# NATIONAL UNIVERSITY OF COMPUTER & EMERGING SCIENCES

## KARACHI CAMPUS



### CL1002 - Machine Learning - LAB

#### PROJECT : "CRICKET PREDICTION MODEL"

#### Group Members

| Name | Student ID | Section |
|------|-----------|---------|
| Mudasir | 22K - 8732 | BAI - 5A |
| Nihal Ali | 22K - 4054 | BAI - 5A |
| Irteza | 22K - 8731 | BAI - 5A |

CLASS/SECTION: BAI-5A

DATE: 14/DEC/2024

INSTRUCTOR: SIR USAMA BIN UMAR

# T20I CRICKET PREDICTION MODEL

This project aims to predict the winner of T20 International cricket matches based on historical data and statistical analysis. The model uses a Random Forest Classifier, which outperformed other algorithms, achieving an accuracy of 0.95 or 95%. Various machine learning models were evaluated, and an interactive GUI was developed to enhance user engagement.

## DATASET OVERVIEW

The dataset is structured on a ball-by-ball basis, later aggregated by match date to eliminate duplicate rows. Key attributes include:

- Match ID, Date, Venue
- Batting and bowling statistics (runs, wickets, overs, etc.)
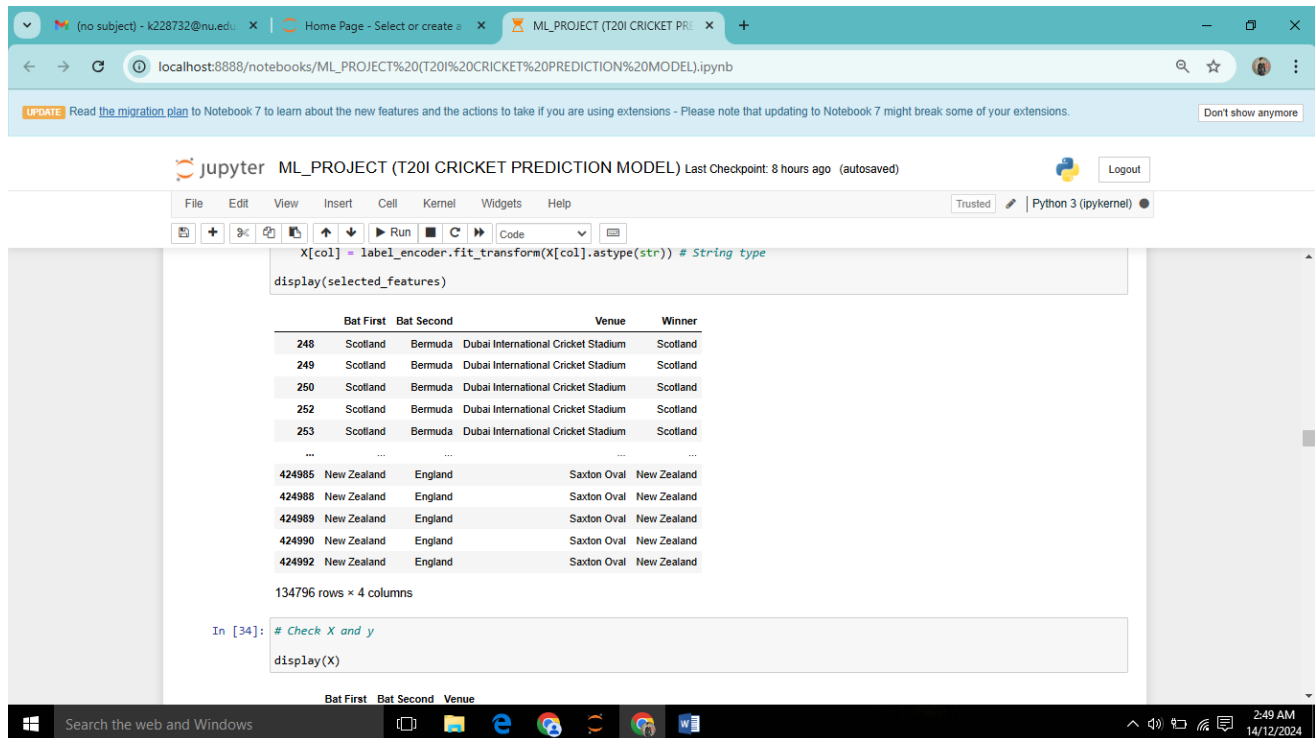- Match outcomes like Winner and Toss Decision (Bat First/Second).

