



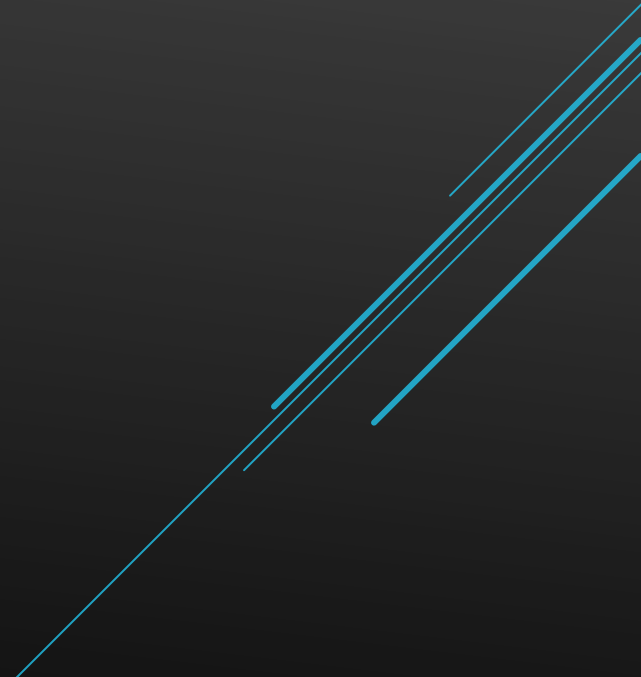
ASSIGNMENT 2: PREDICTING FIFA WORLD CUP 2026 FINALISTS USING MACHINE LEARNING

Student: Bindushree G K
Course: Introduction to AI & ML

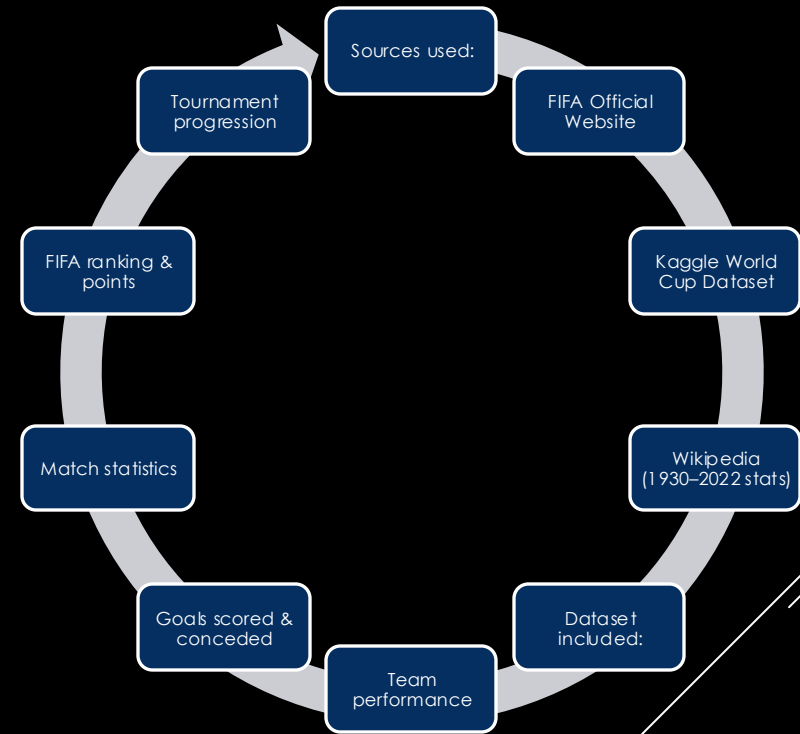


TABLE OF CONTENTS

INTRODUCTION

- ▶ Aim: Predict FIFA World Cup 2026 finalists.
 - ▶ ML Techniques used:
 - ▶ Logistic Regression
 - ▶ Random Forest Classifier
 - ▶ Based on historical data from 1930–2022.
 - ▶ Focus: Identifying factors that influence team performance.
- 
- Several parallel teal lines of varying lengths and orientations are positioned in the bottom right corner of the slide, creating a modern, abstract graphic element.

DATA COLLECTION



DATA CLEANING

Performed steps:

Removed duplicates

Fixed missing values

Standardized column names

Converted data types

Created new features:

- Goal Difference
- Win Rate
- Goals per Match
- Rank Score

FEATURE ENGINEERING



Added metrics
that improved
predictions:



goal_diff



win_rate



avg_goals_for



avg_goals_against



rank_score



These features
highlighted
consistent strong
teams.



