

# **IPL Win Probability Predictor**

## **A PROJECT REPORT**

*In partial fulfilment of the requirements for the award of the degree*

### **BACHELOR OF TECHNOLOGY**

**Electrical Engineering**

*Under the guidance of*

**MAHENDRA DATTA**

**BY**

**MD SERIF**



**KALYANI GOVERNMENT ENGINEERING COLLEGE**

**In association with**



**(ISO9001:2015)**

SDF Building, Module #132, Ground Floor, Salt Lake City, GP Block, Sector V, Kolkata, West Bengal 700091

***(Note: All entries of the proforma of approval should be filled up with appropriate and complete information. Incomplete proforma of approval in any respect will be summarily rejected.)***

1. Title of the Project: **IPL WIN PROBABILITY PREDICTOR**
2. Project Members: **MD SERIF**
3. Name of the guide: **Mr. MAHENDRA DUTTA**
4. Address: Ardent Computech Pvt. Ltd  
(An ISO 9001:2015 Certified)  
SDF Building, Module #132, Ground Floor, Salt  
Lake City, GP Block, Sector V, Kolkata, West  
Bengal, 700091

**Project Version Control History**

Version	Primary Author	Description of Version	Date Completed
Final	MD SERIF	Project Report	4th AUGUST,2022

*Md Serif*

Signature of Team Member

Date: 04/08/2022

*Mahendra Dutta*

Signature of Approver

Date: 04/08/2022

For Office Use Only

**MR. MAHENDRA**

**DUTTA**

**Approved**

**Not Approved**

## **DECLARATION**

We hereby declare that the project work being presented in the project proposal entitled “**IPL WIN PROBABILITY PREDICTOR**” in partial fulfilment of the requirements for the award of the degree of **BACHELOR OF TECHNOLOGY** at **ARDENT COMPUTECH PVT. LTD, SALT LAKE, KOLKATA, WEST BENGAL**, is an authentic work carried out under the guidance of **MR. MAHENDRA DUTTA**. The matter embodied in this project work has not been submitted elsewhere for the award of any degree of our knowledge and belief.

Date: 04/08/2022

Name of the Student: MD SERIF

*Md Serif*

*Signature of the students:*



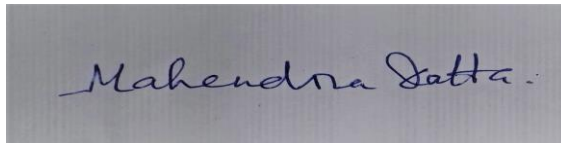
**Ardent Computech Pvt. Ltd (An ISO 9001:2015 Certified)**

SDF Building, Module #132, Ground Floor, Salt Lake City, GP Block, Sector V, Kolkata, West Bengal 700091

## **CERTIFICATE**

This is to certify that this proposal of minor project entitled “**IPL WIN PROBABILITY PREDICTOR**” is a record of Bonafide work, carried out by *MD SERIF* under my guidance at **ARDENT COMPUTECH PVT LTD**. In my opinion, the report in its present form is in partial fulfilment of the requirements for the award of the degree of **BACHELOR OF TECHNOLOGY** and as per regulations of the **ARDENT®**. To the best of my knowledge, the results embodied in this report, are original in nature and worthy of incorporation in the present version of the report.

**Guide / Supervisor**

A rectangular box containing a handwritten signature in blue ink that reads "Mahendra Dutta".

-----  
**MR. MAHENDRA DUTTA**

Project Engineer

**Ardent Computech Pvt. Ltd (An ISO 9001:2015 Certified)**

SDF Building, Module #132, Ground Floor, Salt Lake City, GP Block, Sector V, Kolkata, West Bengal 700091

## **ACKNOWLEDGEMENT**

Success of any project depends largely on the encouragement and guidelines of many others. I take this sincere opportunity to express my gratitude to the people who have been instrumental in the successful completion of this project work.

I would like to show our greatest appreciation to **Mr. MAHENDRA DUTTA**, Project Engineer at Ardent, Kolkata. I always feel motivated and encouraged every time by his valuable advice and constant inspiration; without his encouragement and guidance this project would not have materialized.

Words are inadequate in offering our thanks to the other trainees, project assistants and other members at Ardent Computech Pvt. Ltd. for their encouragement and cooperation in carrying out this project work. The guidance and support received from all the members and who are contributing to this project, was vital for the success of this project.

# CONTENTS

- Overview
- History of Python
- Environment Setup
- Basic Syntax
- Variable Types
- Functions
- Modules
- Packages
- Machine Learning
  - Supervised and Unsupervised Learning
  - NumPy
  - SciPy
  - Scikit-learn
  - Pandas
  - Regression Analysis
  - Matplotlib
  - Clustering
- IPL WIN PROBABILITY PREDICTOR
  1. Introduction
  2. Problem Statement
  3. Advantages & Disadvantages
  4. Future Scope

## **OVERVIEW**

Python is a high-level, interpreted, interactive and object-oriented scripting language. Python is designed to be highly readable. It uses English keywords frequently where as other languages use punctuation, and has fewer syntactical constructions than other languages.

**Python is interpreted:** Python is processed at runtime by the interpreter. You do not need to compile your program before executing it. This is similar to Perl and PHP.

**Python is Interactive:** You can actually sit at a Python prompt and interact with the interpreter directly to write your programs.

**Python is Object-Oriented:** Python supports Object-Oriented style or technique of programming that encapsulates code within objects.

**Python is a Beginner's Language:** Python is a great language for the beginner-level programmers and supports the development of a wide range of applications from simple text processing to WWW browsers to games.