Selected Features for the Model:

- Bat First (toss winner)
- Bat Second
- Venue
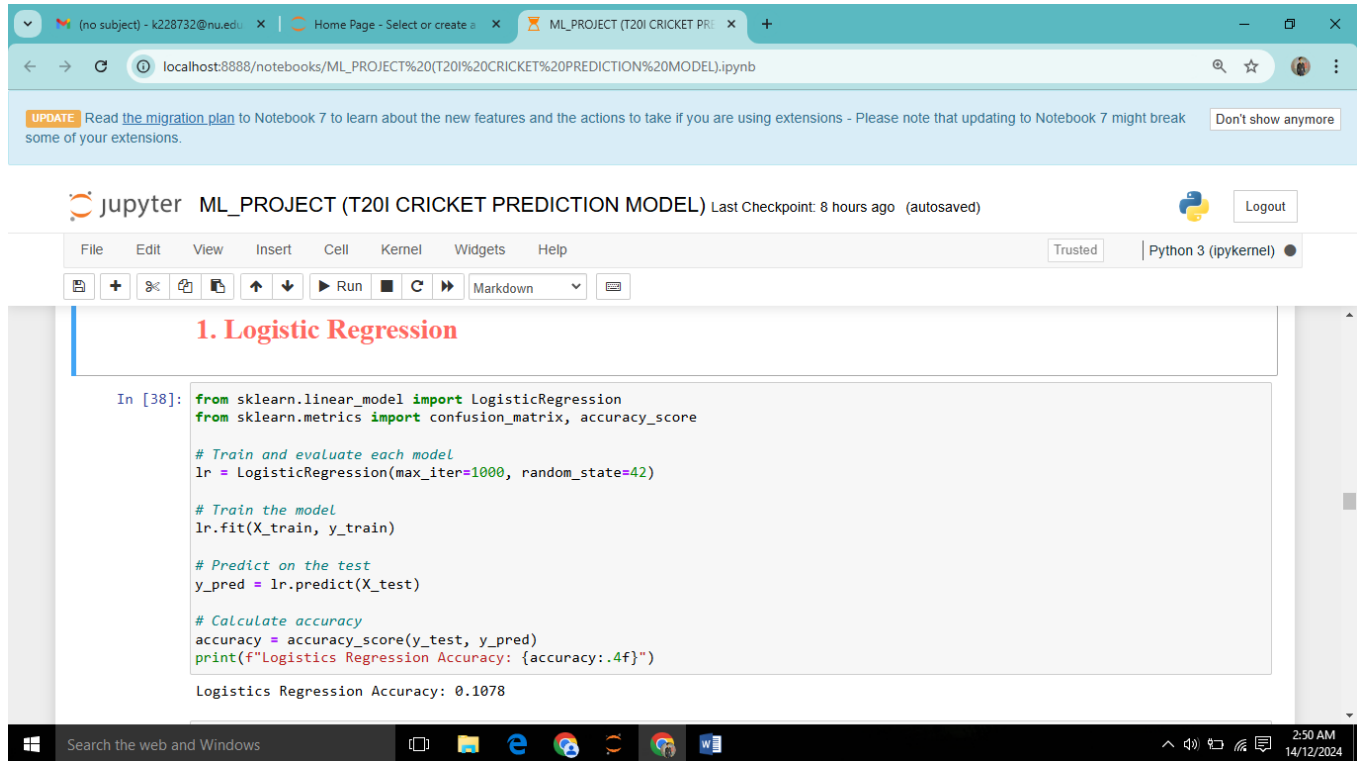- Winner (label/target column)



The dataset considers the toss-winning team as the team batting first.

# MACHINE LEARNING MODELS EVALUATED

## 1. LOGISTIC REGRESSION

- Testing Accuracy: 0.1078 or 10.78%

- Reason for Low Accuracy

Logistic Regression is primarily suited for binary classification. Since the target variable includes multiple teams as possible winners, the algorithm fails to generalize.