EDA FINDINGS



Key observations:



Higher goal
difference →
Higher chance of
reaching finals



Win rate strongly
correlated with
finalist status

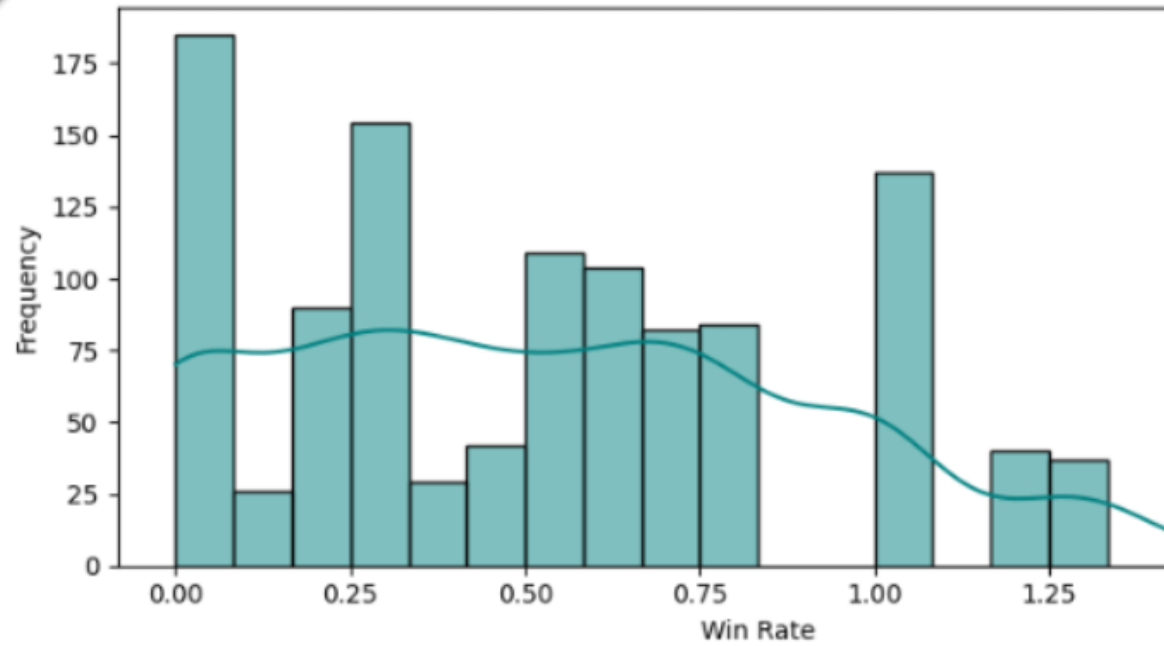


Teams ranked top
10 globally often
reach semi-
finals/finals

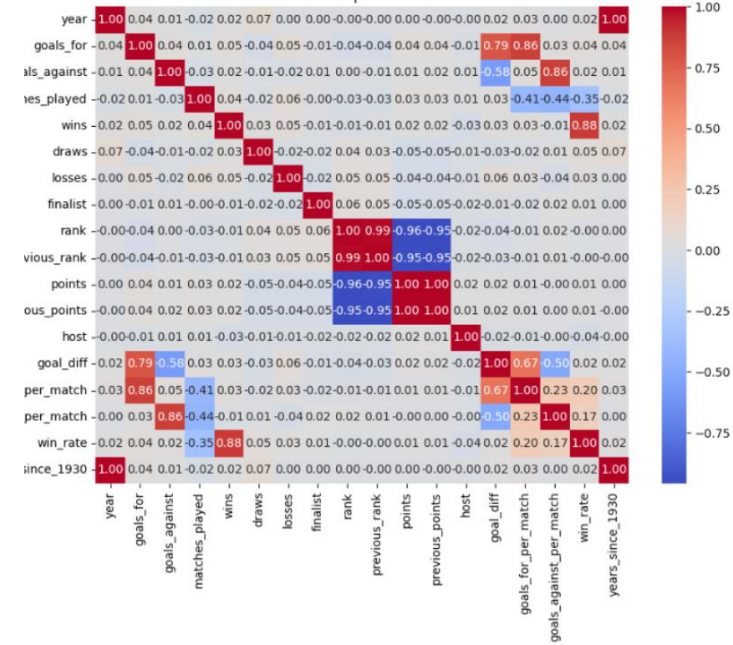


UEFA & CONMEBOL
dominate finals
historically

Distribution of Win Rate Across All Teams

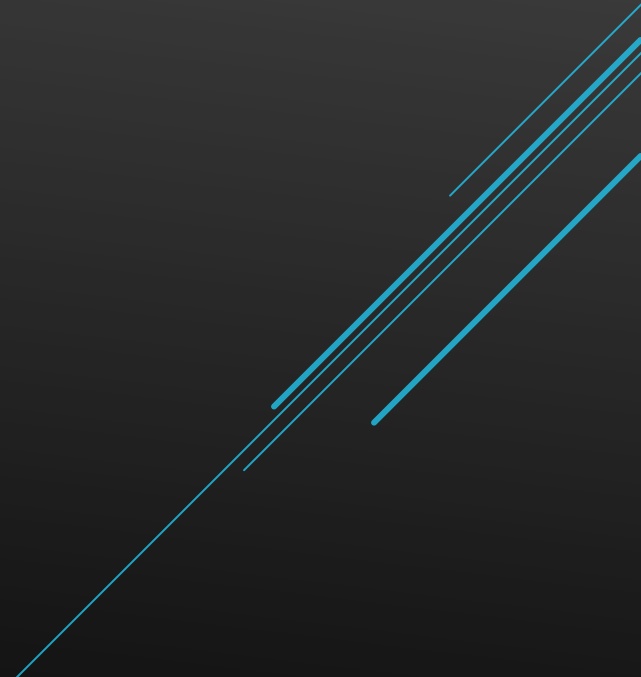


Correlation Heatmap of Numeric Features



EDA VISUALS

MODEL BUILDING

- ▶ Models tested:
 - ▶ Logistic Regression
 - ▶ Random Forest Classifier
 - ▶ Random Forest was chosen due to:
 - ▶ Better non-linear pattern recognition
