

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

```
In [2]: from sklearn.linear_model import LogisticRegression
```

```
In [3]: df=pd.read_csv("detection.csv").dropna()  
df
```

Out[3]:

	User ID	Username	Tweet	Retweet Count	Mention Count	Follower Count	Verified	Bot Label	Loc
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	Camact
5	570928	james00	See wonder travel this suffer less yard office...	41	4	3792	True	1	Che
...	...	...	...	...	...	...	...	...	...
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	Nova

41659 rows × 11 columns

```
In [4]: df.head()
```

Out[4]:

	User ID	Username	Tweet	Retweet Count	Mention Count	Follower Count	Verified	Bot Label	Location	
1	289683	hinesstephanie	Authority research natural life material staff...	55	5	9617	True	0	Sanderston	
2	779715	roberttran	Manage whose quickly especially foot none to g...	6	2	4363	True	0	Harrisonfurt	
3	696168	pmason	Just cover eight opportunity strong policy which.	54	5	2242	True	1	Martinezberg	
4	704441	noah87	Animal sign six data good or.	26	3	8438	False	1	Camachoville	
5	570928	james00	See wonder travel this suffer less yard office...	41	4	3792	True	1	West Cheyenne	

```
In [5]: 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 [6]: `df.describe()`

Out[6]:

	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 [7]: `df.columns`

Out[7]: Index(['User ID', 'Username', 'Tweet', 'Retweet Count', 'Mention Count', 'Follower Count', 'Verified', 'Bot Label', 'Location', 'Created At', 'Hashtags'], dtype='object')

In [8]: `feature_matrix = df[['User ID', 'Retweet Count', 'Mention Count', 'Follower Count']]`  
`target_vector = df[["Verified"]]`

In [9]: `fs=StandardScaler().fit_transform(feature_matrix)`  
`logr=LogisticRegression()`  
`logr.fit(fs,target_vector)`

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples, ), for example using ravel().  
return f(\*args, \*\*kwargs)

Out[9]: LogisticRegression()

In [10]: `observation=[[1,2,3,4,5]]`

In [11]: `prediction=logr.predict(observation)`  
`print(prediction)`

[False]

In [12]: `logr.classes_`

Out[12]: array([False, True])

```
In [13]: logr.predict_proba(observation)[0][0]
```

```
Out[13]: 0.504915130281248
```

```
In [14]: logr.predict_proba(observation)[0][1]
```

```
Out[14]: 0.49508486971875193
```

## Random Forest

```
In [15]: df['Verified'].value_counts()
```

```
Out[15]: True      20845  
        False    20814  
        Name: Verified, dtype: int64
```

```
In [16]: x=df[['User ID', 'Retweet Count', 'Mention Count', 'Follower Count', 'Bot Label']]  
        y=df['Verified']
```

```
In [17]: g1={'Verified':{'True':1, "False":2}}  
df=df.replace(g1)  
df
```

Out[17]:

	User ID	Username	Tweet	Retweet Count	Mention Count	Follower Count	Verified	Bot Label	Loc
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	Camact
5	570928	james00	See wonder travel this suffer less yard office...	41	4	3792	True	1	Che
...	...	...	...	...	...	...	...	...	...
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	Nova



41659 rows × 11 columns

```
In [18]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.70)
```

```
In [19]: from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier()
rfc.fit(x_train,y_train)
```

```
Out[19]: RandomForestClassifier()
```

```
In [20]: parameters = {'max_depth':[1,2,3,4,5], 'min_samples_leaf':[5,10,15,20,25],
                        'n_estimators': [10,20,30,40,50]}
}
```

```
In [21]: from sklearn.model_selection import GridSearchCV
grid_search = GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

```
Out[21]: GridSearchCV(cv=2, estimator=RandomForestClassifier(),
                      param_grid={'max_depth': [1, 2, 3, 4, 5],
                                   'min_samples_leaf': [5, 10, 15, 20, 25],
                                   'n_estimators': [10, 20, 30, 40, 50]},
                      scoring='accuracy')
```

