# Wharton High School Data Science Competition 2025

# Using machine learning to rank teams and predict outcomes of basketball games

# The Bruzz

Ken Lieu Fraser Clyne James Epps Advisor - Mr A McIlroy

Queen Elizabeth's Grammar School Faversham



Gu

Sports Analytics and Business Initiative

#### Introduction and Background Research

- Previously, Stanford won the 2021 Women's NCAA title, beating Arizona 54-53.
- Many high level players competing
- Analysts predict NCAA basketball outcomes using statistical models, machine learning, expert polls, and tools like KenPom for team efficiency.
- Bracket simulations also help, but the unpredictable nature of the tournament makes predictions difficult.



# **Research Question - Can Machine Learning effectively be used in sports?**

- Binary classification problem (two outcomes/labels)
  - In basketball, games can only result in a win or loss due to overtime rules.
- We decided that a machine learning classification algorithm can be employed
  - These algorithms also produce, quantitative results which can be interpreted as probabilities.

#### LLMs

- Used ChatGPT to develop 4 models
  - o XGBoost
  - Logistic regression
  - o Naive Bayes
  - o Random Forest

**Definition:** Classification Algorithms try to predict the correct label of given input data.

### **Data Cleaning**

- Aggregate match data
- Non-D1 team data was excluded
- NBA level statistics
- Lead Retention Rate

 ${\rm LRR} = \frac{{\rm Final\ Score\ Differential\ (team\ score\ -\ opponent\ score)}}{{\rm Largest\ Lead}}$ 

#### **Data Transformation**

- Machine Learning models perform better
  - with strong linear
  - correlations
- Weaker correlations can
  - be transformed using
  - polynomial and log

#### functions

# **Correlation Analysis**

- Analyse and evaluate the relationship between variables.
- Spearman's rank correlation coefficient
- Select the most correlated variables to use as inputs in our models.



Results

	SOUTH		WEST		NORTH	
Rank	Team	Win Loss History	Теат	Win Loss History	Team	Win Loss History
1	Louisville Cardinals	0.8621	BYU Cougars	0.8889	South Carolina Gamecocks	0.9355
2	Iowa State Cyclones	0.8125	Stanford Cardinal	0.9032	Jackson State Lady Tigers	0.7692
3	Taledo Rockets	0.8967	Baylor Bears	0.8125	Stephen F Austin Ladyjacks	0.8571
4	Ohio State Buckeyes	0.7931	Nebraska Conthuskers	0.7500	Lsu Tigers	0.8148
5	Iowa Hawkeyes	0.7667	Texas Longhorns	0.8125	Ucf Knights	0.8889
6	IU Indianapolis Jaguars	0.8571	South Dakota State Jackrabbits	0.7097	Beimont Bruns	0.7586
7	Virginia Tech Hokies	0.7188	Gonzaga Buildogs	0.8125	Florida Gulf Coast Eagles	0.9231
	Dayton Flyers	0.8462	South Dakota Coyotes	0.8276	Tennessee Lady Volunteers	0.7419
9	Michigan Wolvennes	0.8548	UNLV Lady Rebels	0.8125	Ole Mss Rebels	0.7241
10	Notre Dame Fighting Irish	0 7333	Arizona Wildcats	0 7308	Mercer Bears	0.8125
11	Missouri State Lady Bears	0.7967	Creighton Bluejays	0.6897	Troy Trojans	0.7586
12	Indiana Hoosiers	0.7333	New Mexico Lobos	0.7333	Georgia Lady Buildogs	0.6786
13	Depaul Blue Demon	0.6875	Utah Utes	0.6452	Middle Tennessee Blue Raiders	0.7667
14	Kentucky Wildcats	0.6333	Oklahoma Sooners	0.7500	Georgia Tech Yellow Jackets	0.6774
15	Murray State Racers	0.6786	Colorado Buffaloes	0.7333	South Florida Bulls	0.7742
16	Cleveland State Vikings	0.7063	Oregon Ducks	0.6207	Arkansas Razorbacks	0.5806

#### **Results**



Probability the higher seeded team wins

Game ID

Unweighted average of probabilities
Final ranks did not reflect win loss history

Model Type	Accuracy	AUC	F1	
XGboost	0.7677	0.8499	0.8137	
Logistic Regression	0.7768	0.8572	0.8181	
Naive Bayes	0.7609	0.8414	0.7926	
Random Forest	0.7732	0.8515	0.8166	

# Can Machine Learning effectively be used in sports?

Answering the research question created at the beginning of our methodology

**Conclusion:** Machine learning method is effective with some drawbacks

Pros:

- High model accuracy
  - Mean accuracy 0.7727
  - Peak accuracy 0.8047

Cons:

- Sports is affected by other factors
  - Luck, sentiment, condition etc
- Limited dataset

With the right resources, experience and knowledge, machine learning is an effective tool for sports

Our model is trained on this league data

• Insights are specific to this league



## Discussion

#### Through SHAP values, coaches can view areas of focus such as:

- Expected Win Loss
- Average Lead
- Net Rating
- Turnovers
- Free throw scoring

#### which are high value features in our model