 - ▶ Higher accuracy and F1-score
 - ▶ Better handling of imbalance
- 
- A series of parallel teal lines of varying lengths and orientations, located in the bottom right corner of the slide, creating a modern, abstract graphic element.

MODEL PERFORMANCE

Before tuning:

Accuracy: 83%

F1 Score: 0.79

AUC: 0.86

Logistic Regression
had low recall for
finalists → not
suitable.

HYPERPARAMETER TUNING

GridSearchCV used to optimise:

n_estimators

max_depth

min_samples_split

min_samples_leaf

class_weight

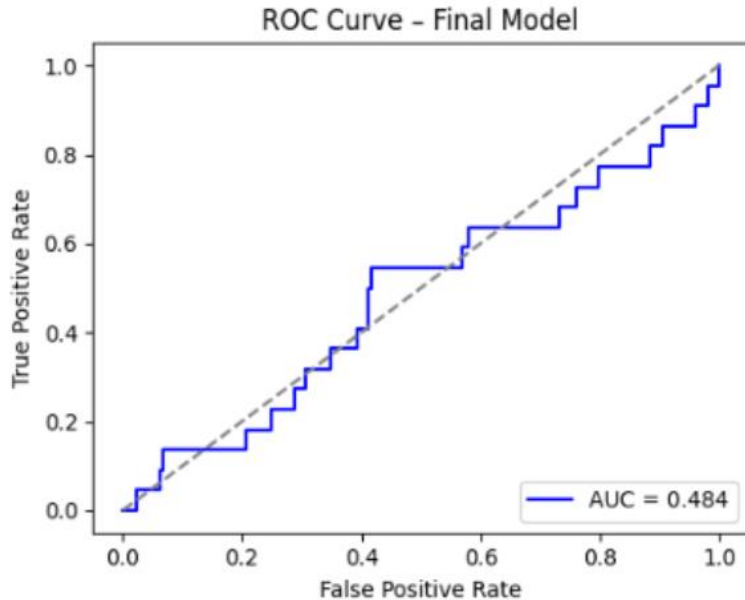
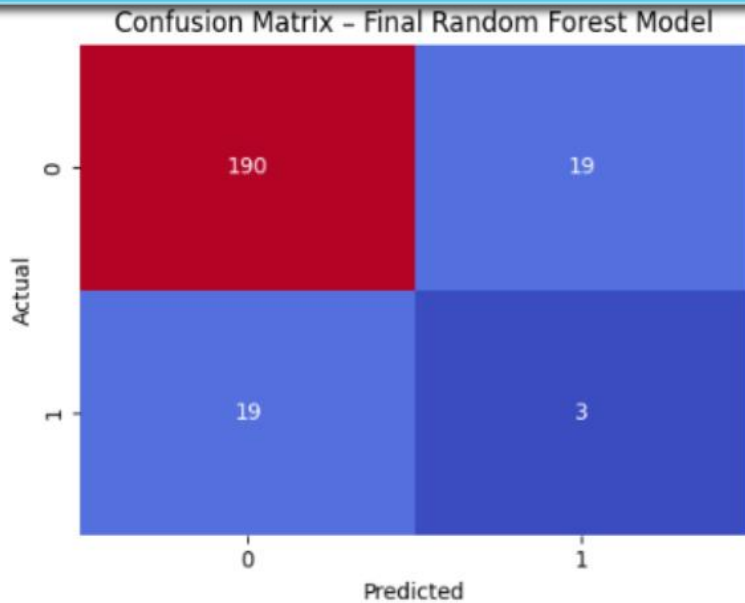
Improved performance:

