

INT247: MACHINE LEARNING FOUNDATION



LOVELY PROFESSIONAL UNIVERSITY

PHAGWARA, PUNJAB

Submitted in partial fulfillment of the requirements for the award of degree of

BTech Computer Science Engineering

PROJECT REPORT

Cricket prediction using ML

SUBMITTED BY

Name of student : Hemanth syam Siddardha

Registration Number : 11904722

Section & Roll number : RKM055A75

**Submitted To:
Dr. Sagar Pande**

Student Declaration

To whom so ever it may concern

I, Hemanth syam, **11904722**, hereby declare that the project on “**Cricket prediction using ML**” is based on my own work and carried by my own

Fact and observation during the course of my study under the guidance of

Dr. Sagar Pande sir

-

Name of the Student :- Hemanth syam

Registration Number :- 11904722

Roll Number :- RKM055A75

ACKNOWLEDGEMENT

Primarily I would like to thank you my college and my teacher **Dr. Sagar Pande sir** for guiding continuously throughout the project .Then I would like to express my special thanks for college who provided such an opportunity for students to gain hands on experience on machine learning by learning and building projects which helps to get career ready.

I would like to again thank my own college Lovely Professional University for offering such a opportunity which not only improve my programming skill but also taught me other new technology.

Then I would like to thank my parents and friends who have helped me with their valuable suggestions and guidance in completion of this project

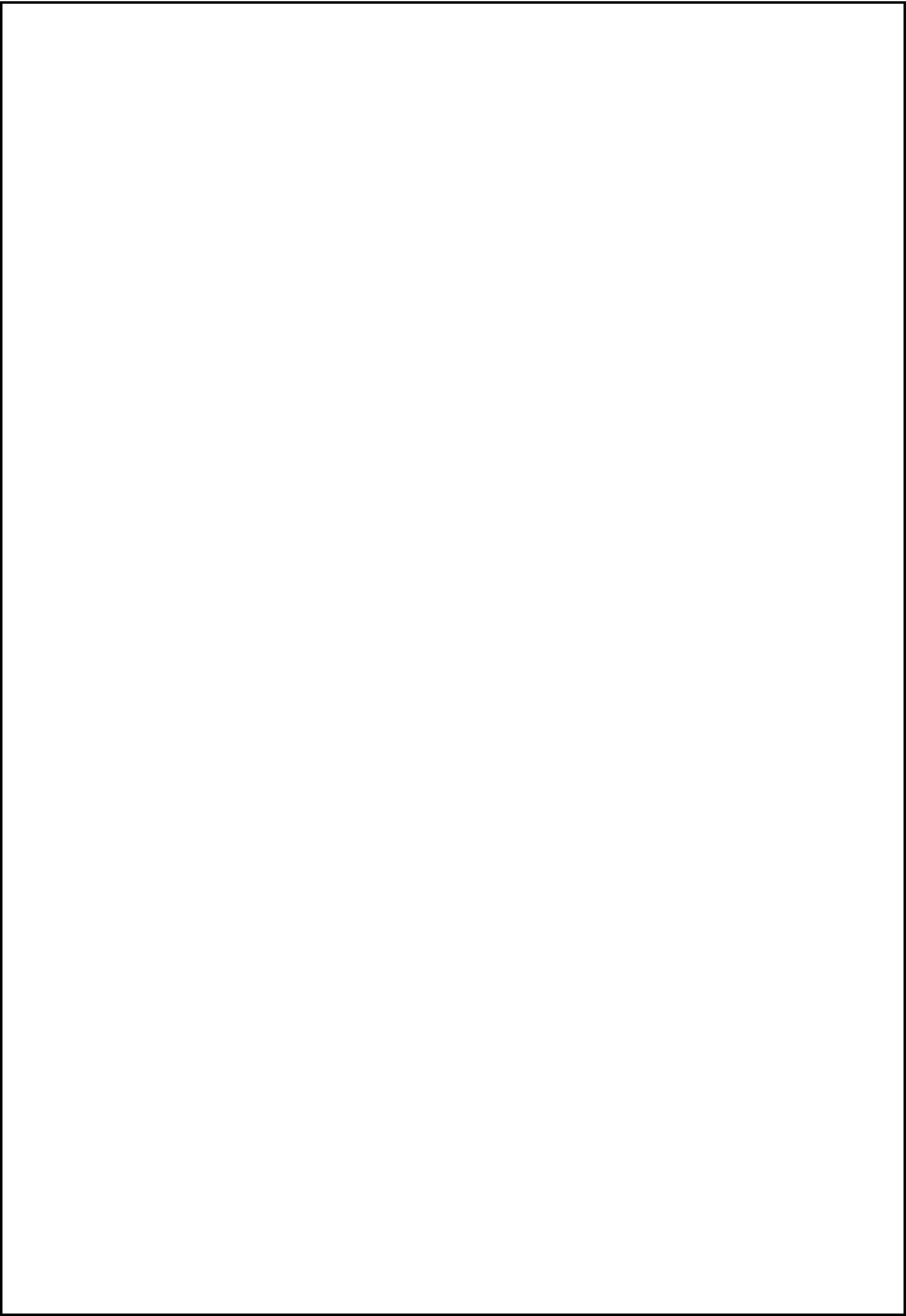
Date: 27/09/2021

Hemanth syam

Reg no: 11904722

Table of Contents

S. No.	Title	Page
1	Cover Page	1
2	Declaration of the student	2
3	Acknowledgement	3
4	Index	4
5	Abstract	5
6	Introduction	6-10
7	Literature Review	11-
8	Software Tools	8
9	Modules/Libraries used	18
10	Learning Outcome	19
11	Bibliography	22



ABSTRACT

In today's date data analysis is need for every data analytics to examine the sets of data to extract the useful information from it and to draw conclusion according to the information. Data analytics techniques and algorithms are more used by the commercial industries which enables them to take precise business decisions. It is also used by the analysts and the experts to authenticate or negate experimental layouts, assumptions and conclusions.

In recent years the analytics is being used in the field of sports to predict and draw various insights. Due to the involvement of money, team spirit, city loyalty and a massive fan following, the outcome of matches is very important for all stake holders. In this paper, the past seven year's data of worldcup containing the player's details, match venue details, teams, ball to ball details, is taken and analyzed to draw various conclusions which help in the improvement of a player's performance.

Various other features like how the venue or toss decision has influenced the winning of the match in last seven years are also predicted. Various machine learning and data extraction models are considered for prediction are Linear regression, Decision tree, K-means, Logistic Regression etc. The cross validation score and the accuracy are also calculated using various machine learning algorithms. Before prediction we have to explore and visualize the data because data exploration and visualization is an important stage of predictive modeling.

INTRODUCTION

Machine Learning is a branch of Artificial Intelligence that aims at solving real-life engineering problems. This technique requires no programming, whereas it depends on only data learning where the machine learns from pre-existing data and predicts the result accordingly. Machine Learning methods have benefit of using decision trees, heuristic learning, knowledge acquisition, and mathematical models. It thus provides controllability, observability, stability and effectiveness. Cricket is being played in many countries around the world. There are a lot of domestic and international cricket tournaments being held in many countries. The cricket game has various forms such as Test Matches, Twenty20 Internationals, Internationals one day, etc. worldcup is also one of them, and has great popularity among them. It's ODICricket game league played to inspire young and talented players in India. The league was conducted once in every 4 years, April or May and has a huge fan base among India. There are eight teams which represent eight cities which are chosen from an auction. These teams compete against each other for the trophy. The whole match depends on the luck for the team, player's performance and lot more parameters that will be taken in to the consideration. The match that is played before the day is also will make a change in the prediction. The stakeholders are much more benefited due to the huge popularity and the huge presence of people at the venue. The accuracy of a data depends on the size of the data we take for analysing and the records that are taken for predicting the outcome. Cricket is a game played between two teams comprising of 11 players in each team. The result is either a win, loss or a tie. However, sometimes due to bad weather conditions the game is also washed out as Cricket is a game which cannot be played in rain. Moreover, this game is also extremely unpredictable because at every stage of the game the momentum shifts to one of the teams between the two. A lot of times the result gets decided on the last ball of the match where the game gets really close. Considering all these unpredictable scenarios of this unpredictable game, there is a huge interest among the spectators to do some prediction either at the start of the game or during the game. Many spectators also play betting games to win money

LITERATURE REVIEW

The Author of [1] discusses the analysis role of machine learning in the improvement of performances of players and the team in different sports and how wearable technology helps the players to know their performance levels and further improvements. The paper describes a system that collects raw data for each sport, team and player, and it is processed into statistical data. These data sets are clustered and stored as the data to be stored is very large. The paper also discusses wearable sensors. These wearable sensors are used in recognizing real-time tasks in sports. The devices are also helping the coaches in the transformation of Decision making. The paper concludes Machine Learning, along with Wearable devices can make a great impact on the players by making patterns, strategies, planning, reducing the risk of injury, and improving their performances. In [2] CricAI Tool is discussed. It can help adjust certain factors to maximize the chances of winning the real game. The paper addressed the problem of predicting the chances of victory in a One Day International cricket match. The paper also shows a comparative evaluation of the classifiers. The systems described have been able to predict the winning criteria formulated using attributes from the dataset. The

CricAI tool can be used in real-world applications by teams playing cricket. It can help adjust certain factors to maximize the chances of winning the real game [2]. Weighted Association Rule Mining algorithm for analyzing the Indian cricket team in one-day international cricket matches against Sri Lanka and South Africa is performed and mentioned in [3]. This analysis is used by the team for framing game-winning plans. In [3] the author shows that a bigger data set can improve the accuracy of the prediction.

SOFTWARE TOOLS

1. JUPYTER NOTEBOOK:-

Project Jupyter is a project and community whose goal is to "develop open-source software, open-standards, and services for interactive computing across dozens of programming languages".^[2] It was spun off from IPython in 2014 by Fernando Pérez^[3] and Brian Granger.^[4] Project Jupyter's name is a reference to the three core programming languages supported by Jupyter, which are Julia, Python and R, and also a homage to Galileo's notebooks recording the discovery of the moons of Jupiter. Project Jupyter has developed and supported the interactive computing products Jupyter Notebook, JupyterHub, and JupyterLab. Jupyter is financially sponsored by NumFOCUS

2. VS CODE:

Visual Studio Code is a source-code editor that can be used with a variety of programming languages, including Java, JavaScript, Go, Node.js, Python and C++. It is based on the Electron framework,^[19] which is used to develop Node.js Web applications that run on the Blink layout engine. Visual Studio Code employs the same editor component (codenamed "Monaco") used in Azure DevOps (formerly called Visual Studio Online and Visual Studio Team Services).

3. GITHUB:-

GitHub, Inc. is a provider of Internet hosting for software development and version control using Git. It offers the distributed version control and source code management (SCM) functionality of Git, plus its own features. It provides access control and several collaboration features such as bug tracking, feature requests, task management, continuous integration and wikis for every project. Headquartered in California, it has been a subsidiary of Microsoft since 2018.

It is commonly used to host open-source projects. As of November 2021, GitHub reports having over 73 million developers and more than 200 million repositories (including at least 28 million public repositories). It is the largest source code host as of November 2021.

4. HEROKU :-

Heroku is a cloud platform as a service (PaaS) supporting several programming languages. One of the first cloud platforms, Heroku has been in development since June 2007, when it supported only the Ruby programming language, but now supports Java, Node.js, Scala, Clojure, Python, PHP, and Go. For this reason, Heroku is said to be a polyglot platform as it has features for a developer to build, run and scale applications in a similar manner across most languages. Heroku was acquired by Salesforce in 2010 for \$212 million.

5. SCIKIT-LEARN :-

Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistent interface in Python. Scikit-learn provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python. Extensions or modules for SciPy are conventionally named SciKits. As such, the module provides learning algorithms and is named scikit-learn.

The library is built upon the SciPy (Scientific Python) that must be installed before you can use scikit-learn. This stack that includes:

- **NumPy**: Base n-dimensional array package
- **SciPy**: Fundamental library for scientific computing
- **Matplotlib**: Comprehensive 2D/3D plotting
- **IPython**: Enhanced interactive console
- **Sympy**: Symbolic mathematics
- **Pandas**: Data structures and analysis

MODULES / LIBRARIES USED

1. FLASK:-

Flask is a **micro web framework written in Python**. It is classified as a microframework because it does not require particular tools or libraries. It has no database abstraction layer, form validation, or any other components where preexisting third-party libraries provide common functions. Armin Ronacher, who leads an international group of Python enthusiasts named Pocco, develops it. Flask is based on Werkzeug WSGI toolkit and Jinja2 template engine. Both are Pocco projects.

However, Flask supports extensions that can add application features as if they were implemented in Flask itself. Extensions exist for object-relational mappers, form validation, upload handling, various open authentication technologies and several common framework related tools.^[3]

Applications that use the Flask framework include Pinterest and LinkedIn.

2. NumPy :-

NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices. NumPy was created in 2005 by Travis Oliphant. It is an open source project and you can use it freely. NumPy stands for Numerical Python. NumPy targets the CPython reference implementation of Python, which is a nonoptimizing bytecode interpreter. Mathematical algorithms written for this version of Python often run much slower than compiled equivalents. NumPy addresses the slowness problem partly by providing multidimensional arrays and functions and operators that operate efficiently on arrays; using these requires rewriting some code, mostly inner loops, using NumPy.

3. OS :-

The main purpose of the OS module is to interact with your operating system. The primary use I find for it is to create folders, remove folders, move folders, and sometimes change the working directory. You can also access the names of files within a file path by doing `listdir()`. We do not cover that in this video, but that's an option. The `os` module is a part of the standard library, or `stdlib`, within Python 3

4. MATPLOTLIB :-

Matplotlib is a plotting library for the Python programming language and its numerical mathematics extension NumPy. It provides an object-oriented API for embedding plots into applications using general-purpose GUI toolkits like Tkinter, wxPython, Qt, or GTK. Matplotlib is a low level graph plotting library in python that serves as a visualization utility. Matplotlib was created by John D. Hunter. Matplotlib is open source and we can use it freely. Matplotlib is mostly written in python, a few segments are written in C, Objective-C and Javascript for Platform compatibility

5. PICKLE :-

Python pickle module is used for serializing and de-serializing a Python object structure. Any object in Python can be pickled so that it can be saved on disk. What pickle does is that it “serializes” the object first before writing it to file. Pickling is a way to convert a python object (list, dict, etc.) into a character stream. The idea is that this character stream contains all the information necessary to reconstruct the object in another python script

Advantages of using Pickle Module:

- 1. Recursive objects (objects containing references to themselves):**
- 2. Object sharing (references to the same object in different places):**
- 3. User-defined classes and their instances**

It is advisable not to unpickle data received from an untrusted source as they may pose security threat. However, the pickle module has no way of knowing or raise alarm while pickling malicious data

4. JOBLIB :-

Joblib is a set of tools to provide lightweight pipelining in Python. In particular:

1. transparent disk-caching of functions and lazy re-evaluation (memoize pattern)
2. easy simple parallel computing

Joblib is optimized to be fast and robust on large data in particular and has specific optimizations for *numpy* arrays

METHODOLOGY

1. DATASET DESCRIPTION :-

I have not scrapped the web pages to prepare the dataset. I have downloaded the dataset from cricsheet. The site gives us ball by ball details of matches. I then wrote a custom code to only include some of the features which I will be using.

The dataset contains ball by ball coverage of:

- 1188 ODI matches: data/odi.csv
- 1474 T-20 matches: data/t20.csv
- 617 IPL matches: data/ipl.csv
- mid: Each match is given a unique number
- date: When the match happened
- venue: Stadium where match is being played
- bat_team: Batting team name
- bowl_team: Bowling team name
- batsman: Batsman name who faced that ball

- bowler: Bowler who bowled that ball
- runs: Total runs scored by team at that instance
- wickets: Total wickets fallen at that instance
- overs: Total overs bowled at that instance
- runs_last_5: Total runs scored in last 5 overs
- wickets_last_5: Total wickets that fell in last 5 overs
- striker: $\max(\text{runs scored by striker}, \text{runs scored by non-striker})$
- non-striker: $\min(\text{runs scored by striker}, \text{runs scored by non-striker})$
- total: Total runs scored by batting team after first innings

