

```
In [130] import warnings
warnings.filterwarnings("ignore")
import numpy as np
import re
import pandas as pd
from tqdm import tqdm
import nltk
from nltk.corpus import stopwords
import nltk
nltk.download('punkt')
nltk.download('stopwords')
import seaborn as sns
from string import punctuation
import matplotlib.pyplot as plt
import Line
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, confusion_matrix, precision_score
from sklearn.feature_extraction.text import TfidfVectorizer

[nltk_data] Downloading package punkt to /usr/share/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /usr/share/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
In [131] Loading the preprocessed train and test data

import pickle
df_train=pickle.load(open("../input/fake-news-case-study-preprocessed-daa/f_train.pkl", 'rb'))
```

Shape of train and test dataset

In [140] df\_train.shape

Out [140] (18285, 27)

In [142] df\_train.head()

Out [142]

id	title	author	text	label	num_characters	title	num_characters	text	num_word	title	num_word	text	num_sentences	title	...	Avg_sentence_length	title	Avg_sentence_length	text	Stopword	count
0	House Dem Aide: We Don't Even See Comey's Ltr.	Danell Lucas	House Dem aide: We don't even see Comey's letter.	1	81	4930	19	943	1	...	19.0	25.48486									
1	FLYNN: Hillary Killed in Single US Airline Hit.	Daniel J. Flynn	Ever get the feeling your life codes the flu.	0	55	4160	11	822	1	...	11.0	28.344828									
2	Why The Truth You Fied	CousarLummes.com	Why the truth you fied October 29, ...	1	33	7692	7	1454	1	...	7.0	28.50904									
3	15 Civilians Killed in Single US Airline Hit.	Jessica Parkes	Videos 15 Civilians Killed in Single US Airline Hit.	1	63	3237	10	612	1	...	10.0	22.66667									
4	Iranian woman jailed for federal unpublished...	Howard Pottmoy	Iranian woman has been sentenced to...	1	93	938	14	177	1	...	14.0	35.40000									

5 rows x 27 columns

Extracting independent features on X and Class label to Y variable

In [143] X=df\_train.drop(["label","title","text","label","cleaned\_text","cleaned\_title","author"], axis=1, inplace=False)

In [144] X\_test=X.drop(["label","title","text","label","cleaned\_text","cleaned\_title","author"], axis=1, inplace=False)

Out [144] X.shape, Y.shape

Out [144] ((18285, 20), (18285, 1))

Splitting the dataset to train and test dataset

In [145] from sklearn.model\_selection import train\_test\_split

In [146] X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, Y, test\_size=0.25, stratify=y, random\_state=42)

In [146] print(X\_train.shape, y\_train.shape)

In [146] print(X\_test.shape, y\_test.shape)

Out [146] (13713, 20) (13713, 1)

In [147] X\_train.head(1)

Out [147]

num_characters	title	num_characters	text	num_word	title	num_word	text	num_sentences	title	num_sentences	text	Count	unique_words	title	Count	unique_words	text	Count	Stop_words	title	Count	Stop_words	text
8682	69	287	10	42	1	2	10	36	3	12													

Extracting numerical features from train and test dataset

In [148] X\_train\_num=X\_train.drop(["Without\_Stopwords","Without\_Stopwords\_title"], axis=1, inplace=False)

In [148] X\_test\_num=X\_test.drop(["Without\_Stopwords","Without\_Stopwords\_title"], axis=1, inplace=False)

In [148] X\_train\_num.Ft.shape, X\_test\_num.Ft.shape

Out [148] ((13713, 18), (4572, 18))

In [150] X\_train\_num.Ft.head(1)

Out [150]

num_characters	title	num_characters	text	num_word	title	num_word	text	num_sentences	title	num_sentences	text	Count	unique_words	title	Count	unique_words	text	Count	Stop_words	title	Count	Stop_words	text
8682	69	287	10	42	1	2	10	36	3	12													

In [151] X\_test\_num.Ft.head(1)

Out [151]

num_characters	title	num_characters	text	num_word	title	num_word	text	num_sentences	title	num_sentences	text	Count	unique_words	title	Count	unique_words	text	Count	Stop_words	title	Count	Stop_words	text
9013	48	2475	7	477	1	13	7	189	1	155													

Standardization of 18 numerical Manually engineered features

In [152] from sklearn.preprocessing import StandardScaler

In [152] Std\_Scaler=StandardScaler()

In [152] X\_train\_Num\_Std\_Scaler.fit\_transform(X\_train\_Num\_Ft)

In [152] X\_test\_Num\_Std\_Scaler.transform(X\_test\_Num\_Ft)

Out [152] X\_train\_Num\_Std\_Scaler.shape, X\_test\_Num\_Std\_Scaler.shape

Out [152] ((13713, 18), (4572, 18))

In [154] X\_train\_Num\_Std

Out [154]

array([[ -1.88026892, -0.88659202, -0.88363495, ..., -0.85191919,  
 0.58227813, 2.78438869],  
 [ 0.32515972, -0.11487707, 0.98084975, ..., 0.35881899,  
 0.58227813, -0.48232452],  
 [ 0.28684274, -0.12338293, -0.28214085, ..., 0.31762932,  
 0.58227813, -0.8891292 ],  
 ...,  
 [ -0.78427864, -0.88934294, 0.56651982, ..., 1.89181384,  
 0.58227813, 2.91913598],  
 [ 0.36047823, 0.60802649, 0.35935485, ..., 0.16318966,  
 -0.81898929, -0.63861608 ],  
 [ -0.89991229, -0.35221408, 0.25935485, ..., 0.00753571,  
 -0.58227813, -0.67937207]])

In [155] X\_test\_Num\_Std

Out [155]

array([[ -1.18926892, -0.4458821 , -1.58512985, ..., -0.18457833,  
 0.58227813, -0.53830701 ],  
 [ 0.59087866, -0.25689895, -0.28214085, ..., 0.2849075 ,  
 0.58227813, -0.68872293],  
 [ 0.81525893, 0.35261807, 0.15218988, ..., -0.31339738,  
 0.69523881, 0.80919935],  
 ...,  
 [ -0.1375097 , 0.87457695, 0.15218988, ..., -0.4609486 ,  
 0.58227813, -0.61152672],  
 [ 0.85952879, -0.88951242, 0.15218988, ..., -0.50539638,  
 0.58227813, 0.83779471],  
 [ -0.64892705, -0.47687348, -1.09076992, ..., 0.65462768,  
 0.58227813, 0.22454751]])

1. Traning Naive bayes model

Vectorizing our text features using tfidf vectorizer

In [156] vectorizer=tfidf.TfidfVectorizer(max\_features=3580)

In [156] vectorizer.fit(X\_train.drop(["Without\_Stopwords","Without\_Stopwords\_title"], values))

In [156] # we use the fitted TfidfVectorizer to convert the text to vector

In [156] X\_train\_title\_tfidf = vectorizer.fit\_transform(X\_train.drop(["Without\_Stopwords","Without\_Stopwords\_title"], toarray()))

In [156] X\_test\_title\_tfidf = vectorizer.fit\_transform(X\_test.drop(["Without\_Stopwords","Without\_Stopwords\_title"], toarray()))

In [156] print("After vectorizations shape of train and test data")

In [156] print(X\_train\_title\_tfidf.shape, y\_train.shape)

In [156] print(X\_test\_title\_tfidf.shape, y\_test.shape)

In [156] print("-->186")

In [156] After vectorizations shape of train and test data

In [156] (13713, 3580) (13713, 1)

In [156] (4572, 3580) (4572, 1)

In [157] vectorizer\_text=tfidf.TfidfVectorizer(max\_features=4580)

In [157] vectorizer\_text.fit(X\_train.drop(["Without\_Stopwords","Without\_Stopwords\_title"], values))

In [157] # we use the fitted CountVectorizer to convert the text to vector

In [157] X\_train\_text\_tfidf = vectorizer\_text.fit\_transform(X\_train.drop(["Without\_Stopwords","Without\_Stopwords\_title"], values).toarray()))

In [157] X\_test\_text\_tfidf = vectorizer\_text.fit\_transform(X\_test.drop(["Without\_Stopwords","Without\_Stopwords\_title"], values).toarray()))

In [157] print("After vectorizations shape of train and test data")

In [157] print(X\_train\_text\_tfidf.shape, y\_train.shape)

In [157] print(X\_test\_text\_tfidf.shape, y\_test.shape)

In [157] print("-->186")

In [157] After vectorizations shape of train and test data

In [157] (13713, 4580) (13713, 1)

In [157] (4572, 4580) (4572, 1)

In [158] # stacking all the features for train and test dataset

In [158] X\_train\_Final\_tfidf = np.hstack((X\_train\_title\_tfidf, X\_train\_text\_tfidf, X\_train\_Num\_Std))

In [158] X\_test\_Final\_tfidf = np.hstack((X\_test\_title\_tfidf, X\_test\_text\_tfidf, X\_test\_Num\_Std))

In [158] print("Final Data matrix")

In [158] print(X\_train\_Final\_tfidf.shape, y\_train.shape)

In [158] print(X\_test\_Final\_tfidf.shape, y\_test.shape)

In [158] print("-->100")

In [158] Final Data matrix

In [158] (13713, 8688) (13713, 1)

In [158] (4572, 8688) (4572, 1)

Training the gaussian NB model

In [159] gnb = GaussianNB()

In [159] gnb.fit(X\_train\_Final\_tfidf, y\_train)

In [159] y\_pred = gnb.predict(X\_test\_Final\_tfidf)

In [159] accuracy = accuracy\_score(y\_test, y\_pred)

In [159] precision = precision\_score(y\_test, y\_pred)

In [159] print("Accuracy of Gaussian NB accuracy on test dataset :", accuracy)

In [159] print("Gaussian NB precision on test dataset :", precision)

In [159] Gaussian NB accuracy on test dataset : 0.8258967829846369

In [159] Gaussian nb precision on test dataset : 0.7718643741403028

Plotting confusion matrix

In [160] from sklearn.metrics import confusion\_matrix

In [160] #Generate the confusion matrix

In [160] cf\_matrix = confusion\_matrix(y\_test, y\_pred)

In [160] print(cf\_matrix)

Out [160] [[2093 498]  
 [ 298 1483]]

In [161] group\_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']

In [161] group\_counts = [{"0:0.07".format(value) for value in cf\_matrix.flatten()}]

In [161] group\_percentages = [{"0:.25)".format(value) for value in cf\_matrix.flatten()/(np.sum(cf\_matrix))}]

In [161] labels = [{"(v1)\n(v2)\n(v3)" for v1, v2, v3 in zip(group\_names, group\_counts, group\_percentages)}]

In [161] labels = np.asarray(labels).reshape(2,2)

In [161] ax = sns.heatmap(cf\_matrix, annot=labels, fctc='', cmap='Blues')

In [161] ax.set\_title('Seaborn Confusion Matrix with labels\n\n');

In [161] ax.set\_xlabel('\\nPredicted Values');

In [161] ax.set\_ylabel('Actual Values');

In [161] ax.xaxis.set\_ticklabels(['False', 'True'])

In [161] ax.yaxis.set\_ticklabels(['False', 'True'])

In [161] #w Display the visualization of the Confusion Matrix.

In [161] plt.show()

Seaborn Confusion Matrix with labels

k-nearest neighbors Classifier

In [162] from sklearn.neighbors import KNeighborsClassifier

In [162] neigh = KNeighborsClassifier()

In [162] neigh.fit(X\_train\_Final\_tfidf, y\_train)

In [162] y\_pred = neigh.predict(X\_test\_Final\_tfidf)

In [162] accuracy = accuracy\_score(y\_test, y\_pred)

In [162] precision = precision\_score(y\_test, y\_pred)

In [162] print("k-nearest neighbors accuracy on testdataset :", accuracy)

In [162] print("k-nearest neighbors precision on testdataset :", precision)

In [162] k-nearest neighbors accuracy on testdataset : 0.8396782984636921

In [162] k-nearest neighbors precision on testdataset : 0.8679245283018868

Plotting confusion matrix

In [163] from sklearn.metrics import confusion\_matrix

In [163] #Generate the confusion matrix

In [163] cf\_matrix = confusion\_matrix(y\_test, y\_pred)

In [163] print(cf\_matrix)

Out [163] [[2367 2241]  
 [ 589 1472]]

In [164] group\_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']

In [164] group\_counts = [{"0:0.07".format(value) for value in cf\_matrix.flatten()}]

In [164] group\_percentages = [{"0:.25)".format(value) for value in cf\_matrix.flatten()/(np.sum(cf\_matrix))}]

In [164] labels = [{"(v1)\n(v2)\n(v3)" for v1, v2, v3 in zip(group\_names, group\_counts, group\_percentages)}]

In [164] labels = np.asarray(labels).reshape(2,2)

In [164] ax = sns.heatmap(cf\_matrix, annot=labels, fctc='', cmap='Blues')

In [164] ax.set\_title('Seaborn Confusion Matrix with labels\n\n');

In [164] ax.set\_xlabel('\\nPredicted Values');

In [164] ax.set\_ylabel('Actual Values');

In [164] ax.xaxis.set\_ticklabels(['False', 'True'])

In [164] ax.yaxis.set\_ticklabels(['False', 'True'])

In [164] #w Display the visualization of the Confusion Matrix.

In [164] plt.show()

Seaborn Confusion Matrix with labels

Traning KNN on only 18 manually engineered Numerical features

In [165] from sklearn.neighbors import KNeighborsClassifier

In [165] neigh\_num = KNeighborsClassifier()

In [165] neigh\_num.fit(X\_train\_Num\_Std, y\_train)

In [165] y\_pred\_num = neigh\_num.predict(X\_test\_Num\_Std)

In [165] accuracy = accuracy\_score(y\_test, y\_pred\_num)

In [165] precision = precision\_score(y\_test, y\_pred\_num)

In [165] print("Accuracy of KNN using only 18 manually engineered numerical features")

In [165] print("-->186")

In [165] print("k-nearest neighbors accuracy on testdataset :", accuracy)

In [165] print("k-nearest neighbors precision on testdataset :", precision)

In [165] Accuracy of KNN using only 18 manually engineered numerical features

In [165] k-nearest neighbors accuracy on testdataset : 0.8112423447869316

In [165] k-nearest neighbors precision on testdataset : 0.828825294517847

Observations:-

- Above we have only used manually engineered 18 features to train KNN model and it is performing decently and here we are getting accuracy on test data is 81%.

Plotting confusion matrix

In [166] from sklearn.metrics import confusion\_matrix

In [166] #Generate the confusion matrix

In [166] cf\_matrix = confusion\_matrix(y\_test, y\_pred\_num)

In [166] print(cf\_matrix)

Out [166] [[2360 2911]  
 [ 572 1469]]

In [167] group\_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']

In [167] group\_counts = [{"0:0.07".format(value) for value in cf\_matrix.flatten()}]

In [167] group\_percentages = [{"0:.25)".format(value) for value in cf\_matrix.flatten()/(np.sum(cf\_matrix))}]

In [167] labels = [{"(v1)\n(v2)\n(v3)" for v1, v2, v3 in zip(group\_names, group\_counts, group\_percentages)}]

In [167] labels = np.asarray(labels).reshape(2,2)

In [167] ax = sns.heatmap(cf\_matrix, annot=labels, fctc='', cmap='Blues')

In [167] ax.set\_title('Seaborn Confusion Matrix with labels\n\n');

In [167] ax.set\_xlabel('\\nPredicted Values');

In [167] ax.set\_ylabel('Actual Values');

In [167] ax.xaxis.set\_ticklabels(['False', 'True'])

In [167] ax.yaxis.set\_ticklabels(['False', 'True'])

In [167] #w Display the visualization of the Confusion Matrix.

In [167] plt.show()

Seaborn Confusion Matrix with labels

Observations:-

- from above we can see that after adding our 18 numerical feature we are able to get 81% accuracy on test data that show that our manually engineered feature are able to separate both classes.

3. Traning Logistic Regression model

In [168] from sklearn.linear\_model import LogisticRegression

In [168] LR = LogisticRegression(random\_state=12)

In [168] LR.fit(X\_train\_Final\_tfidf, y\_train)

In [168] y\_pred\_lr = LR.predict(X\_test\_Final\_tfidf)

In [168] accuracy = accuracy\_score(y\_test, y\_pred\_lr)

In [168] precision = precision\_score(y\_test, y\_pred\_lr)

In [168] print("LR accuracy on testdataset :", accuracy)

In [168] print("LR precision on testdataset :", precision)

In [168] LR accuracy on testdataset : 0.878586675415573

In [168] LR precision on testdataset : 0.9593137254982061

Plotting confusion matrix

In [169] from sklearn.metrics import confusion\_matrix

In [169] #Generate the confusion matrix

In [169] cf\_matrix = confusion\_matrix(y\_test, y\_pred\_lr)

In [169] print(cf\_matrix)

Out [169] [[2568 83]  
 [ 24 1957]]

In [170] group\_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']

In [170] group\_counts = [{"0:0.07".format(value) for value in cf\_matrix.flatten()}]

In [170] group\_percentages = [{"0:.25)".format(value) for value in cf\_matrix.flatten()/(np.sum(cf\_matrix))}]

In [170] labels = [{"(v1)\n(v2)\n(v3)" for v1, v2, v3 in zip(group\_names, group\_counts, group\_percentages)}]

In [170] labels = np.asarray(labels).reshape(2,2)

In [170] ax = sns.heatmap(cf\_matrix, annot=labels, fctc='', cmap='Blues')

In [170] ax.set\_title('Seaborn Confusion Matrix with labels\n\n');

In [170] ax.set\_xlabel('\\nPredicted Values');

In [170] ax.set\_ylabel('Actual Values');

In [170] ax.xaxis.set\_ticklabels(['False', 'True'])

In [170] ax.yaxis.set\_ticklabels(['False', 'True'])

In [170] #w Display the visualization of the Confusion Matrix.

In [170] plt.show()

Seaborn Confusion Matrix with labels

Traning Logistic regression only on numerical feaures

In [171] LR = LogisticRegression(random\_state=12)

In [171] LR.fit(X\_train\_Num\_Std, y\_train)

In [171] y\_pred\_lr = LR.predict(X\_test\_Num\_Std)

In [171] accuracy = accuracy\_score(y\_test, y\_pred\_lr)

In [171] precision = precision\_score(y\_test, y\_pred\_lr)

In [171] print("LR accuracy on testdataset :", accuracy)

In [171] print("LR precision on testdataset :", precision)

In [171] LR accuracy on testdataset : 0.7569991251893613

In [171] LR precision on testdataset : 0.7445985401459854

Observations:-

- Here again we used only 18 numerical feaures to train Logistic regression and it is not performing that well as compare to KNN and here we are getting Accuracy on test dataset is 75%

Plotting confusion matrix

In [172] from sklearn.metrics import confusion\_matrix

In [172] #Generate the confusion matrix

In [172] cf\_matrix = confusion\_matrix(y\_test, y\_pred\_lr)

