

## Fake News Classifier

Dataset: <https://www.kaggle.com/c/fake-news/data#> (<https://www.kaggle.com/c/fake-news/data>)

In [1]: 1 `import pandas as pd`

In [2]: 1 `df=pd.read_csv("D:/fakenews_data/train.csv")`

In [3]: 1 `df.head()`

Out[3]:

	id	title	author	text	label
0	0	House Dem Aide: We Didn't Even See Comey's Let...	Darrell Lucas	House Dem Aide: We Didn't Even See Comey's Let...	1
1	1	FLYNN: Hillary Clinton, Big Woman on Campus - ...	Daniel J. Flynn	Ever get the feeling your life circles the rou...	0
2	2	Why the Truth Might Get You Fired	Consortiumnews.com	Why the Truth Might Get You Fired October 29, ...	1
3	3	15 Civilians Killed In Single US Airstrike Hav...	Jessica Purkiss	Videos 15 Civilians Killed In Single US Aistr...	1
4	4	Iranian woman jailed for fictional unpublished...	Howard Portnoy	Print lnAn Iranian woman has been sentenced to...	1

In [4]: 1 `## Get the Independent Features`  
2 `X=df.drop('label',axis=1)`

In [5]: 1 `X.head()`

Out[5]:

	id	title	author	text
0	0	House Dem Aide: We Didn't Even See Comey's Let...	Darrell Lucas	House Dem Aide: We Didn't Even See Comey's Let...
1	1	FLYNN: Hillary Clinton, Big Woman on Campus - ...	Daniel J. Flynn	Ever get the feeling your life circles the rou...
2	2	Why the Truth Might Get You Fired	Consortiumnews.com	Why the Truth Might Get You Fired October 29, ...
3	3	15 Civilians Killed In Single US Airstrike Hav...	Jessica Purkiss	Videos 15 Civilians Killed In Single US Aistr...
4	4	Iranian woman jailed for fictional unpublished...	Howard Portnoy	Print lnAn Iranian woman has been sentenced to...

```
In [6]: 1  ## Get the Dependent features
        2  y=df['label']
```

```
In [7]: 1  y.head()
```

```
Out[7]: 0    1
        1    0
        2    1
        3    1
        4    1
        Name: label, dtype: int64
```

```
In [8]: 1  df.shape
```

```
Out[8]: (20800, 5)
```

```
In [9]: 1  from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer, HashingVectorizer
```

```
In [10]: 1  df=df.dropna()
```

In [11]: 1 df.head(10)

Out[11]:

	id	title	author	text	label
0	0	House Dem Aide: We Didn't Even See Comey's Let...	Darrell Lucus	House Dem Aide: We Didn't Even See Comey's Let...	1
1	1	FLYNN: Hillary Clinton, Big Woman on Campus - ...	Daniel J. Flynn	Ever get the feeling your life circles the rou...	0
2	2	Why the Truth Might Get You Fired	Consortiumnews.com	Why the Truth Might Get You Fired October 29, ...	1
3	3	15 Civilians Killed In Single US Airstrike Hav...	Jessica Purkiss	Videos 15 Civilians Killed In Single US Aistr...	1
4	4	Iranian woman jailed for fictional unpublished...	Howard Portnoy	Print \nAn Iranian woman has been sentenced to...	1
5	5	Jackie Mason: Hollywood Would Love Trump if He...	Daniel Nussbaum	In these trying times, Jackie Mason is the Voi...	0
7	7	Benoît Hamon Wins French Socialist Party's Pre...	Alissa J. Rubin	PARIS — France chose an idealistic, traditi...	0
9	9	A Back-Channel Plan for Ukraine and Russia, Co...	Megan Twohey and Scott Shane	A week before Michael T. Flynn resigned as nat...	0
10	10	Obama's Organizing for Action Partners with So...	Aaron Klein	Organizing for Action, the activist group that...	0
11	11	BBC Comedy Sketch "Real Housewives of ISIS" Ca...	Chris Tomlinson	The BBC produced spoof on the "Real Housewives...	0

In [12]: 1 messages=df.copy()

In [13]: 1 messages.reset\_index(inplace=True)

In [14]:

```
1 messages.head(10)
```

Out[14]:

	index	id	title	author	text	label
0	0	0	House Dem Aide: We Didn't Even See Comey's Let...	Darrell Lucus	House Dem Aide: We Didn't Even See Comey's Let...	1
1	1	1	FLYNN: Hillary Clinton, Big Woman on Campus - ...	Daniel J. Flynn	Ever get the feeling your life circles the rou...	0
2	2	2	Why the Truth Might Get You Fired	Consortiumnews.com	Why the Truth Might Get You Fired October 29, ...	1
3	3	3	15 Civilians Killed In Single US Airstrike Hav...	Jessica Purkiss	Videos 15 Civilians Killed In Single US Aistr...	1
4	4	4	Iranian woman jailed for fictional unpublished...	Howard Portnoy	Print \nAn Iranian woman has been sentenced to...	1
5	5	5	Jackie Mason: Hollywood Would Love Trump if He...	Daniel Nussbaum	In these trying times, Jackie Mason is the Voi...	0
6	7	7	Benoît Hamon Wins French Socialist Party's Pre...	Alissa J. Rubin	PARIS — France chose an idealistic, traditi...	0
7	9	9	A Back-Channel Plan for Ukraine and Russia, Co...	Megan Twohey and Scott Shane	A week before Michael T. Flynn resigned as nat...	0
8	10	10	Obama's Organizing for Action Partners with So...	Aaron Klein	Organizing for Action, the activist group that...	0
9	11	11	BBC Comedy Sketch "Real Housewives of ISIS" Ca...	Chris Tomlinson	The BBC produced spoof on the "Real Housewives...	0

In [15]:

```
1 messages['title'][0:5]
```

Out[15]:

```
0 House Dem Aide: We Didn't Even See Comey's Let...
1 FLYNN: Hillary Clinton, Big Woman on Campus - ...
2 Why the Truth Might Get You Fired
3 15 Civilians Killed In Single US Airstrike Hav...
4 Iranian woman jailed for fictional unpublished...
Name: title, dtype: object
```

```
In [16]: 1 from nltk.corpus import stopwords
2 from nltk.stem.porter import PorterStemmer
3 import re
4 ps = PorterStemmer()
5 corpus = []
6 for i in range(0, len(messages)):
7     review = re.sub('[^a-zA-Z]', ' ', messages['title'][i])
8     review = review.lower()
9     review = review.split()
10
11     review = [ps.stem(word) for word in review if not word in stopwords.words('english')]
12     review = ' '.join(review)
13     corpus.append(review)
```

```
In [17]: 1 corpus[0:10]
```

```
Out[17]: ['hous dem aid even see comey letter jason chaffetz tweet',
'flynn hillari clinton big woman campu breitbart',
'truth might get fire',
'civilian kill singl us airstrik identifi',
'iranian woman jail fiction unpublish stori woman stone death adulteri',
'jacki mason hollywood would love trump bomb north korea lack tran bathroom exclus video breitbart',
'beno hamon win french socialist parti presidenti nomin new york time',
'back channel plan ukrain russia courtesi trump associ new york time',
'obama organ action partner soro link indivis disrupt trump agenda',
'bbc comedi sketch real housew isi caus outrag']
```

```
In [18]: 1 ## TFidf Vectorizer
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 tfidf_v=TfidfVectorizer(max_features=5000,ngram_range=(1,3))
4 X=tfidf_v.fit_transform(corpus).toarray()
5
```

```
In [19]: 1 X.shape
```

```
Out[19]: (18285, 5000)
```

```
In [20]: 1 y=messages['label']
```

```
In [21]: 1 ## Divide the dataset into Train and Test
          2 from sklearn.model_selection import train_test_split
          3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=0)
```

```
In [22]: 1 tfidf_v.get_feature_names()[ :20]
```

```
Out[22]: ['abandon',
          'abc',
          'abc news',
          'abduct',
          'abe',
          'abedin',
          'abl',
          'abort',
          'abroad',
          'absolut',
          'abstain',
          'absurd',
          'abus',
          'abus new',
          'abus new york',
          'academi',
          'accept',
          'access',
          'access pipelin',
          'access pipelin protest']
```

```
In [23]: 1 tfidf_v.get_params()
```

```
Out[23]: {'analyzer': 'word',  
         'binary': False,  
         'decode_error': 'strict',  
         'dtype': numpy.float64,  
         'encoding': 'utf-8',  
         'input': 'content',  
         'lowercase': True,  
         'max_df': 1.0,  
         'max_features': 5000,  
         'min_df': 1,  
         'ngram_range': (1, 3),  
         'norm': 'l2',  
         'preprocessor': None,  
         'smooth_idf': True,  
         'stop_words': None,  
         'strip_accents': None,  
         'sublinear_tf': False,  
         'token_pattern': '(?u)\\b\\w+\\b',  
         'tokenizer': None,  
         'use_idf': True,  
         'vocabulary': None}
```

```
In [24]: 1 count_df = pd.DataFrame(X_train, columns=tfidf_v.get_feature_names())
```

```
In [25]: 1 count_df.head()
```

Out[25]:

	abandon	abc	abc news	abduct	abe	abedin	abl	abort	abroad	absolut	...	zero	zika	zika viru	zionist	zone	zone new	zone new york	zoo	zu	zu
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	...	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	0.0	0.0	0.0	0.0	0.0	0.000000	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.305244	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	

5 rows × 5000 columns



```
In [26]: 1 import matplotlib.pyplot as plt
```



```
In [27]: 1 def plot_confusion_matrix(cm, classes,
2                               normalize=False,
3                               title='Confusion matrix',
4                               cmap=plt.cm.Blues):
5
6     plt.imshow(cm, interpolation='nearest', cmap=cmap)
7     plt.title(title)
8     plt.colorbar()
9     tick_marks = np.arange(len(classes))
10    plt.xticks(tick_marks, classes, rotation=45)
11    plt.yticks(tick_marks, classes)
12
13    if normalize:
14        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
15        print("Normalized confusion matrix")
16    else:
17        print('Confusion matrix, without normalization')
18
19    thresh = cm.max() / 2.
20    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
21        plt.text(j, i, cm[i, j],
22                 horizontalalignment="center",
23                 color="white" if cm[i, j] > thresh else "black")
24
25    plt.tight_layout()
26    plt.ylabel('True label')
27    plt.xlabel('Predicted label')
```

## MultinomialNB Algorithm

```
In [28]: 1 from sklearn.naive_bayes import MultinomialNB
2 classifier=MultinomialNB()
```

```
In [29]: 1 from sklearn import metrics
2 import numpy as np
3 import itertools
```

```
In [30]: 1 classifier.fit(X_train, y_train)
2         pred = classifier.predict(X_test)
3         score = metrics.accuracy_score(y_test, pred)
4         print("accuracy:  %0.3f" % (score*100))
5         cm = metrics.confusion_matrix(y_test, pred)
6         plot_confusion_matrix(cm, classes=['FAKE', 'REAL'])
```

accuracy: 88.020

Confusion matrix, without normalization