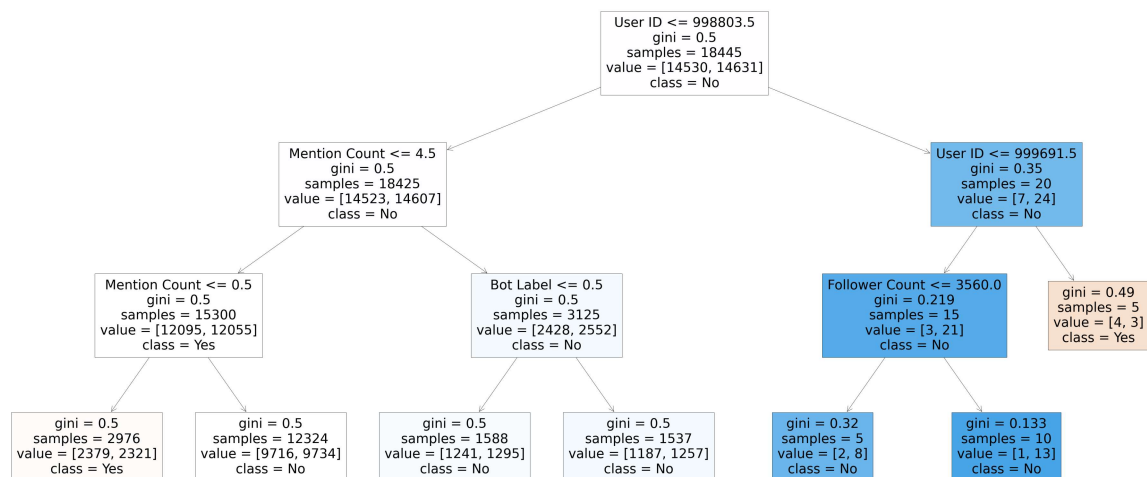
```
In [22]: grid_search.best_score_
```

```
Out[22]: 0.5045781034548127
```

```
In [23]: rfc_best = grid_search.best_estimator_
```

```
In [24]: from sklearn.tree import plot_tree
plt.figure(figsize=(89,40))
plot_tree(rfc_best.estimators_[5], feature_names=x.columns, class_names=['Yes',
```

```
Out[24]: [Text(2865.1153846153848, 1902.6000000000001, 'User ID <= 998803.5\nngini = 0.5\nnsamples = 18445\nnvalue = [14530, 14631]\nnclass = No'),
Text(1528.0615384615385, 1359.0, 'Mention Count <= 4.5\nngini = 0.5\nnsamples = 18425\nnvalue = [14523, 14607]\nnclass = No'),
Text(764.0307692307692, 815.4000000000001, 'Mention Count <= 0.5\nngini = 0.5\nnsamples = 15300\nnvalue = [12095, 12055]\nnclass = Yes'),
Text(382.0153846153846, 271.79999999999995, 'gini = 0.5\nnsamples = 2976\nnvalue = [2379, 2321]\nnclass = Yes'),
Text(1146.0461538461539, 271.79999999999995, 'gini = 0.5\nnsamples = 12324\nnvalue = [9716, 9734]\nnclass = No'),
Text(2292.0923076923077, 815.4000000000001, 'Bot Label <= 0.5\nngini = 0.5\nnsamples = 3125\nnvalue = [2428, 2552]\nnclass = No'),
Text(1910.076923076923, 271.79999999999995, 'gini = 0.5\nnsamples = 1588\nnvalue = [1241, 1295]\nnclass = No'),
Text(2674.1076923076926, 271.79999999999995, 'gini = 0.5\nnsamples = 1537\nnvalue = [1187, 1257]\nnclass = No'),
Text(4202.169230769231, 1359.0, 'User ID <= 999691.5\nngini = 0.35\nnsamples = 20\nnvalue = [7, 24]\nnclass = No'),
Text(3820.153846153846, 815.4000000000001, 'Follower Count <= 3560.0\nngini = 0.219\nnsamples = 15\nnvalue = [3, 21]\nnclass = No'),
Text(3438.1384615384613, 271.79999999999995, 'gini = 0.32\nnsamples = 5\nnvalue = [2, 8]\nnclass = No'),
Text(4202.169230769231, 271.79999999999995, 'gini = 0.133\nnsamples = 10\nnvalue = [1, 13]\nnclass = No'),
Text(4584.184615384615, 815.4000000000001, 'gini = 0.49\nnsamples = 5\nnvalue = [4, 3]\nnclass = Yes')]
```



In [ ]: