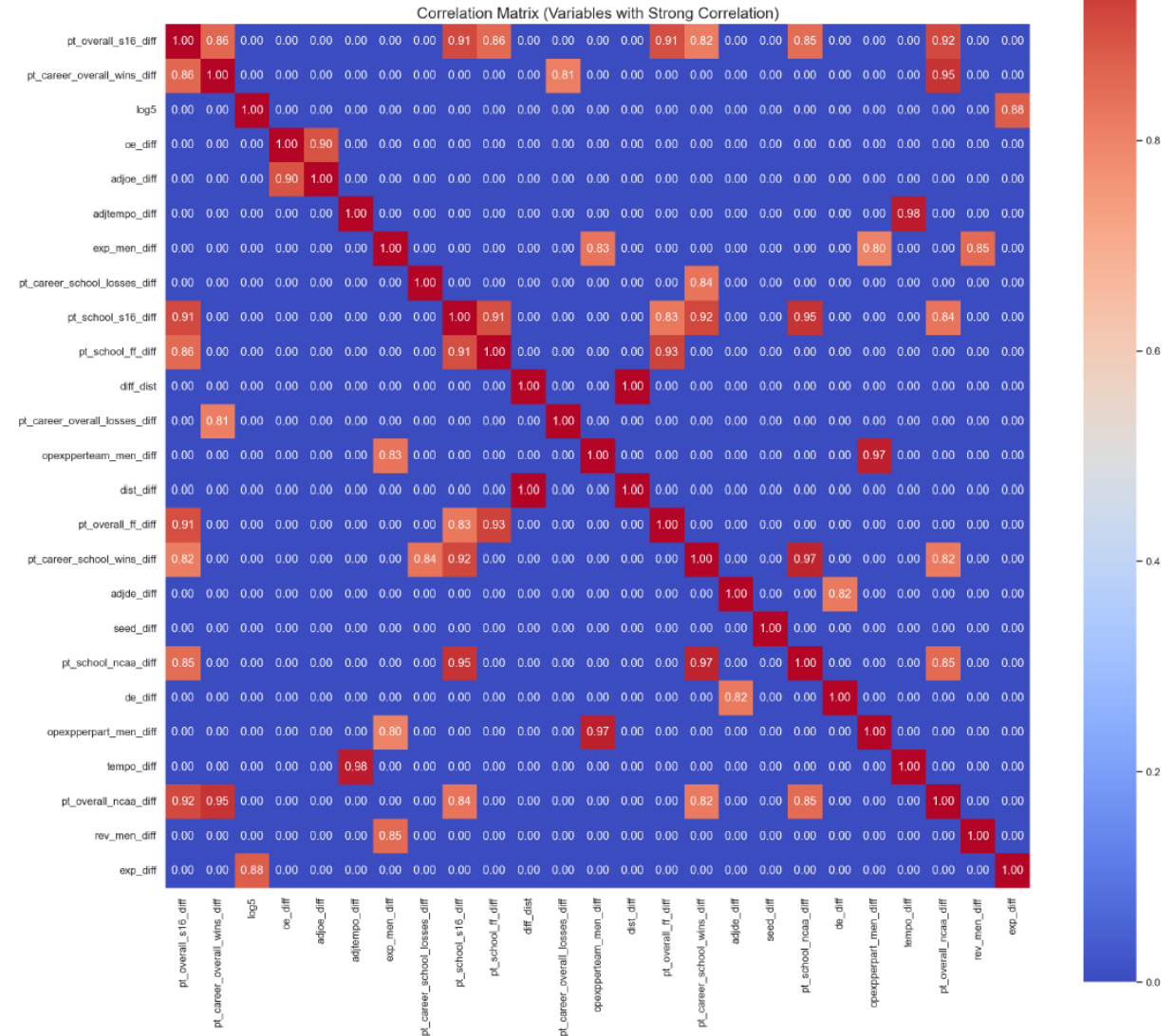
## Enhanced Historical Training Data:

- Historical data for new features only available for 2011-2023, so 2002-2010 season data removed from training set
- Predicting value from new features outweighs loss of smaller training data, and removes potentially outdated seasons