# 2. SUPPORT VECTOR CLASSIFIER (SVC)

➢ Testing Accuracy: 0.11 or 11%
➢ Reason for Low Performance
SVC struggles with multi-class classification when the dataset is large and complex. The need for significant parameter tuning also impacted its performance.

### 3. Applying SVC (for Classification)

```python
import numpy as np
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

model = SVC(kernel='rbf', C=1.0, gamma='scale', random_state=42)

# Train the model
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy for SVR Model: {accuracy:.2f}")
```
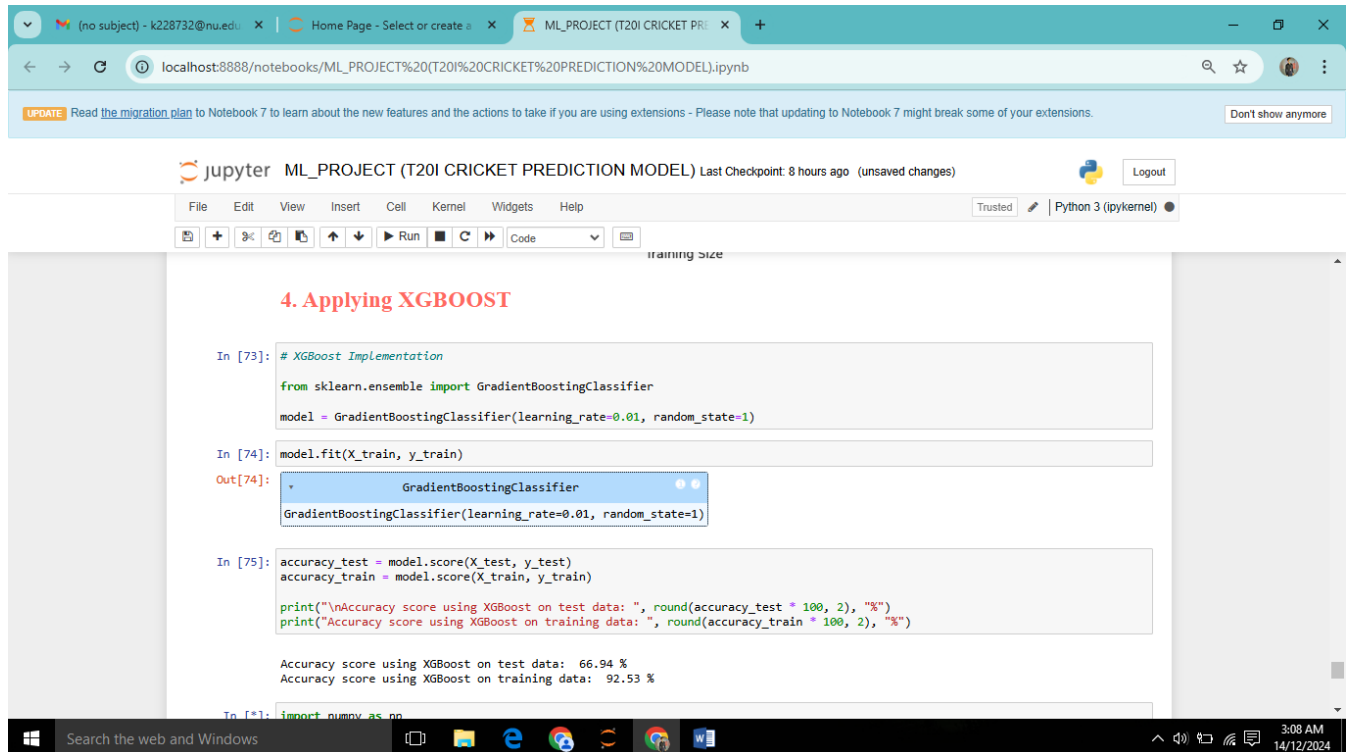
```
Accuracy for SVR Model: 0.11
```

```python
# Classification Report

report = classification_report(y_test, y_pred, zero_division=0)
print(report)
```

```
                precision    recall  f1-score   support
```



# 3. XGBOOST CLASSIFIER

➢ Testing Accuracy: 0.6694 or 66.94%
➢ Reason for Partial Success:
XGBoost handles classification tasks better but requires fine-tuned hyperparameters and specific handling of categorical features for optimal performance.