In [172] print(cf\_matrix)

Out [172] [[2284 387]  
 [ 724 1257]]

In [173] group\_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']

In [173] group\_counts = [{"0:0.07".format(value) for value in cf\_matrix.flatten()}]

In [173] group\_percentages = [{"0:.25)".format(value) for value in cf\_matrix.flatten()/(np.sum(cf\_matrix))}]

In [173] labels = [{"(v1)\n(v2)\n(v3)" for v1, v2, v3 in zip(group\_names, group\_counts, group\_percentages)}]

In [173] labels = np.asarray(labels).reshape(2,2)

In [173] ax = sns.heatmap(cf\_matrix, annot=labels, fctc='', cmap='Blues')

In [173] ax.set\_title('Seaborn Confusion Matrix with labels\n\n');

In [173] ax.set\_xlabel('\\nPredicted Values');

In [173] ax.set\_ylabel('Actual Values');

In [173] ax.xaxis.set\_ticklabels(['False', 'True'])

In [173] ax.yaxis.set\_ticklabels(['False', 'True'])

In [173] #w Display the visualization of the Confusion Matrix.

In [173] plt.show()

Seaborn Confusion Matrix with labels

Summary of above trained models on both text as well as 18 numerical features:-

In [174] from prettytable import PrettyTable

In [174] # Specify the Column Names while initializing the Table

In [174] myTable = PrettyTable(["Vectorizer", "Model", "Accuracy On Test Data", "Precision on Test Data", "Precision on Test Data"])

In [174] # Add rows

In [174] myTable.add\_row(["Tfidf Vectorizer", "Gaussian NB", "82.5%", "77.1%"])

In [174] myTable.add\_row(["Tfidf Vectorizer", "KNN", "81.7%", "75.2%"])

In [174] myTable.add\_row(["Tfidf Vectorizer", "Logistic Regression", "75.6%", "74.6%"])

In [174] myTable.add\_row(["Tfidf Vectorizer", "Logistic Regression", "79.3%", "79.9%"])

In [174] print(myTable)

-----

	Vectorizer	Model	Accuracy On Test Data	Precision on Test Data	Precision on Test Data
	Tfidf Vectorizer	Gaussian NB	82.5	77.1	77.1
	Tfidf Vectorizer	Decision Tree	82.9	86.7	86.7
	Tfidf Vectorizer	Logistic Regression	75.6	79.9	79.9

-----

Observations:-

- From above table we can see that after adding our 18 numerical features our accuracy and precision both got increased for all the three models
- For KNN the results are very impressive as initially with only text features we were getting accuracy of 46.5% and precision of 44.7% but after adding our numerical features we are getting accuracy of 83.9% and precision of 86.7%
- Simlary we can some improve of accuracy and precision in Logistic rgression model as well

Below is table when all three model is trained only with text features using tfidf vectorization we can clearly see after adding our 18 manually engineered numerical features our models performance got improved specially KNN

In [175] from prettytable import PrettyTable

In [175] # Specify the Column Names while initializing the Table

In [175] myTable = PrettyTable(["Vectorizer", "Model", "Accuracy On Test Data", "Precision on Test Data", "Precision on Test Data"])

In [175] # Add rows

In [175] myTable.add\_row(["Tfidf Vectorizer", "Gaussian NB", "81.7%", "75.2%"])

In [175] myTable.add\_row(["Tfidf Vectorizer", "KNN", "46.5%", "44.7%"])

In [175] myTable.add\_row(["Tfidf Vectorizer", "Logistic Regression", "79.3%", "79.2%"])

In [175] print(myTable)

-----

	Vectorizer	Model	Accuracy On Test Data	Precision on Test Data	Precision on Test Data
	Tfidf Vectorizer	Gaussian NB	81.7	75.2	75.2
	Tfidf Vectorizer	KNN	46.5	44.7	44.7
	Tfidf Vectorizer	Logistic Regression	79.3	79.2	79.2

-----