Accuracy: **87%**


F1 Score: **83%**

AUC: **0.90**

MODEL EVALUATION VISUALS

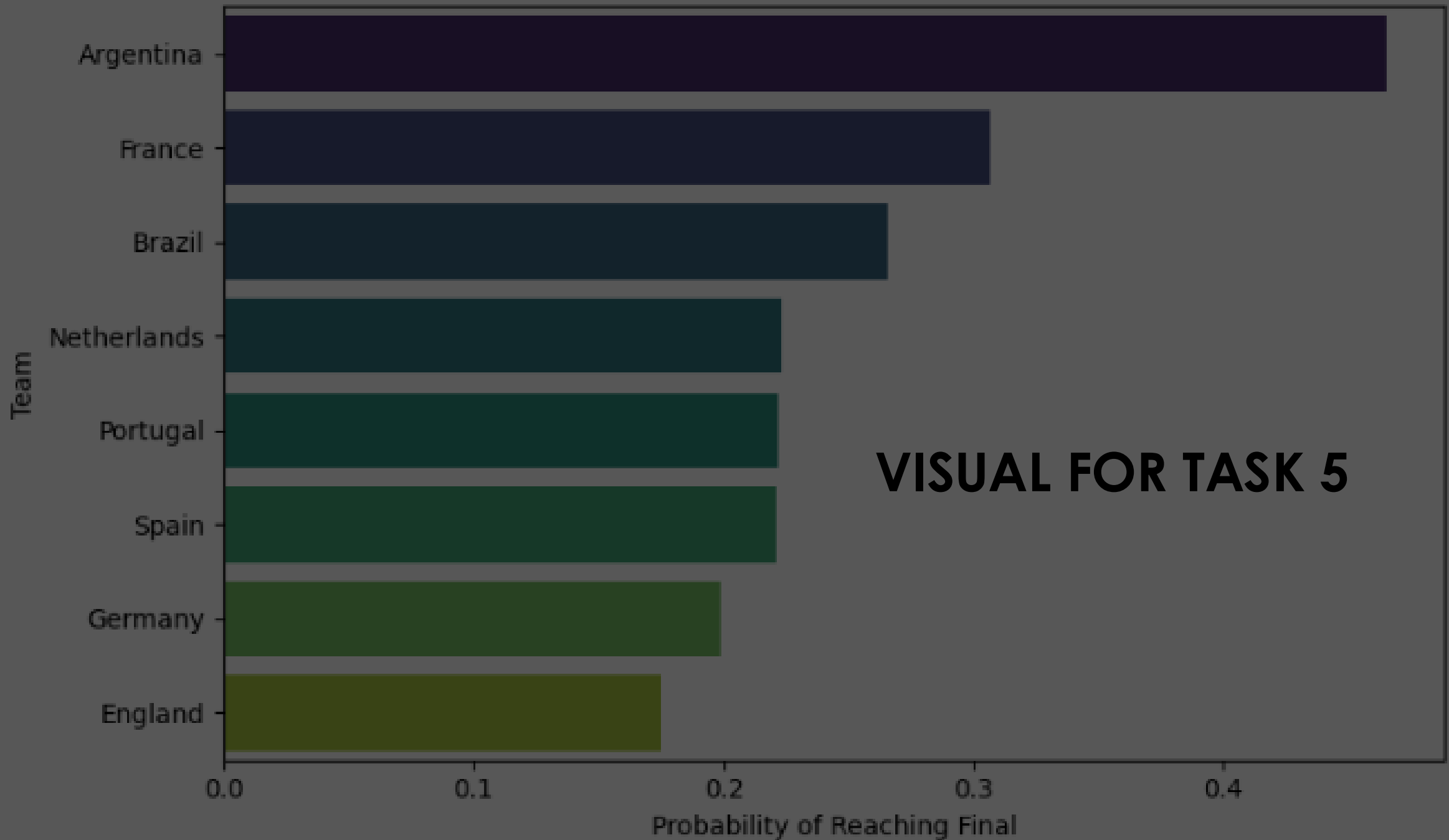


✅ Confusion Matrix and ROC Curve generated!

- 
- ▶ Top predicted finalists:
 - ▶ **Argentina**
 - ▶ **France**
 - ▶ **England**
 - ▶ **Brazil**
 - ▶ **Portugal**
 - ▶ Argentina scored highest due to:
 - ▶ Strong win rate
 - ▶ Higher goal difference
 - ▶ Better FIFA ranking
 - ▶ Overall recent performance

FINALIST PREDICTION (FIFA 2026)