## 4. ADABOOST CLASSIFIER

- ➤ Testing Accuracy: 0.0921 or 9.21%
- ➤ Reason for Lower Performance:
  AdaBoost relies heavily on weak learners and performs well for simpler datasets. The complexity and scale of the dataset limited its effectiveness.

## 5. Applying AdaBoost

```python
# AdaBoost

from sklearn.ensemble import AdaBoostClassifier

model = AdaBoostClassifier(n_estimators=n_estimators, algorithm='SAMME', random_state=42)
model.fit(X_train, y_train)
```
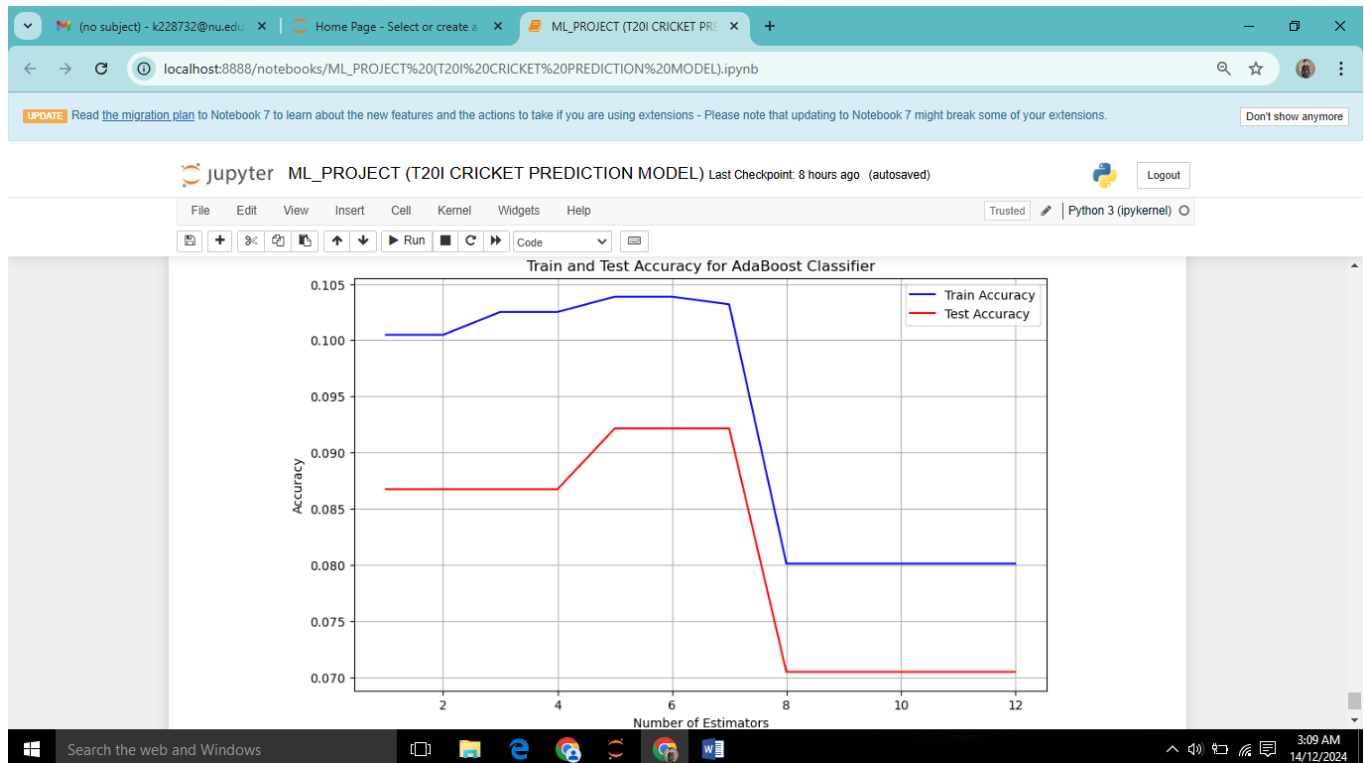
Out[78]:
AdaBoostClassifier
AdaBoostClassifier(algorithm='SAMME', n_estimators=9, random_state=42)

```python
accuracy_test = model.score(X_test, y_test)
accuracy_train = model.score(X_train, y_train)

print("\nAccuracy score using AdaBoost on test data: ", round(accuracy_test * 100, 2), "%")
print("Accuracy score using AdaBoost on training data: ", round(accuracy_train * 100, 2), "%")
```

```
Accuracy score using AdaBoost on test data:  9.21 %
Accuracy score using AdaBoost on training data:  10.39 %
```

```python
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

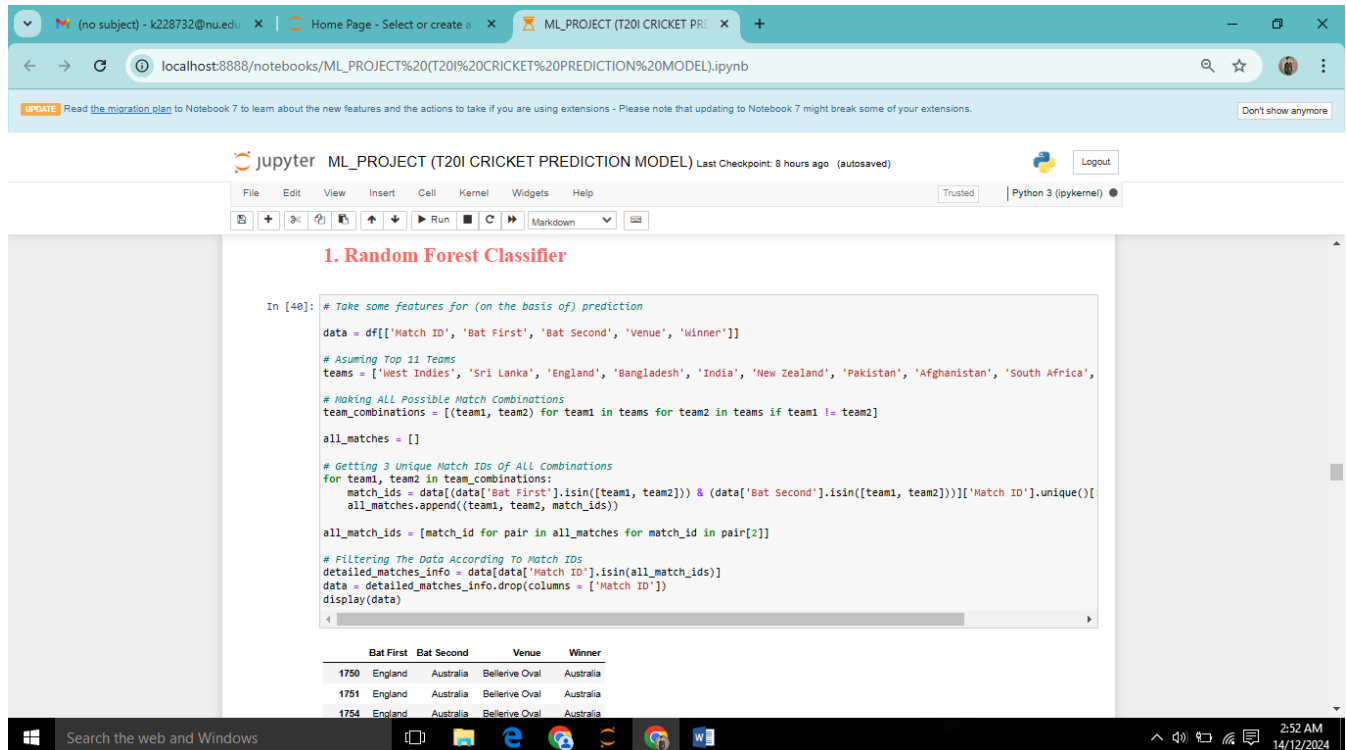

Train and Test Accuracy for AdaBoost Classifier

# 5. RANDOM FOREST CLASSIFIER

➤ Testing Accuracy: 0.95 or 95%

➢ Why It Succeeded

Random Forest effectively handles categorical variables, captures non-linear relationships, and prevents overfitting with its ensemble approach.

# CREATING GUI FOR PREDICTION AND STATISTICS OF MATCHES

## 1. MATCH OUTCOME PREDICTION

- Functionality: Users can select:
  - Team 1
  - Team 2
  - Venue
  - Toss Prediction (the team bat first)


Result: The model predicts the winning team based on past match winning statistics on this specific venue.