## Feature Selection

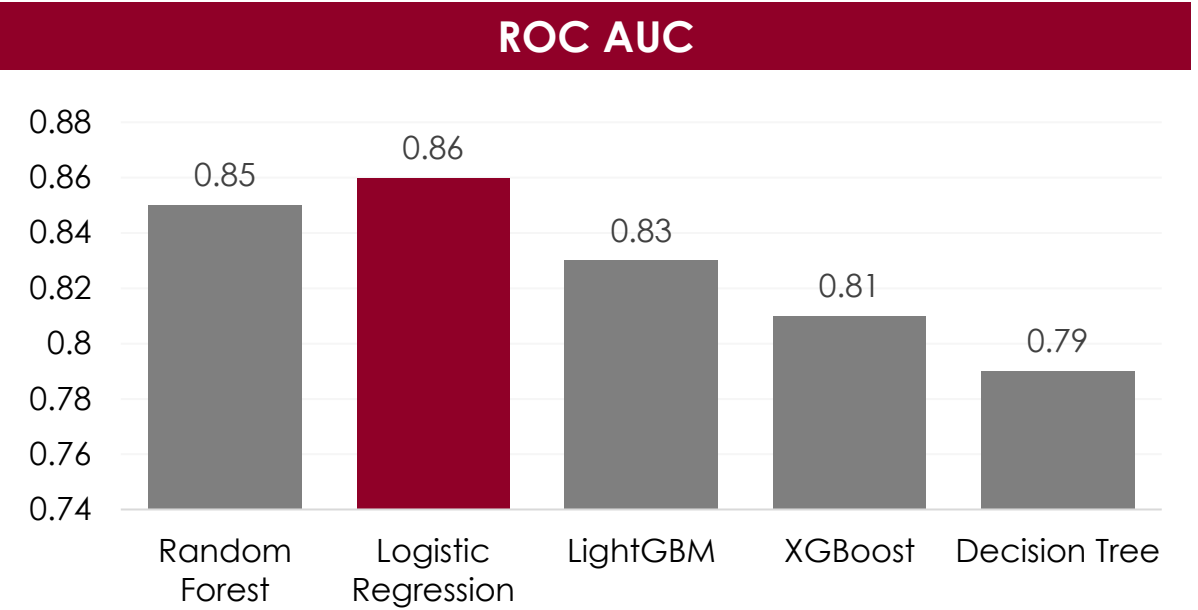
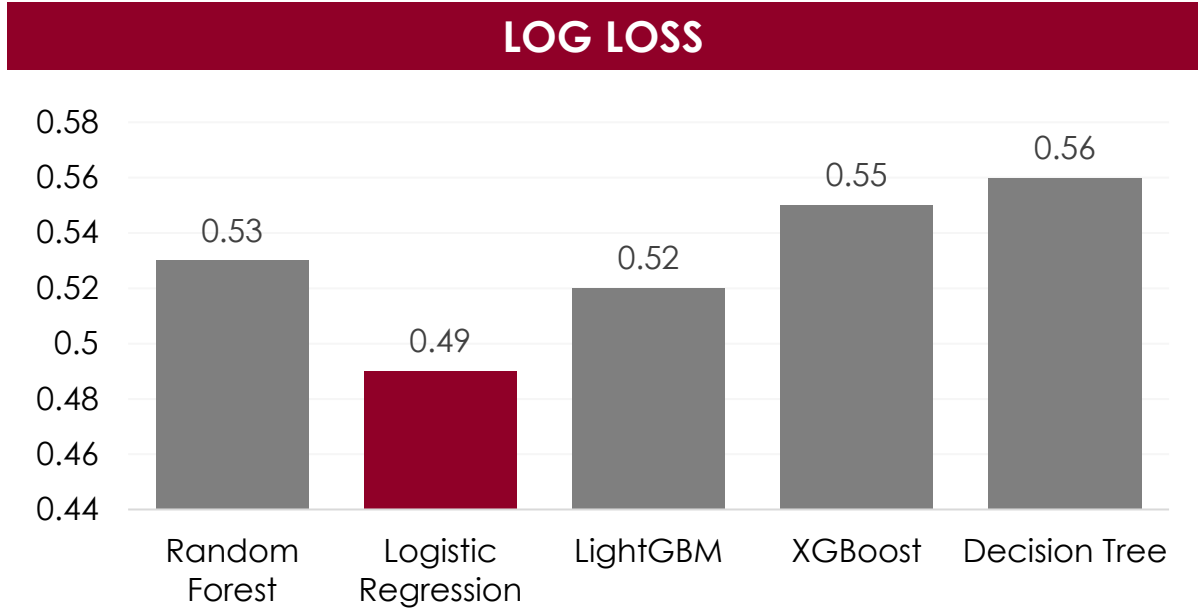
Correlation analysis performed post feature selection and new feature addition refining feature list to 14