## HISTORY OF PYTHON

Python was developed by Guido van Rossum in the late eighties and early nineties at the National Research Institute for Mathematics and Computer Science in the Netherlands. Python is derived from many other languages, including ABC, Modula-3, C, C++, Algol-68, Small Talk, UNIX shell, and other scripting languages. Python is copyrighted. Like Perl, Python source code is now available under the GNU General Public License (GPL). Python is now maintained by a core development team at the institute, although Guido van Rossum still holds a vital role in directing its progress.

## FEATURES OF PYTHON

Easy-to-learn: Python has few Keywords, simple structure and clearly defined syntax. This allows a student to pick up the language quickly.

Easy-to-Read: Python code is more clearly defined and visible to the eyes.

Easy -to-Maintain: Python's source code is fairly easy-to-maintain.

A broad standard library: Python's bulk of the library is very portable and cross platform compatible on UNIX, Windows, and Macintosh.

Interactive Mode: Python has support for an interactive mode which allows interactive testing and debugging of snippets of code.

Portable: Python can run on the wide variety of hardware platforms and has the same interface on all platforms.

Extendable: You can add low level modules to the python interpreter. These modules enables programmers to add to or customize their tools to be more efficient.

Databases: Python provides interfaces to all major commercial databases.

GUI Programming: Python supports GUI applications that can be created and ported to many system calls, libraries, and windows systems, such as Windows MFC, Macintosh, and the X Window system of Unix.

Scalable: Python provides a better structure and support for large programs than shell scripting.

Apart from the above-mentioned features, Python has a big list of good features, few are listed below:



- It support functional and structured programming methods as well as OOP.
- It can be used as a scripting language or can be compiled to byte code for building large applications.
- It provides very high level dynamic datatypes and supports dynamic type checking.
- It supports automatic garbage collections.
- It can be easily integrated with C, C++, COM, ActiveX, CORBA and JAVA.

## **ENVIRONMENT SETUP**

Open a terminal window and type "python" to find out if it is already installed and which version is installed.

- UNIX (Solaris, Linux, FreeBSD, AIX, HP/UX, SunOS, IRIX, etc.)
- Win 9x/NT/2000
- Macintosh (Intel, PPC, 68K)
- OS/2
- DOS (multiple versions)
- PalmOS
- Nokia mobile phones
- Windows CE
- Acorn/RISC OS

# **BASIC SYNTAX OF PYTHON PROGRAM**

Type the following text at the Python prompt and press the Enter –

```
>>> print "Hello, Python!"
```

*If you are running new version of Python, then you would need to use print statement with parenthesis as in **print ("Hello, Python!");**.*

However in Python version 2.4.3, this produces the following result –

Hello, Python!

## **Python Identifiers**

A Python identifier is a name used to identify a variable, function, class, module or other object. An identifier starts with a letter A to Z or a to z or an underscore (\_) followed by zero or more letters, underscores and digits (0 to 9).

Python does not allow punctuation characters such as @, \$, and % within identifiers. Python is a case sensitive programming language.

## **Python Keywords**

The following list shows the Python keywords. These are reserved words and you cannot use them as constant or variable or any other identifier names. All the Python keywords contain lowercase letters only.

**And, exec, not**  
**Assert, finally, or**  
**Break, for, pass**  
**Class, from, print**  
**continue, global, raise**  
**def, if, return**  
**del, import, try**  
**elif, in, while**  
**else, is, with**  
**except, lambda, yield**

## **Lines & Indentation**

Python provides no braces to indicate blocks of code for class and function definitions or flow control. Blocks of code are denoted by line indentation, which is rigidly enforced.

The number of spaces in the indentation is variable, but all statements within the block must be indented the same amount. For example –

```
if True:
    print "True"
else:
    print "False"
```

## **Command Line Arguments**

Many programs can be run to provide you with some basic information about how they should be run. Python enables you to do this with -h –

```
$ python -h
usage: python [option]...[-c cmd|-m mod | file | -][arg]...
```

Options and arguments (and corresponding environment variables):

- c cmd: program passed in as string(terminates option list)
- d : debug output from parser (also PYTHONDEBUG=x)
- E : ignore environment variables (such as PYTHONPATH)
- h : print this help message and exit [ etc.]

# VARIABLE TYPES

Variables are nothing but reserved memory locations to store values. This means that when you create a variable you reserve some space in memory.

## Assigning Values to Variables

Python variables do not need explicit declaration to reserve memory space. The declaration happens automatically when you assign a value to a variable. The equal sign (=) is used to assign values to variables.

```
counter=10      # An integer assignment
weight=10.60    # A floating point
name="Ardent"   # A string
```

## Multiple Assignment

Python allows you to assign a single value to several variables simultaneously. For example –

```
a = b = c = 1
a,b,c = 1,2,"hello"
```

## Standard Data Types

The data stored in memory can be of many types. For example, a person's age is stored as a numeric value and his or her address is stored as alphanumeric characters. Python has five standard data types –

- ☐ String
- ☐ List
- ☐ Tuple
- ☐ Dictionary
- ☐ Number

## Data Type Conversion

Sometimes, you may need to perform conversions between the built-in types. To convert between types, you simply use the type name as a function.

There are several built-in functions to perform conversion from one data type to another.

Sr.No.	Function & Description
1	<b>int(x [,base])</b> Converts x to an integer. base specifies the base if x is a string
2	<b>long(x [,base] )</b> Converts x to a long integer. base specifies the base if x is a string.
3	<b>float(x)</b> Converts x to a floating-point number.
4	<b>complex(real [,imag])</b> Creates a complex number.
5	<b>str(x)</b> Converts object x to a string representation.
6	<b>repr(x)</b> Converts object x to an expression string.
7	<b>eval(str)</b> Evaluates a string and returns an object.
8	<b>tuple(s)</b> Converts s to a tuple.
9	<b>list(s)</b> Converts s to a list.

