

fake-junu

June 27, 2024

```
[37]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split as ttp
from sklearn.metrics import classification_report
import re
import string
import matplotlib.pyplot as plt
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression, LinearRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
```

```
[34]: data_true = pd.read_csv('/content/drive/MyDrive/True.csv')
data_fake = pd.read_csv('/content/drive/MyDrive/Fake.csv')
```

```
[35]: data_true['label'] = 1
data_fake['label'] = 0
```

```
[36]: data_true_manual_testing = data_true.tail(10)
data_true = data_true.iloc[:-10]
```

```
[38]: data_fake_manual_testing = data_fake.tail(10)
data_fake = data_fake.iloc[:-10]
```

```
[39]: data_manual_testing = pd.concat([data_true_manual_testing,
↪data_fake_manual_testing], axis=0)
data_manual_testing.to_csv('manual_testing.csv', index=False)
```

```
[40]: data_merge = pd.concat([data_true, data_fake], axis=0)
```

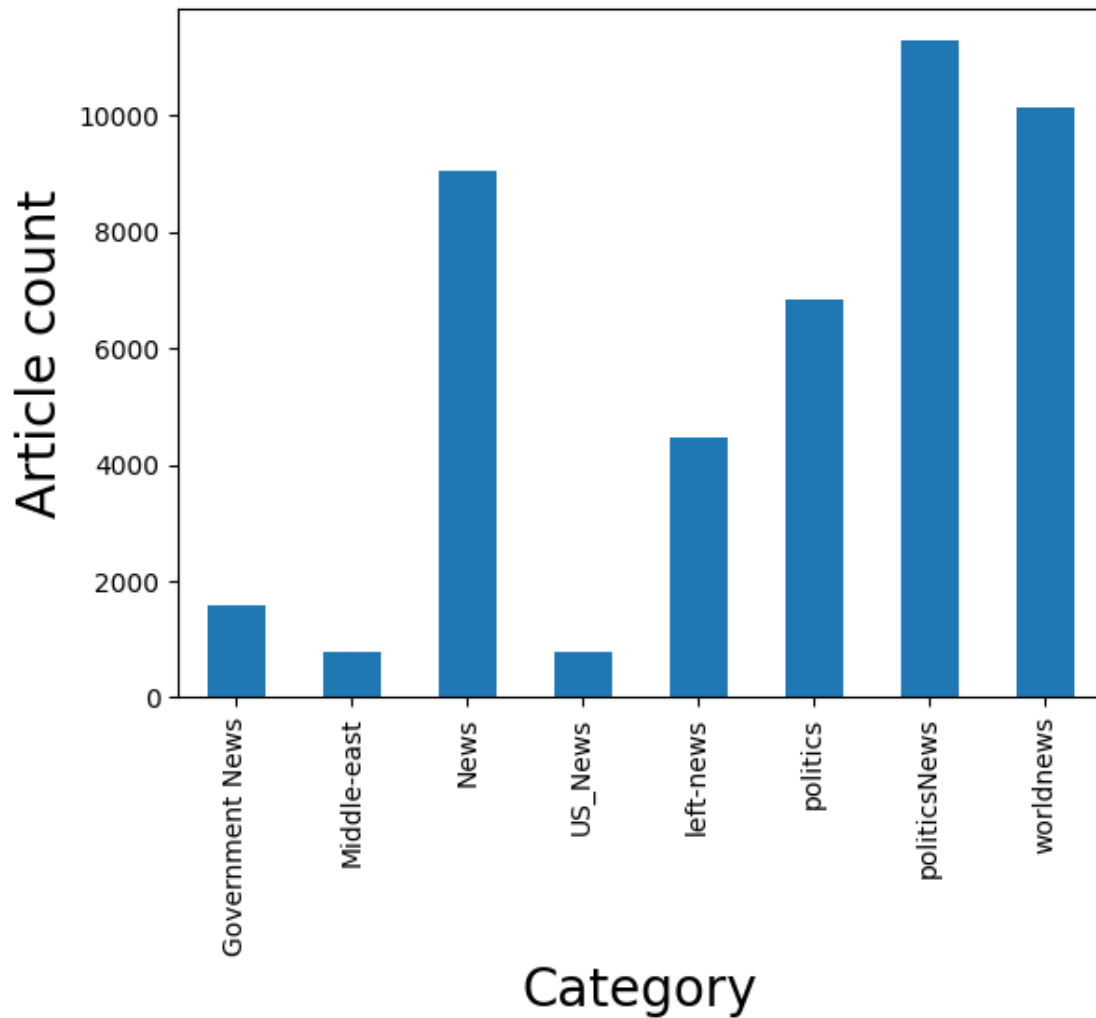
```
[41]: data_true.shape, data_fake.shape
```

```
[41]: ((21407, 5), (23471, 5))
```

```
[47]: print(data_merge.groupby(['subject'])['text'].count())
data_merge.groupby(['subject'])['text'].count().plot(kind='bar')
```

