

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
```

```
In [4]: df=pd.read_csv(r"C3_bot_detection_data.csv")  
df
```

Out[4]:

	User ID	Username	Tweet	Retweet Count	Mention Count	Follower Count	Verified	Bot Label	Location
0	132131	flong	Station activity person against natural majori...	85	1	2353	False	1	Adk
1	289683	hinesstephanie	Authority research natural life material staff...	55	5	9617	True	0	Sand
2	779715	roberttran	Manage whose quickly especially foot none to g...	6	2	4363	True	0	Harris
3	696168	pmason	Just cover eight opportunity strong policy which.	54	5	2242	True	1	Martine
4	704441	noah87	Animal sign six data good or.	26	3	8438	False	1	Camac
...	...	...	...	...	...	...	...	...	...
49995	491196	uberg	Want but put card direction know miss former h...	64	0	9911	True	1	Kimberly
49996	739297	jessicamunoz	Provide whole maybe agree church respond most ...	18	5	9900	False	1	Gree
49997	674475	lynncunningham	Bring different everyone international capital...	43	3	6313	True	1	Debor
49998	167081	richardthompson	Than about single generation itself seek sell ...	45	1	6343	False	0	Stephe
49999	311204	daniel29	Here morning class various room human true bec...	91	4	4006	False	0	Novæ

50000 rows × 11 columns

In [7]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 41659 entries, 1 to 49999
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   User ID                41659 non-null  int64
1   Username               41659 non-null  object
2   Tweet                  41659 non-null  object
3   Retweet Count          41659 non-null  int64
4   Mention Count          41659 non-null  int64
5   Follower Count         41659 non-null  int64
6   Verified               41659 non-null  bool
7   Bot Label              41659 non-null  int64
8   Location               41659 non-null  object
9   Created At             41659 non-null  object
10  Hashtags               41659 non-null  object
dtypes: bool(1), int64(5), object(5)
memory usage: 3.5+ MB
```

In [8]: df=df.dropna()

In [9]: df.describe()

Out[9]:

	User ID	Retweet Count	Mention Count	Follower Count	Bot Label
<b>count</b>	41659.000000	41659.000000	41659.000000	41659.000000	41659.000000
<b>mean</b>	548640.613097	49.950911	2.515207	4990.867928	0.500204
<b>std</b>	259990.806985	29.195286	1.709249	2880.947193	0.500006
<b>min</b>	100025.000000	0.000000	0.000000	0.000000	0.000000
<b>25%</b>	321829.500000	25.000000	1.000000	2493.500000	0.000000
<b>50%</b>	548396.000000	50.000000	3.000000	4997.000000	1.000000
<b>75%</b>	772751.500000	75.000000	4.000000	7475.500000	1.000000
<b>max</b>	999995.000000	100.000000	5.000000	10000.000000	1.000000

In [11]: df1=df[['User ID','Retweet Count','Mention Count','Verified','Follower Count'],

```
In [28]: x=df1[['User ID','Retweet Count','Mention Count','Follower Count','Bot Label']]
y=df1['Verified']
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
lr=LogisticRegression()
lr.fit(x_train,y_train)
```

Out[28]: LogisticRegression()

```
In [16]: lr.predict(x_test)
```

```
Out[16]: array([ True,  True,  True, ...,  True,  True,  True])
```

```
In [17]: lr.score(x_test,y_test)
```

```
Out[17]: 0.5020803328532565
```

```
In [27]: from sklearn.preprocessing import StandardScaler  
fs=StandardScaler().fit_transform(x)  
logr=LogisticRegression()  
logr.fit(fs,y)
```

```
Out[27]: LogisticRegression()
```

```
In [21]: o=[[1,2,3,4,5]]  
prediction=logr.predict(o)  
print(prediction)
```

```
[False]
```

```
In [22]: logr.classes_
```

```
Out[22]: array([False,  True])
```

```
In [26]: logr.predict_proba(o)[0][0]
```

```
Out[26]: 0.5049151302812482
```

```
In [25]: logr.predict_proba(o)[0][1]
```

```
Out[25]: 0.4950848697187518
```

```
In [ ]:
```