# FUNCTIONS

## Defining a Function

- `def function name( parameters ):`  
    `"function_docstring"`  
    function suite  
    `return [expression]`

## Pass by reference vs Pass by value

All parameters (arguments) in the Python language are passed by reference. It means if you change what a parameter refers to within a function, the change also reflects back in the calling function. For example –

*# Function definition is here*

```
def change me(mylist):  
    "This changes a passed list into this function"  
    mylist.append([1,2,3,4]);  
    print"Values inside the function: ",mylist  
    return
```

*# Now you can call changeme function*

```
mylist=[10,20,30];  
change me(mylist);  
print"Values outside the function: ",mylist
```

Here, we are maintaining reference of the passed object and appending values in the same object. So, this would produce the following result –

Values inside the function: [10, 20, 30, [1, 2, 3, 4]]  
Values outside the function: [10, 20, 30, [1, 2, 3, 4]]

## Global vs. Local variables

Variables that are defined inside a function body have a local scope, and those defined outside have a global scope. For Example-

```
total=0;           # This is global variable.
```

```
# Function definition is here
```

```
def sum( arg1, arg2 ):
```

```
# Add both the parameters and return them."
```

```
total= arg1 + arg2;    # Here total is local variable.  
print"Inside the function local total: ", total  
return total;
```

```
# Now you can call sum function
```

```
sum(10,20);  
Print"Outside the function global total: ", total
```

When the above code is executed, it produces the following result –

```
Inside the function local total: 30  
Outside the function global total: 0
```

## MODULES

A module allows you to logically organize your Python code. Grouping related code into a module makes the code easier to understand and use. A module is a Python object with arbitrarily named attributes that you can bind and reference.

The Python code for a module named *aname* normally resides in a file named *aname.py*. Here's an example of a simple module, support.py

```
def print_func(par):  
    print"Hello : ", par  
    return
```

### The *import* Statement

You can use any Python source file as a module by executing an import statement in some other Python source file. The *import* has the following syntax

—

```
Import module1 [, module2 [... moduleN]
```



## PACKAGES

A package is a hierarchical file directory structure that defines a single Python application environment that consists of modules and sub packages and sub-subpackages, and so on.

Consider a file *Pots.py* available in *Phone* directory. This file has following line of source code –

```
def Pots ():  
    print "I'm Pots Phone"
```

Similar way, we have another two files having different functions with the same name as above –

- ☐ *Phone/Isdn.py* file having function *Isdn ()*
- ☐ *Phone/G3.py* file having function *G3 ()*

Now, create one more file *\_\_init\_\_.py* in *Phone* directory –

- ☐ *Phone/\_\_init\_\_.py*

To make all of your functions available when you've imported *Phone*, you need to put explicit import statements in *\_\_init\_\_.py* as follows –

```
from Pots import Pots  
from Isdn import Isdn  
from G3 import
```

# INTRODUCTION TO MACHINE LEARNING

**Machine learning** is a field of computer science that gives computers the ability to learn without being explicitly programmed.

**Arthur Samuel**, an American pioneer in the field of computer gaming and artificial intelligence, coined the term "Machine Learning" in 1959 while at IBM. Evolved from the study of pattern recognition and computational learning theory in artificial intelligence, machine learning explores the study and construction of algorithms that can learn from and make predictions on data

Machine learning tasks are typically classified into two broad categories, depending on whether there is a learning "signal" or "feedback" available to a learning system:-

## SUPERVISED LEARNING

**Supervised learning** is the machine learning task of inferring a function from *labelled training data*.<sup>[1]</sup> The training data consist of a set of *training examples*. In supervised learning, each example is a *pair* consisting of an input object (typically a vector) and a desired output value.

A supervised learning algorithm analyses the training data and produces an inferred function, which can be used for mapping new examples. An optimal scenario will allow for the algorithm to correctly determine the class labels for unseen instances. This requires the learning algorithm to generalize from the training data to unseen situations in a "reasonable" way.

## UNSUPERVISED LEARNING

**Unsupervised learning** is the machine learning task of inferring a function to describe hidden structure from "unlabelled" data (a classification or categorization is not included in the observations). Since the examples given to the learner are unlabelled, there is no evaluation of the accuracy of the structure that is output by the relevant algorithm—which is one way of distinguishing unsupervised learning from supervised learning and reinforcement learning.

A central case of unsupervised learning is the problem of density estimation in statistics, though unsupervised learning encompasses many other problems (and solutions) involving summarizing and explaining key features of the data.

## 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. The ancestor of NumPy, Numeric, was originally created by Jim Hugunin.

NumPy targets the CPython reference implementation of Python, which is a non-optimizing bytecode interpreter. Mathematical algorithms written for this version of Python often run much slower than compiled equivalents.

Using NumPy in Python gives functionality comparable to MATLAB since they are both interpreted, and they both allow the user to write fast programs as long as most operations work on arrays or matrices instead of scalars.

## NUMPY ARRAY

NumPy's main object is the homogeneous multidimensional array. It is a table of elements (usually numbers), all of the same type, indexed by a tuple of positive integers. In NumPy dimensions are called *axes*. The number of axes is *rank*.

For example, the coordinates of a point in 3D space [1, 2, 1] is an array of rank 1, because it has one axis. That axis has a length of 3. In the example pictured below, the array has rank 2 (it is 2-dimensional). The first dimension (axis) has a length of 2, the second dimension has a length of 3.

```
[[1., 0., 0.],  
 [ 0., 1., 2.]]
```

NumPy's array class is called *ndarray*. It is also known by the alias.