I followed the general machine learning workflow step-by-step:

- 1. Data cleaning and formatting.**
- 2. Exploratory data analysis.**
- 3. Feature engineering and selection.**
- 4. Compare several machine learning models on a performance metric.**
- 5. Perform hyper-parameter tuning on the best model.**
- 6. Evaluate the best model on the testing set.**
- 7. Interpret the model results.**
- 8. Draw conclusions and document work.**

Then I loaded the csv file containing the details of each team's history in previous world cups. I also loaded the csv file containing the results of matches played between 2010 and 2017.

Importing necessary libraries :

```
In [60]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib.ticker as ticker
import matplotlib.ticker as plticker
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import make_classification
```

```
world_cup.head()
```

	Team	Group	Previous \r\nappearances	Previous \r\nTitles	Previous\r\n finals	Previous\r\n semifinals	Current \r\n rank
0	England	A	11	0	3	5	1
1	South Africa	A	6	0	0	4	3
2	West Indies	A	11	2	3	4	8
3	Pakistan	A	11	1	2	6	6
4	New Zealand	A	11	0	1	7	4

```
In [50]: results.head()
```

```
Out[50]:
```

	date	Team_1	Team_2	Winner	Margin	Ground
0	4-Jan-10	Bangladesh	Sri Lanka	Sri Lanka	7 wickets	Dhaka
1	5-Jan-10	India	Sri Lanka	Sri Lanka	5 wickets	Dhaka
2	7-Jan-10	Bangladesh	India	India	6 wickets	Dhaka
3	8-Jan-10	Bangladesh	Sri Lanka	Sri Lanka	9 wickets	Dhaka
4	10-Jan-10	India	Sri Lanka	India	8 wickets	Dhaka

1. Data cleaning and formatting

Next, let's display the details of matches played by India.

```
In [51]: df = results[(results['Team_1'] == 'India') | (results['Team_2'] == 'India')]  
india = df.iloc[:]  
india.head()
```

```
Out[51]:
```

	date	Team_1	Team_2	Winner	Margin	Ground
1	5-Jan-10	India	Sri Lanka	Sri Lanka	5 wickets	Dhaka
2	7-Jan-10	Bangladesh	India	India	6 wickets	Dhaka
4	10-Jan-10	India	Sri Lanka	India	8 wickets	Dhaka
5	11-Jan-10	Bangladesh	India	India	6 wickets	Dhaka
6	13-Jan-10	India	Sri Lanka	Sri Lanka	4 wickets	Dhaka

I continued by creating a column to display the details of matches played in 2010 and taking it as a reference for future work.

```
In [52]: year = []
         for row in india['date']:
             year.append(int(row[7:]))
         india['match_year'] = year
         india_2010 = india[india.match_year >= 10]
         india_2010.count()
```

```
Out[52]: date          151
         Team_1         151
         Team_2         151
         Winner         151
         Margin         142
         Ground         151
         match_year     151
         dtype: int64
```

2. Exploratory data analysis

After that, I merged the details of the teams participating this year with their past results.

```
In [53]: worldcup_teams = ['England', 'South Africa', '', 'West Indies',
                           'Pakistan', 'New Zealand', 'Sri Lanka', 'Afghanistan',
                           'Australia', 'Bangladesh', 'India']
         df_teams_1 = results[results['Team_1'].isin(worldcup_teams)]
         df_teams_2 = results[results['Team_2'].isin(worldcup_teams)]
         df_teams = pd.concat((df_teams_1, df_teams_2))
         df_teams.drop_duplicates()
         df_teams.count()
```

```
Out[53]: date          1445
         Team_1         1445
         Team_2         1445
         Winner         1445
         Margin         1368
         Ground         1445
         dtype: int64
```