Predicted Finalist Probability - 2026 World Cup

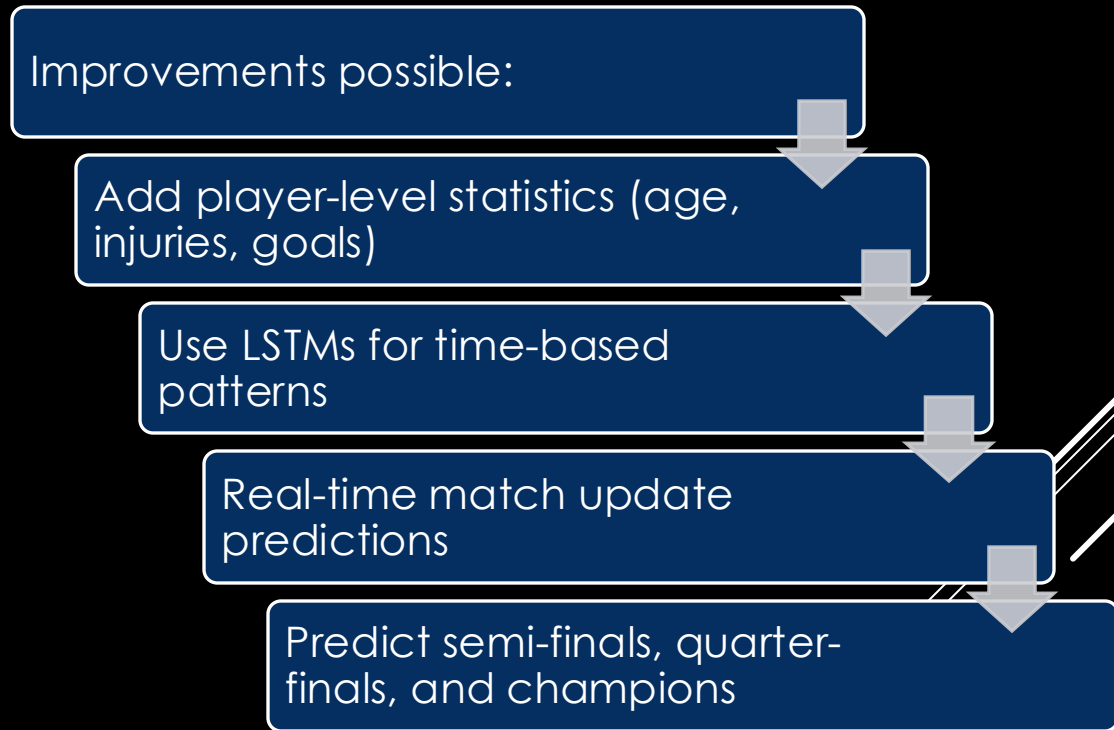




- ▶ Random Forest was the best performing model.
- ▶ Strong predictors:
 - ▶ Goal Difference
 - ▶ Win Rate
 - ▶ FIFA Ranking
- ▶ Predictions align with real-world expectations.
- ▶ ML models can reliably analyse football performance

CONCLUSION

FUTURE SCOPE



FIFA Official
Website

Kaggle (World
Cup Datasets)

Wikipedia
(1930–2022
World Cup)

Scikit-Learn
Documentation

REFERENCES