## SLICING NUMPY ARRAY

**Import numpy as np**

```
a = np.array ([[1, 2, 3],[3,4,5],[4,5,6]])
```

```
print 'Our array is:'
```

```
Print a
```

```
print '\n'
```

```
print 'The items in the second column are:'
```

```
print a[:,1]
```

```
print '\n'
```

```
print 'The items in the second row are:'  
print a[1...]  
print '\n'
```

```
print 'The items columns 1 onwards are:'  
print a [...,1:]
```

### **OUTPUT**

Our array is:

```
[[1 2 3]  
 [3 4 5]  
 [4 5 6]]
```

The items in the second column are:

```
[2 4 5]
```

The items in the second row are:

```
[3 4 5]
```

The items column 1 onwards are:

```
[[2 3]  
 [4 5]  
 [5 6]]
```

## **SCIPY**

modules for optimization, linear algebra, integration, interpolation, special functions, FFT, signal and image processing, ODE solvers and other tasks common in science and engineering.

SciPy builds on the NumPy array object and is part of the NumPy stack which includes tools like Matplotlib, pandas and SymPy, and an expanding set of scientific computing libraries. This NumPy stack has similar users to other applications such as MATLAB, GNU Octave, and Scilab. The NumPy stack is also sometimes referred to as the SciPy stack.

## The SciPy Library/Package

The SciPy package of key algorithms and functions core to Python's scientific computing capabilities. Available sub-packages include:

- **constants:** physical constants and conversion factors (since version 0.7.0)
- **cluster:** hierarchical clustering, vector quantization, K-means
- **fftpack:** Discrete Fourier Transform algorithms
- **integrate:** numerical integration routines
- **interpolate:** interpolation tools
- **io:** data input and output
- **lib:** Python wrappers to external libraries
- **linalg:** linear algebra routines
- **misc:** miscellaneous utilities (e.g. image reading/writing)
- **ndimage:** various functions for multi-dimensional image processing
- **optimize:** optimization algorithms including linear programming
- **signal:** signal processing tools
- **sparse:** sparse matrix and related algorithms
- **spatial:** KD-trees, nearest neighbours, distance functions
- **special:** special functions
- **stats:** statistical functions
- **weave:** tool for writing C/C++ code as Python multiline strings

## Data Structures

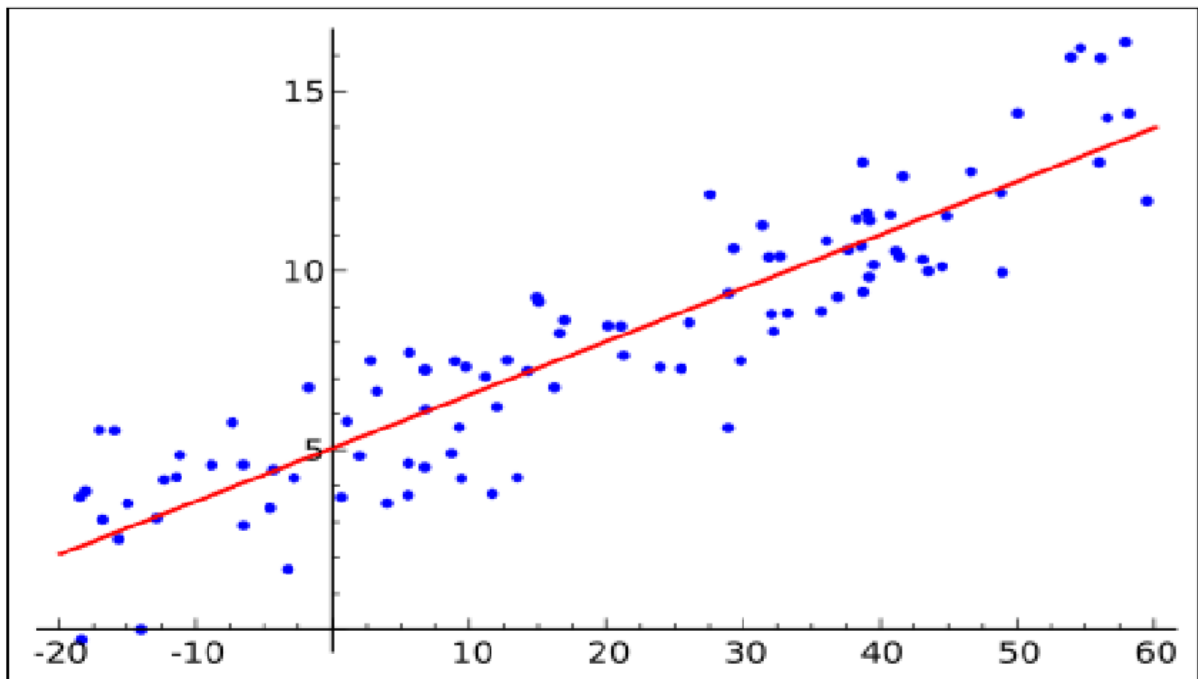
The basic data structure used by SciPy is a multidimensional array provided by the NumPy module. NumPy provides some functions for linear algebra, Fourier transforms and random number generation, but not with the generality of the equivalent functions in SciPy. NumPy can also be used as an efficient multi-dimensional container of data with arbitrary data-types. This allows NumPy to seamlessly and speedily integrate with a wide variety of databases. Older versions of SciPy used Numeric as an array type, which is now deprecated in favour of the newer NumPy array code.

## SCIKIT-LEARN

**Scikit-learn** is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support vector machines, random forests, gradient boosting, *k*-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy.