I deleted the columns like date of the match, margin of victory, and the ground on which the match was played. These features doesn't look important for our prediction.

```
In [54]: df_teams.head()
```

```
Out[54]:
```

	date	Team_1	Team_2	Winner	Margin	Ground
0	4-Jan-10	Bangladesh	Sri Lanka	Sri Lanka	7 wickets	Dhaka
1	5-Jan-10	India	Sri Lanka	Sri Lanka	5 wickets	Dhaka
2	7-Jan-10	Bangladesh	India	India	6 wickets	Dhaka
3	8-Jan-10	Bangladesh	Sri Lanka	Sri Lanka	9 wickets	Dhaka
4	10-Jan-10	India	Sri Lanka	India	8 wickets	Dhaka

```
In [55]: df_teams_2010 = df_teams.drop(['date', 'Margin', 'Ground'], axis=1)
df_teams_2010.head()
```

```
Out[55]:
```

	Team_1	Team_2	Winner
0	Bangladesh	Sri Lanka	Sri Lanka
1	India	Sri Lanka	Sri Lanka
2	Bangladesh	India	India
3	Bangladesh	Sri Lanka	Sri Lanka
4	India	Sri Lanka	India

3. Feature engineering and selection

This is probably the most important part in the machine learning workflow. Since the algorithm is totally dependent on how we feed data into it, feature engineering should be given topmost priority for every machine learning project.

Advantages of feature engineering

- **Reduces Overfitting:** Less redundant data means less opportunity to make decisions based on noise.
- **Improves Accuracy:** Less misleading data means modeling accuracy improves.
- **Reduces Training Time:** fewer data points reduce algorithm complexity and algorithms train faster.

So continuing with the work, I created the model. If team-1 won the match, I assigned it label 1, else if team-2 won, I assigned it label 2.

```
In [56]: df_teams_2010 = df_teams_2010.reset_index(drop=True)
df_teams_2010.loc[df_teams_2010.Winner == df_teams_2010.Team_1, 'winning_team']=1
df_teams_2010.loc[df_teams_2010.Winner == df_teams_2010.Team_2, 'winning_team']=2
df_teams_2010 = df_teams_2010.drop(['winning_team'], axis=1)

df_teams_2010.head()
```

Out[56]:

	Team_1	Team_2	Winner
0	Bangladesh	Sri Lanka	Sri Lanka
1	India	Sri Lanka	Sri Lanka
2	Bangladesh	India	India
3	Bangladesh	Sri Lanka	Sri Lanka
4	India	Sri Lanka	India

Then I converted team-1 and team-2 from categorical variables to continuous inputs using pandas function *pd.get_dummies*. This variable has only two answer choices: team 1 and team 2. It creates a new dataframe which consists of zeros and ones. The dataframe will have a one depending on the team of a particular game in this case.

Also, I separated training and test sets with 70% and 30% in training and validation sets respectively.

```
In [57]: final = pd.get_dummies(df_teams_2010, prefix=['Team_1', 'Team_2'], columns=['Team_1', 'Team_2'])
X = final.drop(['Winner'], axis=1)
y = final["winner"]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)
```

```
In [58]: final.head()
```

Out[58]:

1_2_Kenya	Team_2_Netherlands	Team_2_New Zealand	Team_2_Pakistan	Team_2_Scotland	Team_2_South Africa	Team_2_Sri Lanka	Team_2_U.A.E.	Team_2_West Indies	Team_2_Zimbabwe
0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	1	0	0	0

4. Compare several machine learning models on a performance metric

I used Logistic Regression, Support Vector Machines, Random Forests and K Nearest Neighbours for training the model.

Random Forest outperformed all other algorithms with 70% training accuracy and 67.5% test accuracy.

RFs train each tree independently, using a random sample of the data. This randomness helps to make the model more robust than a single decision tree, and less likely to overfit on the training data.