- Heat map shows features with correlation greater than 0.8 in red
- High collinearity features removed



# Model Performance Comparison

Out of five models tested, Logistic Regression strongest based on log loss



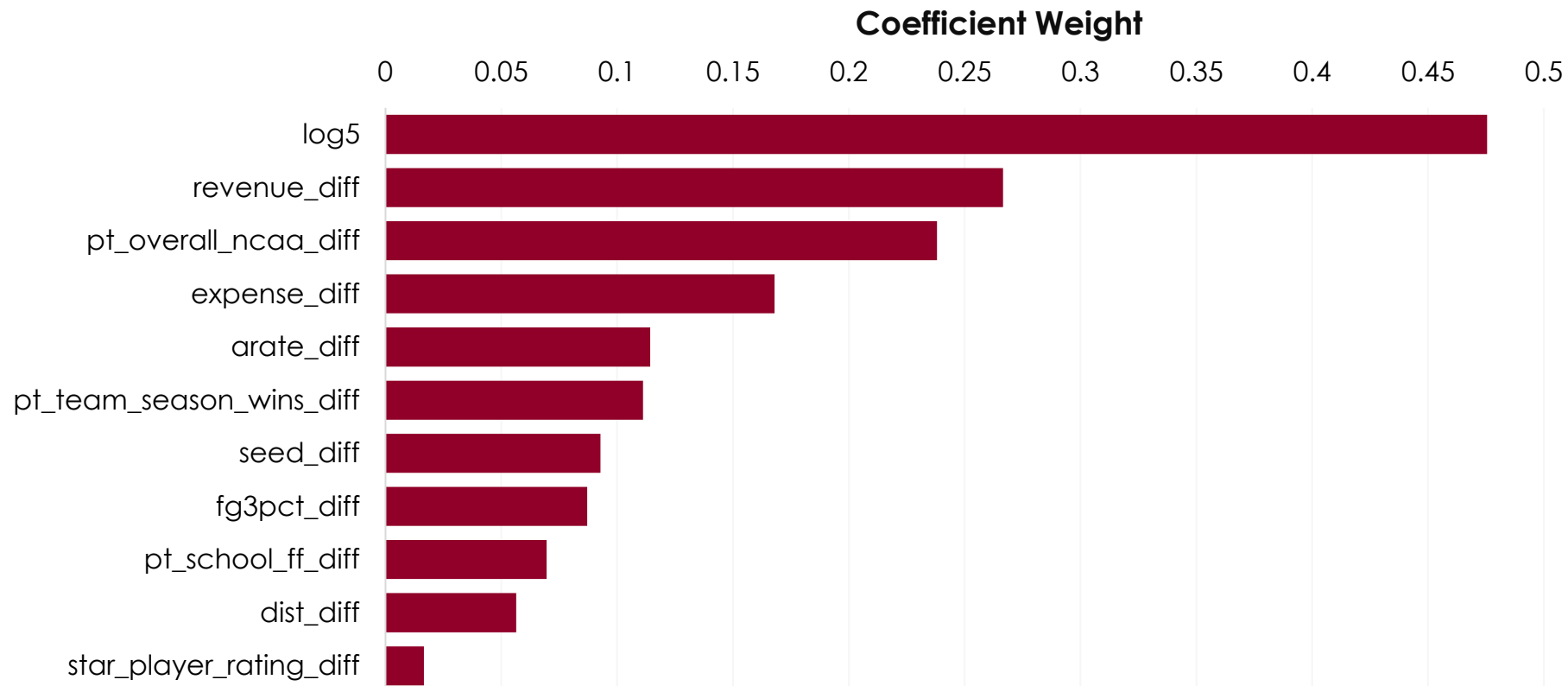
Name	Accuracy	Precision	Recall	F1 Score	Log Loss	ROC AUC
Random Forest	0.82	0.75	0.89	0.81	0.53	0.85
Logistic Regression 🏆	0.80	0.73	0.85	0.79	0.49	0.86
LightGBM	0.72	0.62	0.97	0.76	0.52	0.83
XGBoost	0.74	0.66	0.89	0.76	0.55	0.81
Decision Tree	0.75	0.71	0.75	0.73	0.56	0.79



# Model Evaluation

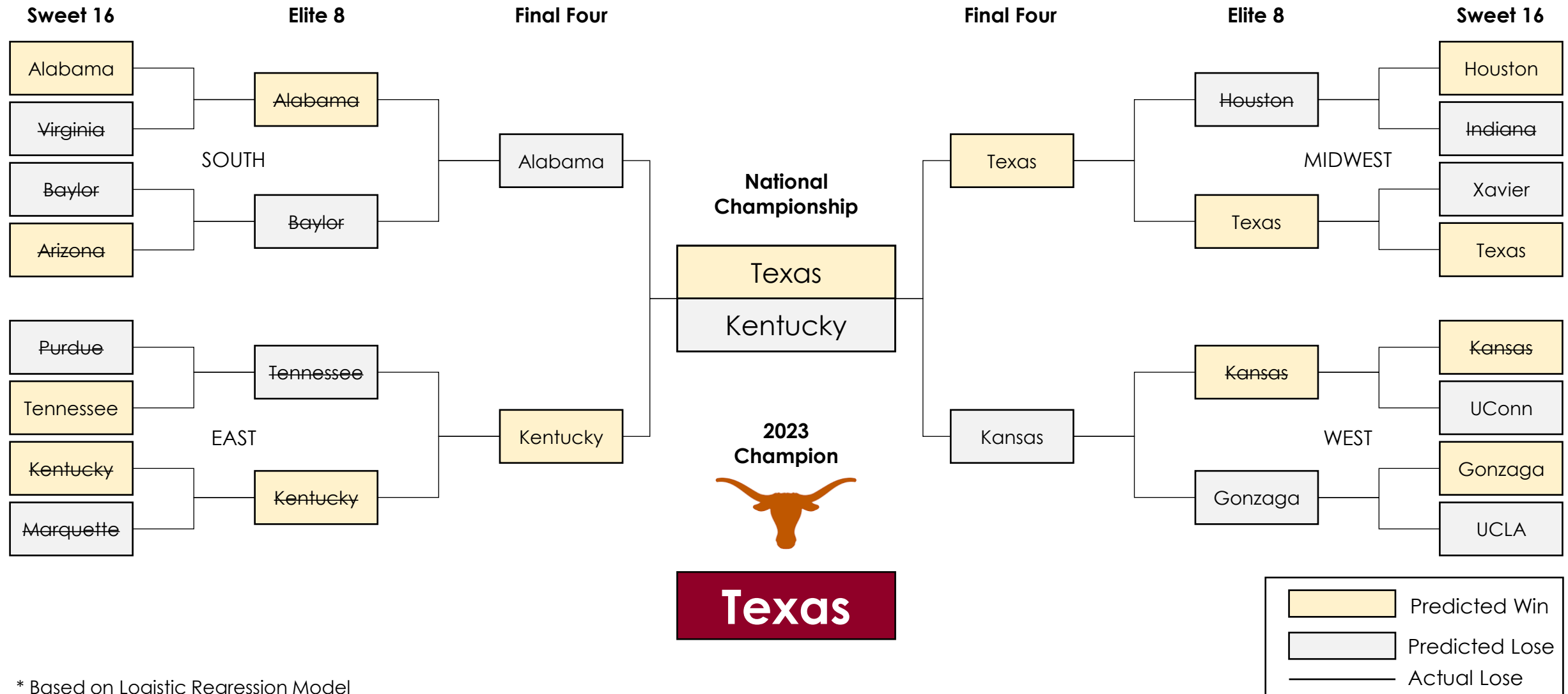
New features Revenue, Team Expense and Star Player Rating increased model accuracy

## Logistic Regression Feature Importance



# Tournament Result

Comparing predictions against Final Four outcome, correctly picked 50% of Sweet 16 and 25% of Elite 8



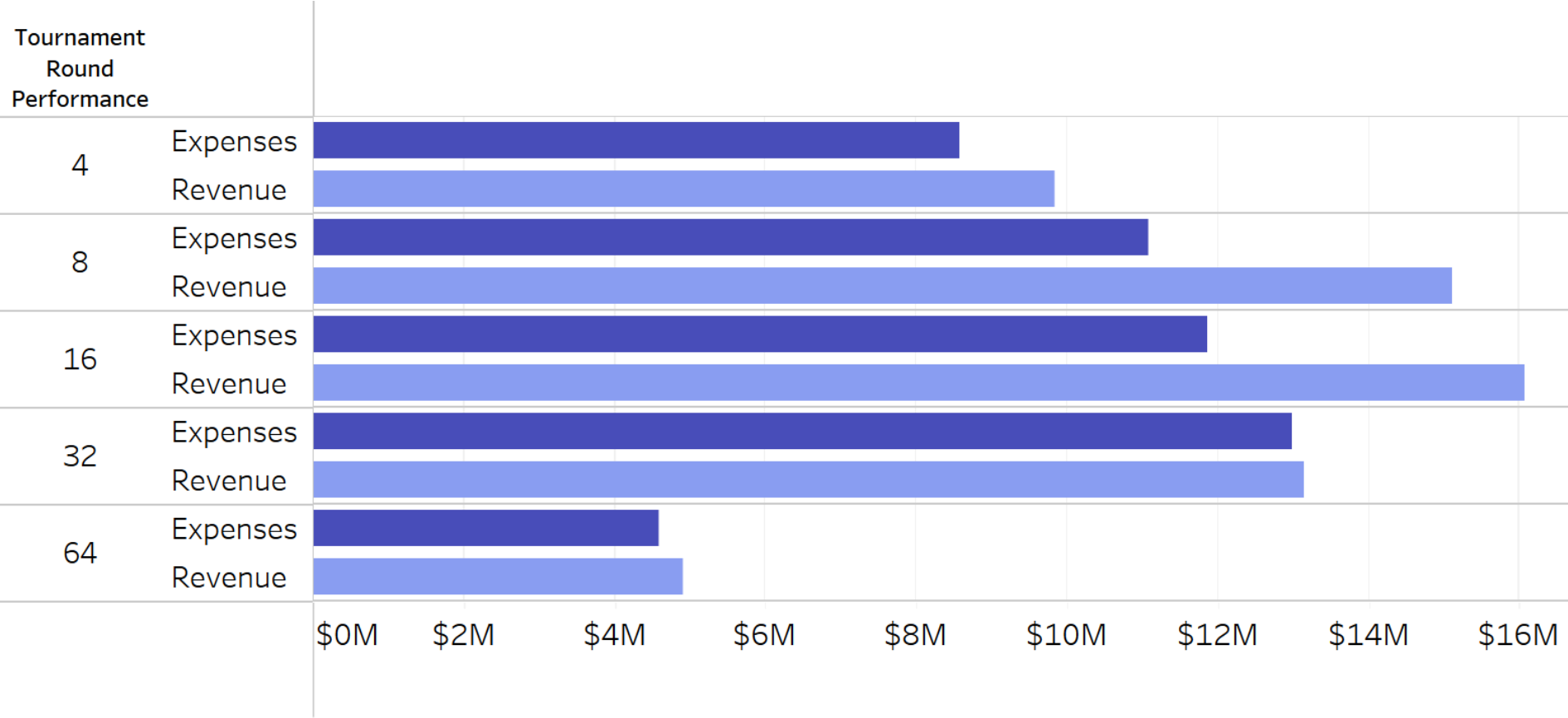
\* Based on Logistic Regression Model

# New Features Result

## Evaluating Team Expenses and Revenue by actual tournament round finish

Higher expense budget and revenue advantageous in early rounds, importance waned in later rounds