## 2. CHECKING OUR PREDICTION BY PAST MATCHES STATISTICS OF THIS SPECIFIC SELECTED FIELDS (GUI)

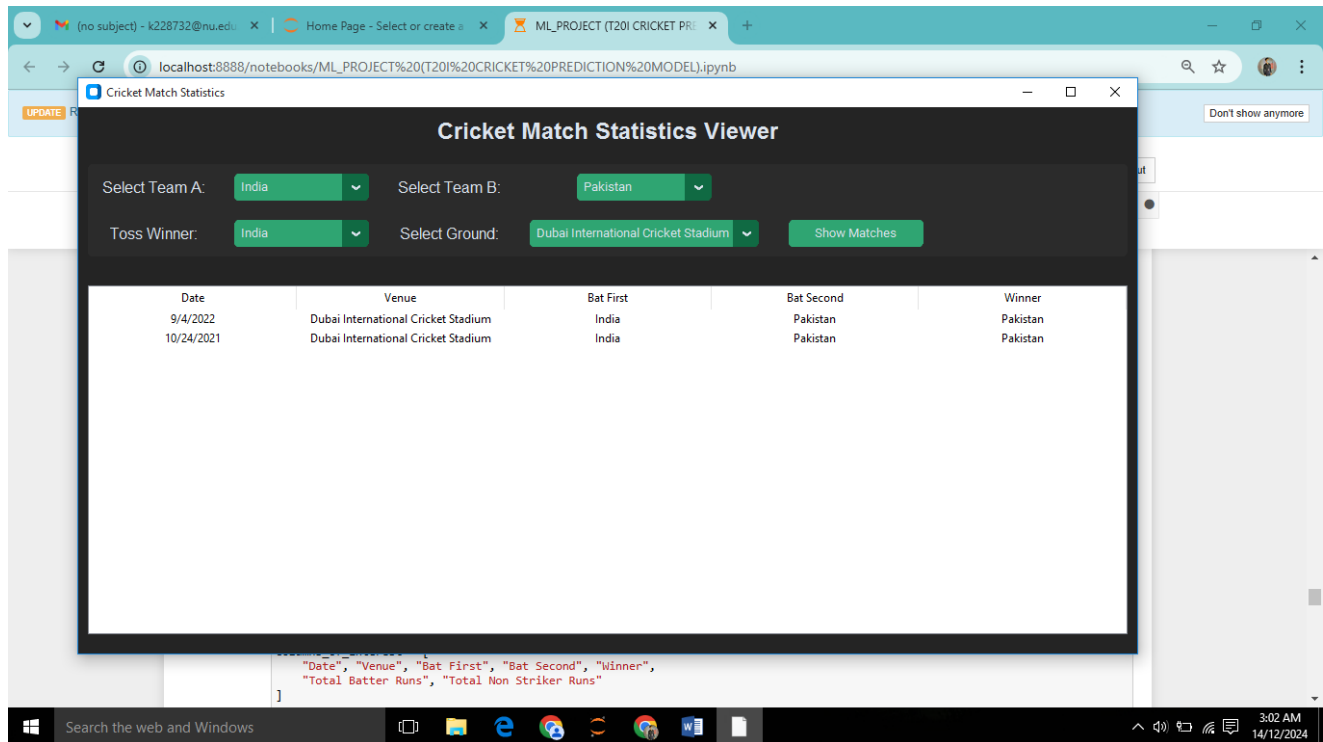Feature: Displays all historical match data for the selected teams.

Purpose: Users can verify the accuracy of model predictions by comparing them with actual outcomes.
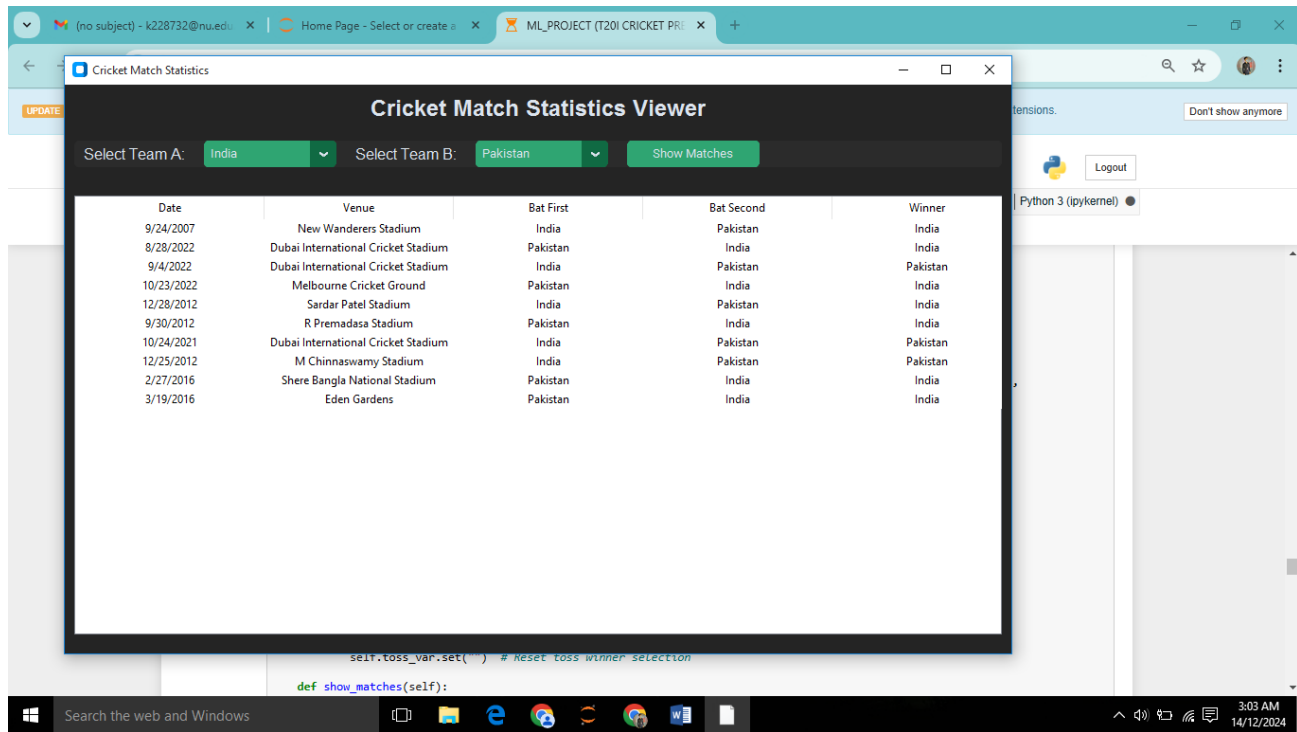
## Testing Prediction:

✓ 2 matches played in the past between India and Pakistan at Dubai International Stadium where India was bat first in both games and Pakistan were won the both matches. Now when I am predicting for future match with same venue and other statistics. Our model is predicting Winner as "Pakistan".

# 3. GROUND SPECIFIC STATISTICS (GUI)

Feature: Users select Team 1, Team 2, and a specific venue.

Output: Displays past match statistics specific to the selected ground to help users understand performance trends.

# MODEL DEVELOPMENT AND ANALYSIS

## 1. DATASET PREPROCESSING (PERFORMED EDA)

- Aggregated the dataset by date to ensure each match is represented as a single entry.
- Eliminated duplicate rows and retained only relevant features.

## 2. FEATURE ENGINEERING AND SCALING

- Encoded categorical variables such as teams and venues using one-hot encoding and label encoding.p
- Implement Standard Scaler to make sure the data is in small range.

## 3. SPLITTING IN TRAINING AND TESTING DATA:

- Split the dataset into training (80%) and testing (20%) sets.

# RESULTS AND CONCLUSION

- The Random Forest Classifier emerged as the best model, achieving a 95% testing accuracy.
- GUI enhancements provide a user-friendly interface for exploring predictions and historical data.