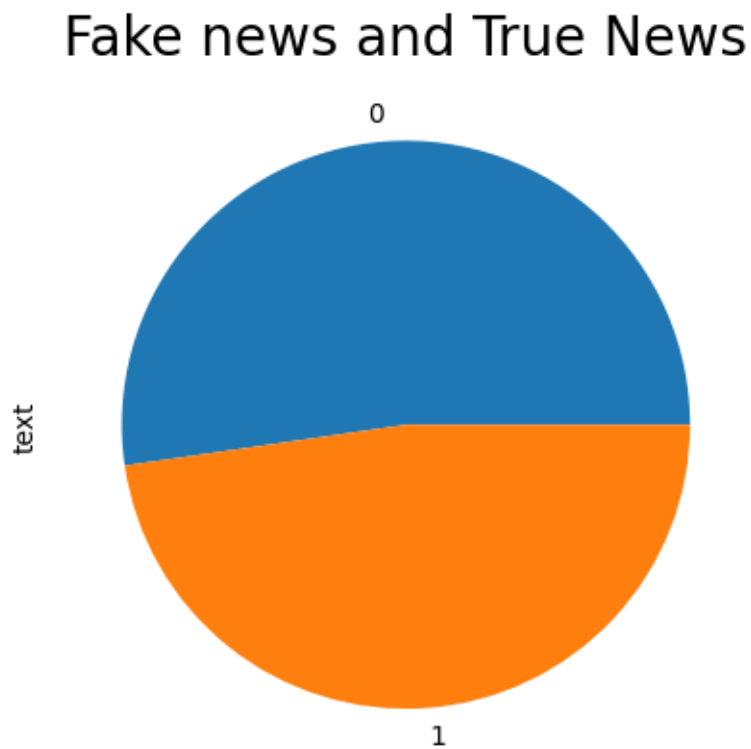
```
plt.xlabel("Category", size=20)
plt.ylabel("Article count", size=20)
plt.show()
```

```
subject
Government News    1570
Middle-east        768
News               9050
US_News            783
left-news          4459
politics           6841
politicsNews       11272
worldnews          10135
Name: text, dtype: int64
```



```
[49]: print(data_merge.groupby(['label'])['text'].count()) # Assuming 'label' is the
      ↪ correct column name
print("0 = fake news\n1 = true news")
data_merge.groupby(['label'])['text'].count().plot(kind='pie')
plt.title("Fake news and True News",size=20)
plt.show()
```

```
label
0    23471
1    21407
Name: text, dtype: int64
0 = fake news
1 = true news
```



```
[11]: data = data_merge.drop(['title', 'subject', 'date'], axis=1)
      data = data.sample(frac=1).reset_index(drop=True)

      print(data.isnull().sum())
```

```
text    0
label    0
dtype: int64
```

```
[50]: data = data.sample(frac=1)
data.head(10)
```

```
[50]:
```

	text	label
15089	president trump addressed the america people t...	0
22722	washington reuters u s senator john mccai...	1
24020	new york reuters u s ambassador to the un...	1
20369	what a crazy group of professional agitators a...	0
36705	seoul reuters china has resumed its ban on...	1
36900	beirut reuters lebanon s prime minister sa...	1
4949	hysterical as always conservative comedian s...	0
6617	mexico city reuters seven latin american g...	1
38954	fbi insiders are spreading the word hillary c...	0
38323	donald trump s past support of hillary clinton...	0

```
[13]: def filtering(text):
    text = text.lower()
    text = re.sub('[.*?\\]', ' ', text)
    text = re.sub("\\W", " ", text)
    text = re.sub('https?://\\S+|www\\.\\S+', ' ', text)
    text = re.sub('<.*?>+', ' ', text)
    text = re.sub('[%s]' % re.escape(string.punctuation), ' ', text)
    text = re.sub('\\n', ' ', text)
    text = re.sub('\\w*\\d\\w*', ' ', text)
    return text

data['text'] = data['text'].apply(filtering)
```

```
[51]: data['text']=data['text'].apply(filtering)
data.head(10)
```

```
[51]:
```

	text	label
15089	president trump addressed the america people t...	0
22722	washington reuters u s senator john mccai...	1
24020	new york reuters u s ambassador to the un...	1
20369	what a crazy group of professional agitators a...	0
36705	seoul reuters china has resumed its ban on...	1
36900	beirut reuters lebanon s prime minister sa...	1
4949	hysterical as always conservative comedian s...	0
6617	mexico city reuters seven latin american g...	1
38954	fbi insiders are spreading the word hillary c...	0
38323	donald trump s past support of hillary clinton...	0

```
[14]: vectorizer = TfidfVectorizer()
x = vectorizer.fit_transform(data['text'])
y = data['label']
```

```
x_train, x_test, y_train, y_test = ttp(x, y, test_size=0.25, random_state=0)
```

```
[16]: LR = LogisticRegression(max_iter=1000)
      LR.fit(x_train, y_train)
```

```
[16]: LogisticRegression(max_iter=1000)
```

```
[17]: DT = DecisionTreeClassifier()
      DT.fit(x_train, y_train)
```

```
[17]: DecisionTreeClassifier()
```

```
[18]: RF = RandomForestClassifier()
      RF.fit(x_train, y_train)
```

```
[18]: RandomForestClassifier()
```

```
[21]: LinR = LinearRegression()
      LinR.fit(x_train, y_train)
```

```
[21]: LinearRegression()
```

```
[22]: KNN = KNeighborsClassifier()
      KNN.fit(x_train, y_train)
```

```
[22]: KNeighborsClassifier()
```

```
[28]: user_input = input("Enter the news article text: ")
      user_input_transformed = vectorizer.transform([user_input])
```

Enter the news article text: MODI IS THE NOT THE PM OF THE COUNTRY

```
[29]: prediction_LR = LR.predict(user_input_transformed)
      if prediction_LR == 1:
          print("Logistic Regression: This news article is real.")
      else:
          print("Logistic Regression: This news article is fake.")

      prediction_DT = DT.predict(user_input_transformed)
      if prediction_DT == 1:
          print("Decision Tree: This news article is real.")
      else:
          print("Decision Tree: This news article is fake.")

      prediction_RF = RF.predict(user_input_transformed)
      if prediction_RF == 1:
          print("Random Forest: This news article is real.")
```

```

else:
    print("Random Forest: This news article is fake.")

prediction_LinR = LinR.predict(user_input_transformed)

prediction_LinR = np.where(prediction_LinR >= 0.5, 1, 0)
if prediction_LinR == 1:
    print("Linear Regression: This news article is real.")
else:
    print("Linear Regression: This news article is fake.")

prediction_KNN = KNN.predict(user_input_transformed)
if prediction_KNN == 1:
    print("K-Nearest Neighbors: This news article is real.")
else:
    print("K-Nearest Neighbors: This news article is fake.")

```

Logistic Regression: This news article is fake.

Decision Tree: This news article is fake.

Random Forest: This news article is fake.

Linear Regression: This news article is fake.

K-Nearest Neighbors: This news article is fake.