**Average team expenses and revenue by tournament round performance** (through final four)

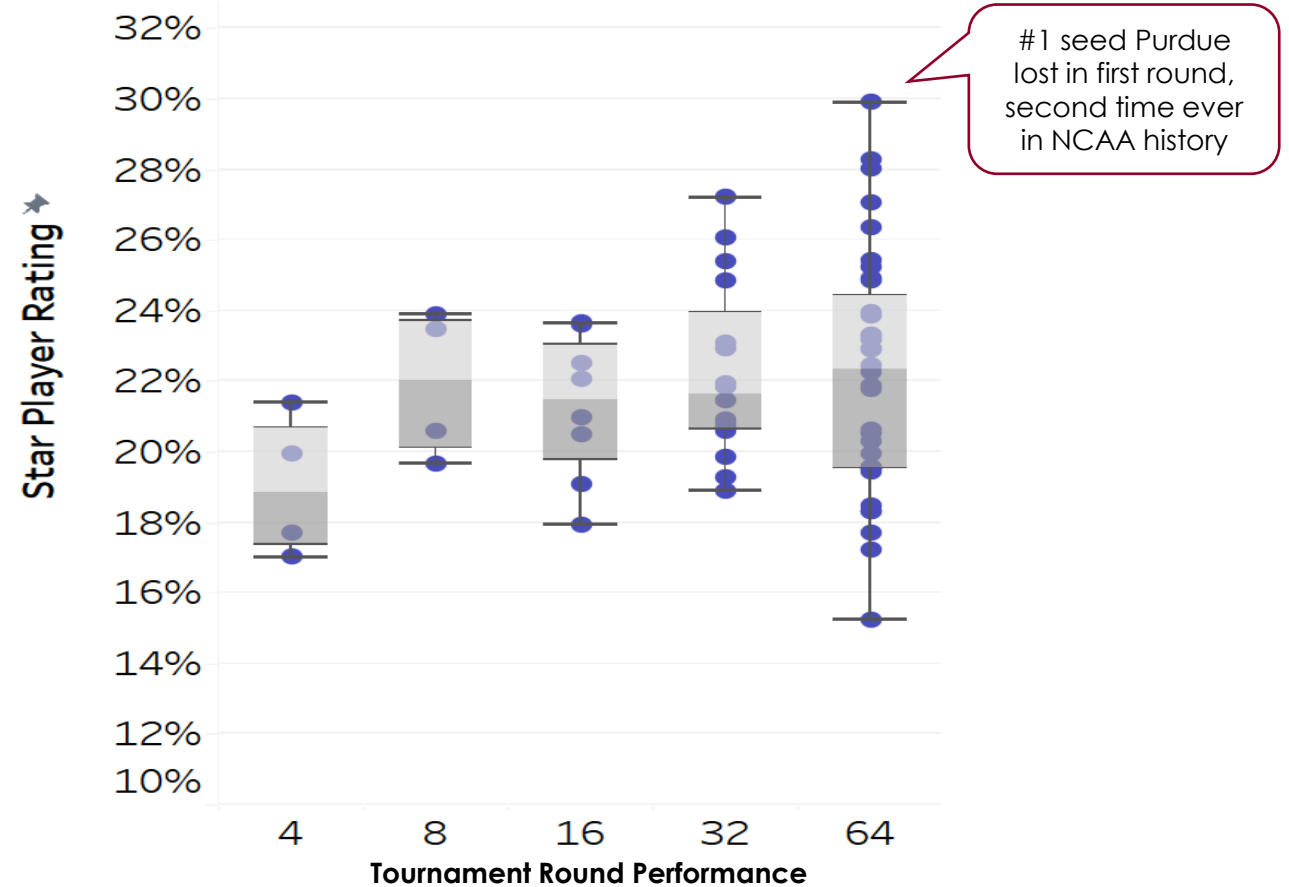


# New Features Result

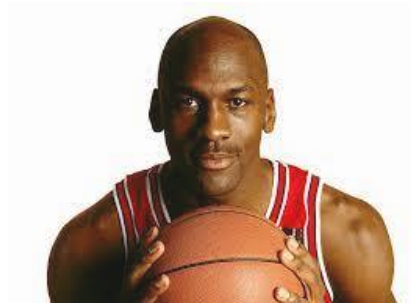
Break down of team's Star Player Ratings according to actual tournament round finish

- Bigger is not necessarily better for star player rating
- Top scorer should account for middle ground of ~20-22% total team points

**Team star player rating by tournament round performance**  
(through final four)



# Thank You!



Sponsored by: **Deloitte.**

**FORDHAM**  
THE JESUIT UNIVERSITY OF NEW YORK