5. Perform hyperparameter tuning on the best model

```
In [61]: rf = RandomForestClassifier(n_estimators=100, max_depth=20, random_state=0)
         rf.fit(X_train, y_train)

         score = rf.score(X_train, y_train)
         score2 = rf.score(X_test, y_test)

         print("Training set accuracy: ", '%.3f'%(score))
         print("Test set accuracy: ", '%.3f'%(score2))
```

```
Training set accuracy:  0.700
Test set accuracy:    0.675
```

The popularity of the Random Forest model is explained by its various advantages:

- Accurate and efficient when running on large databases
- Multiple trees reduce the variance and bias of a smaller set or single tree
- Resistant to overfitting
- Can handle thousands of input variables without variable deletion
- Can estimate what variables are important in classification
- Provides effective methods for estimating missing data
- Maintains accuracy when a large proportion of the data is missing

6. Evaluate the best model on the testing set

Let's continue. I added ICC rankings of teams giving priority to higher ranked team to win this year.

```
In [63]: ranking = pd.read_csv('icc_rankings.csv')
          fixtures = pd.read_csv('fixtures.csv')

          pred_set = []
```

Next, I added new columns with ranking position for each team and slicing the dataset for first 45 games since there are 45 league stage games in total.

```
In [64]: fixtures.insert(1, 'first_position', fixtures['Team_1'].map(ranking.set_index('Team')['Position']))
fixtures.insert(2, 'second_position', fixtures['Team_2'].map(ranking.set_index('Team')['Position']))

fixtures = fixtures.iloc[:45, :]
fixtures.tail()
```

```
Out[64]:
```

	Round Number	first_position	second_position	Date	Location	Team_1	Team_2	Group	Result
40	1	1.0	4.0	3/7/2019	Riverside Ground, Chester-le-Street	England	New Zealand	Group A	NaN
41	1	10.0	8.0	4/7/2019	Headingley, Leeds	Afghanistan	West Indies	Group A	NaN
42	1	6.0	7.0	5/7/2019	Lord's, London	Pakistan	Bangladesh	Group A	NaN
43	1	9.0	2.0	6/7/2019	Headingley, Leeds	Sri Lanka	India	Group A	NaN
44	1	5.0	3.0	6/7/2019	Emirates Old Trafford, Manchester	Australia	South Africa	Group A	NaN

Then I added teams to new prediction dataset based on ranking position of each team.

```
In [65]: for index, row in fixtures.iterrows():
          if row['first_position'] < row['second_position']:
              pred_set.append({'Team_1': row['Team_1'], 'Team_2': row['Team_2'], 'winning_team': None})
          else:
              pred_set.append({'Team_1': row['Team_2'], 'Team_2': row['Team_1'], 'winning_team': None})

pred_set = pd.DataFrame(pred_set)
backup_pred_set = pred_set
pred_set.head()
```

```
Out[65]:
```

	Team_1	Team_2	winning_team
0	England	South Africa	None
1	Pakistan	West Indies	None
2	New Zealand	Sri Lanka	None
3	Australia	Afghanistan	None
4	South Africa	Bangladesh	None

After that, I added scripts for getting dummy variables and added missing columns compared to model training dataset.

```
In [66]: pred_set = pd.get_dummies(pred_set, prefix=['Team_1', 'Team_2'], columns=['Team_1', 'Team_2'])

missing_cols = set(final.columns) - set(pred_set.columns)
for c in missing_cols:
    pred_set[c] = 0
pred_set = pred_set[final.columns]

pred_set = pred_set.drop(['Winner'], axis=1)
pred_set.head()
```

Out[66]:

	Team_1_Afghanistan	Team_1_Australia	Team_1_Bangladesh	Team_1_Canada	Team_1_England	Team_1_India	Team_1_Ireland	Team_1_Kenya	Team_1_Nethe
0	0	0	0	0	1	0	0	0	
1	0	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	0	
3	0	1	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	

5 rows x 35 columns

7. Interpret the model results

Finally, the below code is for getting the results for each and every league stage match.


```
In [67]: predictions = rf.predict(pred_set)
for i in range(fixture.shape[0]):
    print(fixture.iloc[i, 1] + " and " + fixture.iloc[i, 0])
    if predictions[i] == 1:
        print("Winner: " + fixture.iloc[i, 1])

    else:
        print("Winner: " + fixture.iloc[i, 0])
    print("")
```

Bangladesh and India
Winner: India

New Zealand and England
Winner: England

Afghanistan and West Indies
Winner: West Indies

Bangladesh and Pakistan
Winner: Pakistan

Sri Lanka and India
Winner: India

Australia and South Africa
Winner: South Africa

```
In [68]: semi = [('New Zealand', 'India'),
                 ('England', 'South Africa')]
```