The scikit-learn project started as scikits.learn, a Google Summer of Code project by David Cournapeau. Its name stems from the notion that it is a "SciKit" (SciPy Toolkit), a separately-developed and distributed third-party extension to SciPy.[4] The original codebase was later rewritten by other developers. In 2010 Fabian Pedregosa, Gael Varoquaux, Alexandre Gramfort and Vincent Michel, all from INRIA took leadership of the project and made the first public release on February the 1st 2010[5]. Of the various scikits, scikit-learn as well as scikit-image were described as "well-maintained and popular" in November 2012.

## REGRESSION ANALYSIS



In statistical modelling, **regression analysis** is a set of statistical processes for estimating the relationships among variables. It includes many techniques for modelling and analysing several variables, when the focus is on the relationship between a dependent variable and one or more independent variables (or 'predictors'). More specifically, regression analysis helps one understand how the typical value of the dependent variable (or 'criterion variable') changes when any one of the independent variables is varied, while the other independent variables are held fixed.

Regression analysis is widely used for prediction and forecasting, where its use has substantial overlap with the field of machine learning. Regression analysis is also used to understand which among the independent variables are related to the dependent variable, and to explore the forms of these relationships. In restricted circumstances, regression analysis can be used to infer casual relationships between the independent and dependent variables. However this can lead to illusions or false relationships, so caution is advisable

### LINEAR REGRESSION

Linear regression is a linear approach for modelling the relationship between a scalar dependent variable  $y$  and one or more explanatory variables (or independent variables) denoted  $X$ . The case of one explanatory variable is called *simple linear regression*. For more than one explanatory variable, the process is called *multiple linear regression*.

In linear regression, the relationships are modelled using linear predictor functions whose unknown model parameters are estimated from the data. Such models are called *linear models*.

## LOGISTIC REGRESSION

Logistic regression, or logit regression, or logit model<sup>[1]</sup> is a regression model where the dependent variable (DV) is categorical. This article covers the case of a binary dependent variable—that is, where the output can take only two values, "0" and "1", which represent outcomes such as pass/fail, win/lose, alive/dead or healthy/sick. Cases where the dependent variable has more than two outcome categories may be analysed in multinomial logistic regression, or, if the multiple categories are ordered, in ordinal logistic regression. In the terminology of economics, logistic regression is an example of a qualitative response/discrete choice model.

## POLYNOMIAL REGRESSION

Polynomial regression is a form of regression analysis in which the relationship between the independent variable  $x$  and the dependent variable  $y$  is modelled as an  $n^{\text{th}}$  degree polynomial in  $x$ .

Polynomial regression fits a nonlinear relationship between the value of  $x$  and the corresponding conditional mean of  $y$ , denoted  $E(y | x)$ , and has been used to describe nonlinear phenomena such as the growth rate of tissues, the distribution of carbon isotopes in lake sediments, and the progression of disease epidemics.

Although *polynomial regression* fits a nonlinear model to the data, as a statistical estimation problem it is linear, in the sense that the regression function  $E(y | x)$  is linear in the unknown parameters that are estimated from the data.

## MATPLOTLIB

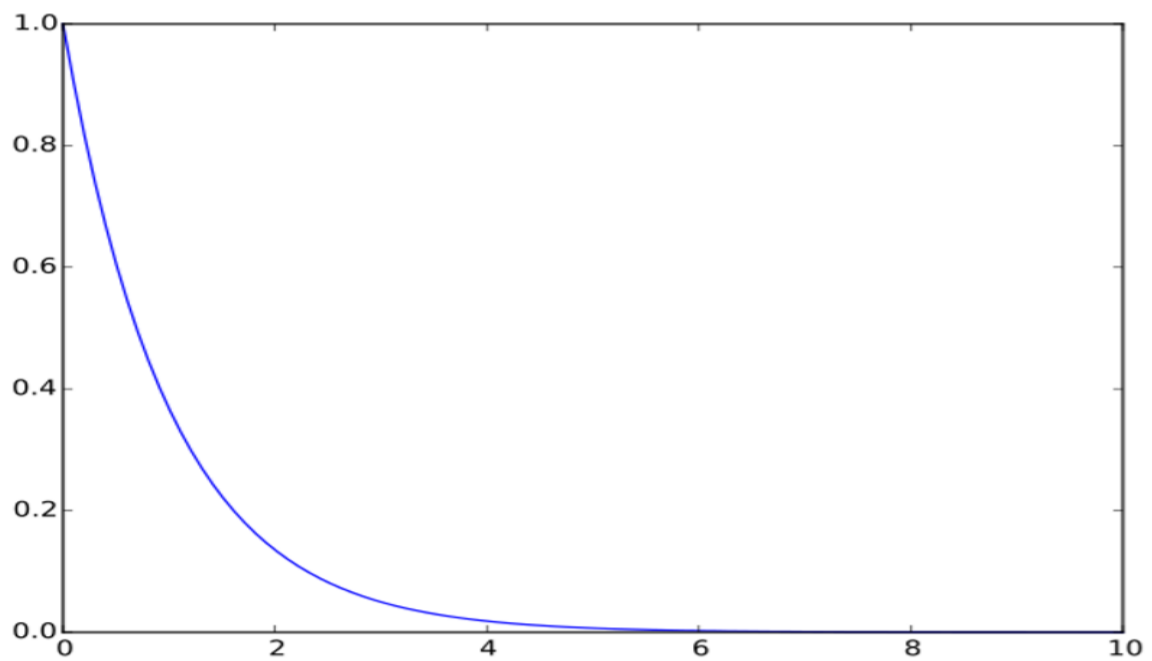
**Matplotlib** is a plotting library for the Python programming language and its numerical mathematics extension NumPy. It provides an object-oriented API for embedding plots into applications using general-purpose GUI toolkits like Tkinter, wxPython, Qt, or GTK+. There is also a procedural "pylab" interface based on a state machine (like OpenGL), designed to closely resemble that of MATLAB, though its use is discouraged. SciPy makes use of matplotlib.

## EXAMPLE



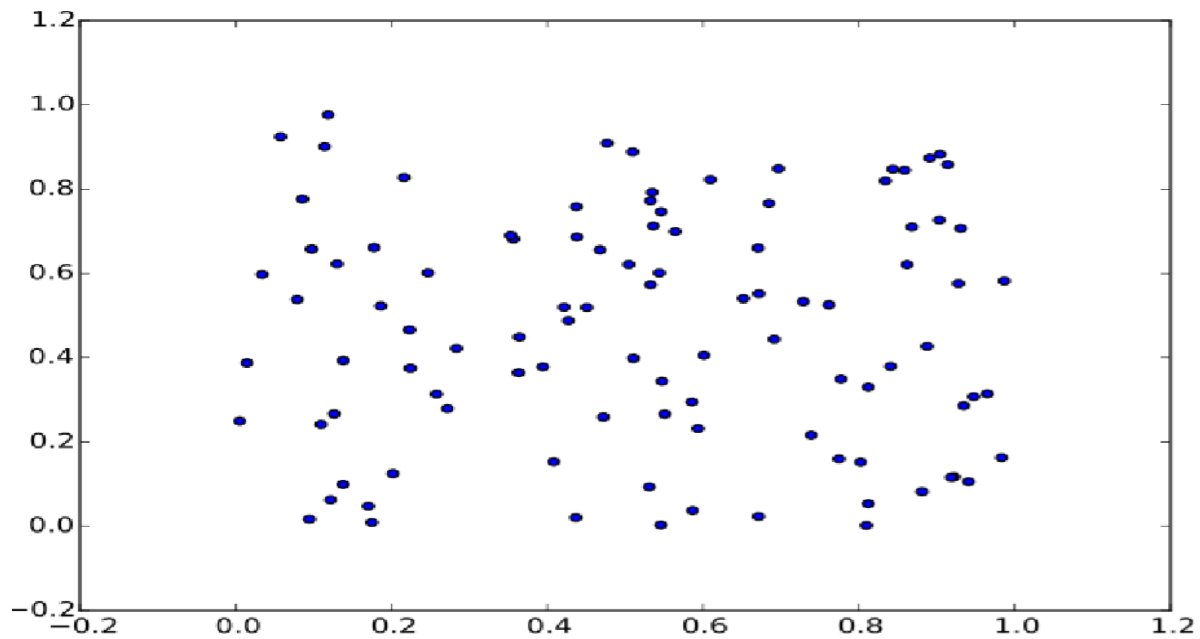
### ➤ LINE PLOT

```
>>>importmatplotlib.pyplotasplt
>>>importnumpyasnp
>>>a=np.linspace(0,10,100)
>>>b=np.exp(-a)
>>>plt.plot(a,b)
>>>plt.show()
```



### ➤ SCATTERPLOT

```
>>>importmatplotlib.pyplotasplt
>>>fromnumpy.randomimportrand
>>>a=rand(100)
>>>b=rand(100)
>>>plt.scatter(a,b)
>>>plt.show()
```



## PANDAS

In computer programming, **pandas** is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series. It is free software released under the three-clause BSD license. "Panel data", an econometrics term for multidimensional, structured data sets.

### LIBRARY FEATURES

- Data Frame object for data manipulation with integrated indexing.
- Tools for reading and writing data between in-memory data structures and different file formats.
- Data alignment and integrated handling of missing data.
- Reshaping and pivoting of data sets.
- Label-based slicing, fancy indexing, and sub setting of large data sets.
- Data structure column insertion and deletion.
- Group by engine allowing split-apply-combine operations on data sets.
- Data set merging and joining.
- Hierarchical axis indexing to work with high-dimensional data in a lower-dimensional data structure.
- Time series-functionality: Date range generation.

## CLUSTERING

**Cluster analysis** or **clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar (in some sense or another) to each other than to those in other groups (clusters). It is a main task of exploratory data mining, and a common technique for statistical data analysis, used in many fields, including machine learning, pattern recognition, image analysis, information retrieval, bioinformatics, data compression, and computer graphics.

Cluster analysis itself is not one specific algorithm, but the general task to be solved. It can be achieved by various algorithms that differ significantly in their notion of what constitutes a cluster and how to efficiently find them. Popular notions of clusters include groups with small distances among the cluster members, dense areas of the data space, intervals or particular statistical distributions. Clustering can therefore be formulated as a multi-objective optimization problem.

The appropriate clustering algorithm and parameter settings (including values such as the distance function to use, a density threshold or the number of expected clusters) depend on the individual data set and intended use of the results. Cluster analysis as such is not an automatic task, but an iterative process of knowledge discovery or interactive multi-objective optimization that involves trial and failure. It is often necessary to modify data pre-processing and model parameters until the result achieves the desired properties.

# ALGORITHM

- Data Collection
- Data Formatting
- Model Selection
- Training
- Testing

**Data Collection**: We have collected data sets of weather from online website. We have downloaded the .csv files in which information was present.

**Data Formatting**: The collected data is formatted into suitable data sets. We check the collinearity with mean temperature. The data sets which have collinearity nearer to 1.0 has been selected.

**Model Selection**: We have selected different models to minimize the error of the predicted value. The different models used are Linear Regression Linear Model, Ridge Linear model, Lasso Linear Model and Bayesian Ridge Linear Model.

**Training**: The data set was divided such that x\_train is used to train the model with corresponding x\_test values and some y\_train kept reserved for testing.

**Testing**: The model was tested with y\_train and stored in y\_predict. Both y\_train and y\_predict was compared.

# Actual Codes For **IPL WIN PROBABILITY PREDICTOR**

Introduction

Problem Statement

Advantages & Disadvantages

Future scope

# Explanation:

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline
from sklearn.metrics import accuracy_score
import pickle
```

```
In [2]: #importing the dataset from kaggle
delivery=pd.read_csv(r"C:\Users\Admin\Desktop\dataset\deliveries.csv")
matches=pd.read_csv(r"C:\Users\Admin\Desktop\dataset\matches.csv")
```

**\*Importing the necessary libraries and Dataset.**

**Datasets are downloaded from Kaggle.**

```
In [35]: delivery_df['player_dismissed'] = delivery_df['player_dismissed'].fillna("0")
delivery_df['player_dismissed'] = delivery_df['player_dismissed'].apply(lambda x:x if x == "0" else "1")
delivery_df['player_dismissed'] = delivery_df['player_dismissed'].astype('int')
wickets = delivery_df.groupby('match_id').cumsum()['player_dismissed'].values
delivery_df['wickets'] = 10 - wickets
delivery_df.head()
```

Out[35]:

	match_id	city	winner	total_runs_x	inning	batting_team	bowling_team	over	ball	batsman	...	batsman_runs	extra_runs	total_runs_y	playe
125	1	Hyderabad	Sunrisers Hyderabad	207	2	Royal Challengers Bangalore	Sunrisers Hyderabad	1	1	CH Gayle	...	1	0	1	
126	1	Hyderabad	Sunrisers Hyderabad	207	2	Royal Challengers Bangalore	Sunrisers Hyderabad	1	2	Mandeep Singh	...	0	0	0	
127	1	Hyderabad	Sunrisers Hyderabad	207	2	Royal Challengers Bangalore	Sunrisers Hyderabad	1	3	Mandeep Singh	...	0	0	0	
128	1	Hyderabad	Sunrisers Hyderabad	207	2	Royal Challengers Bangalore	Sunrisers Hyderabad	1	4	Mandeep Singh	...	2	0	2	
129	1	Hyderabad	Sunrisers Hyderabad	207	2	Royal Challengers Bangalore	Sunrisers Hyderabad	1	5	Mandeep Singh	...	4	0	4	

5 rows × 28 columns

**\*Creating the required dataset**

```
In [49]: from sklearn.model_selection import train_test_split
```

```
In [50]: x=final_df.iloc[:, :-1]
y=final_df.iloc[:, -1]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=1)
```

```
In [51]: x_train
```

Out[51]:

	batting_team	bowling_team	city	runs_left	balls_left	wickets	total_runs_x	current_run_rate	required_run_rate
69272	Chennai Super Kings	Kolkata Knight Riders	Kolkata	128	101	10	158	9.473684	7.603960
29103	Mumbai Indians	Kings XI Punjab	Centurion	69	72	8	119	6.250000	5.750000
98854	Rajasthan Royals	Kings XI Punjab	Chandigarh	89	32	5	179	6.136364	16.687500
25456	Delhi Daredevils	Chennai Super Kings	Johannesburg	143	100	8	163	6.000000	8.580000
44677	Delhi Daredevils	Mumbai Indians	Mumbai	116	85	9	183	11.485714	8.188235
...	...	...	...	...	...	...	...	...	...
29146	Mumbai Indians	Kings XI Punjab	Centurion	7	31	8	119	7.550562	1.354839
77065	Royal Challengers Bangalore	Rajasthan Royals	Bangalore	64	81	10	117	8.153846	4.740741
147841	Delhi Capitals	Rajasthan Royals	Delhi	58	73	7	124	8.425532	4.767123
50822	Delhi Daredevils	Deccan Chargers	Delhi	89	59	7	168	7.770492	9.050847
54967	Rajasthan Royals	Chennai Super Kings	Jaipur	158	90	9	196	7.600000	10.533333

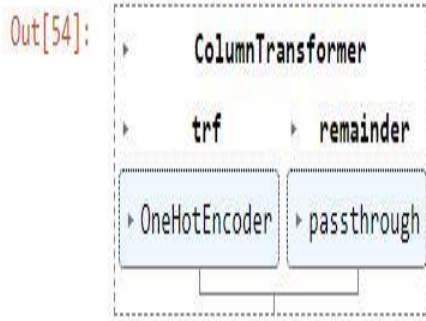
57073 rows × 9 columns

**\*Splitting the dataset**

```
In [52]: from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder
```

```
In [53]: trf=ColumnTransformer([
    ('trf',OneHotEncoder(sparse=True,drop='first'),['batting_team','bowling_team','city'])
],remainder='passthrough')
```

```
In [54]: trf
```

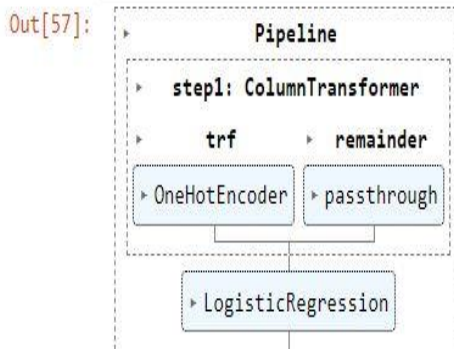


### \*Transforming the data

```
In [55]: from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline
```

```
In [56]: pipe=Pipeline(steps=[
    ('step1',trf),('step2',LogisticRegression(solver='liblinear'))
])
```

```
In [57]: pipe.fit(x_train,y_train)
```



### \*Building ML model



```
In [58]: y_pred=pipe.predict(x_test)
```

```
In [59]: y_pred
```

```
Out[59]: array([0, 0, 1, ..., 1, 1, 1], dtype=int64)
```

```
In [60]: from sklearn.metrics import accuracy_score
```

```
In [61]: accuracy_score(y_test,y_pred)
```

```
Out[61]: 0.8020884434788703
```

```
In [62]: pipe.predict_proba(x_test)[9]
```

```
Out[62]: array([0.70038767, 0.29961233])
```

**\*Predicting the Accuracy of the model.**

I have got 80% accuracy.

```
In [63]: teams
```

```
Out[63]: ['Sunrisers Hyderabad',  
          'Mumbai Indians',  
          'Royal Challengers Bangalore',  
          'Kolkata Knight Riders',  
          'Kings XI Punjab',  
          'Chennai Super Kings',  
          'Rajasthan Royals',  
          'Delhi Capitals']
```

```
In [64]: delivery_df['city'].unique()
```

```
Out[64]: array(['Hyderabad', 'Bangalore', 'Mumbai', 'Indore', 'Kolkata', 'Delhi',  
               'Chandigarh', 'Jaipur', 'Chennai', 'Cape Town', 'Port Elizabeth',  
               'Durban', 'Centurion', 'East London', 'Johannesburg', 'Kimberley',  
               'Bloemfontein', 'Ahmedabad', 'Cuttack', 'Nagpur', 'Dharamsala',  
               'Visakhapatnam', 'Pune', 'Raipur', 'Ranchi', 'Abu Dhabi',  
               'Sharjah', nan, 'Mohali', 'Bengaluru'], dtype=object)
```

```
In [65]: import pickle  
pickle.dump(pipe,open('pipe.pkl','wb'))
```

**\*Creating a write binary file**

```
In [6]: import streamlit as st
```

```
In [7]: st.title("IPL Win Predictor")
```

```
Out[7]: DeltaGenerator(_root_container=0, _provided_cursor=None, _parent=None, _block_type=None, _form_data=None)
```

```
In [ ]:
```

\*Importing 'streamlit' for web app

## Code for Web App

```
web_app.py x  pipe.pkl
web_app.py > ...
1  import streamlit as st
2  import pickle
3  import pandas as pd
4
5  teams = ['Sunrisers Hyderabad',
6          'Mumbai Indians',
7          'Royal Challengers Bangalore',
8          'Kolkata Knight Riders',
9          'Kings XI Punjab',
10         'Chennai Super Kings',
11         'Rajasthan Royals',
12         'Delhi Capitals']
13
14  cities = ['Hyderabad', 'Bangalore', 'Mumbai', 'Indore', 'Kolkata', 'Delhi',
15           'Chandigarh', 'Jaipur', 'Chennai', 'Cape Town', 'Port Elizabeth',
16           'Durban', 'Centurion', 'East London', 'Johannesburg', 'Kimberley',
17           'Bloemfontein', 'Ahmedabad', 'Cuttack', 'Nagpur', 'Dharamsala',
18           'Visakhapatnam', 'Pune', 'Raipur', 'Ranchi', 'Abu Dhabi',
19           'Sharjah', 'Mohali', 'Bengaluru']
20
21  pipe = pickle.load(open('pipe.pkl','rb'))
22  st.title('IPL Win Predictor')
23
```

web\_app.py x pipe.pkl

web\_app.py > ...

```
23
24 col1, col2 = st.columns(2)
25
26 with col1:
27     batting_team = st.selectbox('Select the batting team',sorted(teams))
28 with col2:
29     bowling_team = st.selectbox('Select the bowling team',sorted(teams))
30
31 selected_city = st.selectbox('Select host city',sorted(cities))
32
33 target = st.number_input('Target')
34
35 col3,col4,col5 = st.columns(3)
36
37 with col3:
38     score = st.number_input('Score')
39 with col4:
40     overs = st.number_input('Overs completed')
41 with col5:
42     wickets = st.number_input('Wickets out')
43
```

```
if st.button('Predict Probability'):
    runs_left = target - score
    balls_left = 120 - (overs*6)
    wickets = 10 - wickets
    current_run_rate = score/overs
    required_run_rate = (runs_left*6)/balls_left

    input_df = pd.DataFrame({'batting_team':[batting_team],'bowling_team':[bowling_team],
    'city':[selected_city],'runs_left':[runs_left],'balls_left':[balls_left],
    'wickets':[wickets],'total_runs_x':[target],'current_run_rate':[current_run_rate],'required_run_r
    #st.table(input_df)

    result = pipe.predict_proba(input_df)
    loss = result[0][0]
    win = result[0][1]
    st.header(batting_team + "- " + str(round(win*100)) + "%")
    st.header(bowling_team + "- " + str(round(loss*100)) + "%")
```

## Final GUI of the the web app

# IPL Win Predictor

Select the batting team

Chennai Super Kings

Select the bowling team

Kolkata Knight Riders

Select host city

Bengaluru

Target

190.00

Score

110.00

Overs completed

12.00

Wickets out

4.00

Predict Probability

**Chennai Super Kings- 64%**

**Kolkata Knight Riders- 36%**

## **CONCLUSION**

I have collected the raw data from online sources (Kaggle). Then I take this raw data and format it and also transform it.

Now I have selected a model for predicting the probability of winning an IPL match. I have used Logistic regression model to get the required output and got 80% accuracy using this ML algorithm. I have also deployed this model in a web app using streamlit library in python.

## **FUTURE SCOPE**

The data I have taken into consideration was not that big. The error can be minimized as well using other algorithms. The GUI in this project has scope of improvement in it.

## Bibliography:

- <https://scikit-learn.org/stable/>
- [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)
- <https://www.python.org/>
- <https://www.datacamp.com/tutorial/streamlit>
- <https://www.youtube.com/c/CampusX-official>

# **THANK YOU**