And then I created a function to repeat the above work. This is the final function to predict the winner of ICC Cricket World Cup 2019.

```
In [70]: def clean_and_predict(matches, ranking, final, logreg):

    # Initialization of auxiliary list for data cleaning
    positions = []

    # Loop to retrieve each team's position according to ICC ranking
    for match in matches:
        positions.append(ranking.loc[ranking['Team'] == match[0], 'Position'].iloc[0])
        positions.append(ranking.loc[ranking['Team'] == match[1], 'Position'].iloc[0])

    # Creating the DataFrame for prediction
    pred_set = []

    # Initializing iterators for while loop
    i = 0
    j = 0

    # 'i' will be the iterator for the 'positions' list, and 'j' for the list of matches (list of tuples)
    while i < len(positions):
        dict1 = {}

        # If position of first team is better then this team will be the 'Team_1' team, and vice-versa
        if positions[i] < positions[i + 1]:
            dict1.update({'Team_1': matches[j][0], 'Team_2': matches[j][1]})
        else:
            dict1.update({'Team_1': matches[j][1], 'Team_2': matches[j][0]})

        # Append updated dictionary to the list, that will later be converted into a DataFrame
        pred_set.append(dict1)
        i += 2
        j += 1

    # Convert list into DataFrame
    pred_set = pd.DataFrame(pred_set)
    backup_pred_set = pred_set

    # Get dummy variables and drop winning team column
    pred_set = pd.get_dummies(pred_set, prefix=['Team_1', 'Team_2'], columns=['Team_1', 'Team_2'])

    # Add missing columns compared to the model's training dataset
    missing_cols2 = set(final.columns) - set(pred_set.columns)
    for c in missing_cols2:
        pred_set[c] = 0
    pred_set = pred_set[final.columns]

    pred_set = pred_set.drop(['Winner'], axis=1)

    # Predict!
    predictions = logreg.predict(pred_set)
    for i in range(len(pred_set)):
        print(backup_pred_set.iloc[i, 1] + " and " + backup_pred_set.iloc[i, 0])
        if predictions[i] == 1:
            print("Winner: " + backup_pred_set.iloc[i, 1])
        else:
            print("Winner: " + backup_pred_set.iloc[i, 0])
        print("")
```

I ran the function for semi-finals prediction.

```
In [71]: clean_and_predict(semi, ranking, final, rf)
```

```
New Zealand and India  
Winner: India
```

```
South Africa and England  
Winner: England
```

New Zealand and India

Winner: India

South Africa and England

Winner: England

Hence the two finalists are India and England which is quite evident as they are considered the favourites to win this year. Also, they are first and second ranked team in ICC rankings.

```
In [72]: finals = [('India', 'England')]
```

```
In [73]: clean_and_predict(finals, ranking, final, rf)
```

```
India and England  
Winner: England
```

8. Draw conclusions and document work

Finally on running the main function.

India and England

Winner: England

CONCLUSION

In these project I have used **K - Nearest Neighbour**, **Random Forest**, **Gradient Boost** **Ada Boost With Random Forest** and the highest accuracy achived is **as** follows:

Sr no	Model Name	Accuracy %
0	K - Nearest Neighbor	67.82%
1	Random Forest	85 %
2	Ada Boost With Random Forest	63.47 %

3	Gradient Boosting	69.91%
---	-------------------	--------

The outcome of this analysis shows that the **random forest** algorithm is the most powerful algorithm for heart disease prediction, with an accuracy score of 85

It can be concluded that there is a huge scope for machine learning algorithms in predicting cardiovascular diseases or heart related diseases.

REFERENCES

1. <https://www.kaggle.com/datasets?search=cricket>https://en.wikipedia.org/wiki/Cardiovascular_disease.
2. <https://stackoverflow.com>
3. <https://cricsheet.org/downloads/><https://www.medicalnewstoday.com/articles/237191>
4. <https://flask.palletsprojects.com/en/2.1.x/>

GITHUB LINK: