

MKUGANKUMAR /  
Fack\_news\_detection\_using\_NLP\_phase5

&lt;&gt; Code

Issues

Pull requests

Projects

Wiki

Security

Insights

Settings

Fack\_news\_detection\_using\_NLP\_phase5 / main.py



MKUGANKUMAR Update main.py

4 minutes ago



226 lines (178 loc) · 13.5 KB

Code

Blame

Raw



```
1  # %% [markdown]
2  # *Fake News Detection*
3  #
4  # It has become humanly impossible to identify fake news on the online portals across the g
5  #
6  # The most crucial thing here is data which has been already available in the kaggle. We wi
7
8  # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:19:57.550035Z","iopub.execute_
9  # This Python 3 environment comes with many helpful analytics libraries installed
10 # It is defined by the kaggle/python Docker image: https://github.com/kaggle/docker-python
11 # For example, here's several helpful packages to load
12
13
14
15 import numpy as np # linear algebra
16 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
17
18 # Input data files are available in the read-only "../input/" directory
19 # For example, running this (by clicking run or pressing Shift+Enter) will list all files u
20
21 import os
22 for dirname, __, filenames in os.walk('/kaggle/input'):
23     for filename in filenames:
24         print(os.path.join(dirname, filename))
25
26 # You can write up to 20GB to the current directory (/kaggle/working/) that gets preserved
27 # You can also write temporary files to /kaggle/temp/, but they won't be saved outside of t
28
29 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:19:57.573835Z","iopub.execute_
30 !pip install gensim # Gensim is an open-source library for unsupervised topic modeling and
31 import nltk
32 nltk.download('punkt')
```

```
34 import pandas as pd
35 import numpy as np
36 import matplotlib.pyplot as plt
37 import seaborn as sns
38 from wordcloud import WordCloud, STOPWORDS
39 import nltk
40 import re
41 from nltk.corpus import stopwords
42 import seaborn as sns
43 import gensim
44 from gensim.utils import simple_preprocess
45 from gensim.parsing.preprocessing import STOPWORDS
46
47 import plotly.express as px
48 from sklearn.model_selection import train_test_split
49 from sklearn.feature_extraction.text import CountVectorizer
50 from sklearn.linear_model import LogisticRegression
51 from sklearn.metrics import roc_auc_score
52 from sklearn.metrics import confusion_matrix
53
54 # %% [markdown]
55 # *Import the data & Clean ups*
56
57 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:12.47781Z","iopub.execute_i
58 #importing data
59 fake_data = pd.read_csv('/kaggle/input/fake-and-real-news-dataset/Fake.csv')
60 print("fake_data", fake_data.shape)
61
62 true_data= pd.read_csv('/kaggle/input/fake-and-real-news-dataset/True.csv')
63 print("true_data", true_data.shape)
64
65 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.043858Z","iopub.execute_
66 fake_data.head(5)
67
68 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.066765Z","iopub.execute_
69 true_data.head(5)
70
71 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.0866Z","iopub.execute_in
72 #adding additional column to separate between true & fake data
73 # true =1, fake =0
74 true_data['target'] = 1
75 fake_data['target'] = 0
76 df = pd.concat([true_data, fake_data]).reset_index(drop = True)
77 df['original'] = df['title'] + ' ' + df['text']
78 df.head()
79
80 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.737583Z","iopub.execute_
81 df.isnull().sum()
82
```

```

83 # %% [markdown]
84 # *Data Clean up*
85 # - create a function here that will be responsible to remove any unnecessary words (Stopwo
86
87 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.781215Z","iopub.execute_
88 stop_words = stopwords.words('english')
89 stop_words.extend(['from', 'subject', 're', 'edu', 'use'])
90 ✓ def preprocess(text):
91     result = []
92     for token in gensim.utils.simple_preprocess(text):
93         if token not in gensim.parsing.preprocessing.STOPWORDS and len(token) > 2 and token
94             result.append(token)
95
96     return result
97
98 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.79379Z","iopub.execute_i
99 # Transforming the unmatched subjects to the same notation
100 df.subject=df.subject.replace({'politics':'PoliticsNews','politicsNews':'PoliticsNews'})
101
102 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.819181Z","iopub.execute_
103 sub_tf_df=df.groupby('target').apply(lambda x:x['title'].count()).reset_index(name='Counts'
104 sub_tf_df.target.replace({0:'False',1:'True'},inplace=True)
105 fig = px.bar(sub_tf_df, x="target", y="Counts",
106             color='Counts', barmode='group',
107             height=350)
108 fig.show()
109
110
111
112 # %% [markdown]
113 # - The data looks balanced and no issues on building the model
114
115 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:14.938929Z","iopub.execute_
116 sub_check=df.groupby('subject').apply(lambda x:x['title'].count()).reset_index(name='Counts'
117 fig=px.bar(sub_check,x='subject',y='Counts',color='Counts',title='Count of News Articles by
118 fig.show()
119
120 # %% [markdown]
121 #
122
123 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:15.056981Z","iopub.execute_
124 df['clean_title'] = df['title'].apply(preprocess)
125 df['clean_title'][0]
126
127 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:18.393978Z","iopub.execute_
128 df['clean_joined_title']=df['clean_title'].apply(lambda x:" ".join(x))
129
130 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:18.444948Z","iopub.execute_
131 plt.figure(figsize = (20,20))

```

```
132 wc = WordCloud(max_words = 2000 , width = 1600 , height = 800 , stopwords = stop_words).gen
133 plt.imshow(wc, interpolation = 'bilinear')
134
135 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:40.057751Z","iopub.execute_
136 maxlen = -1
137 for doc in df.clean_joined_title:
138     tokens = nltk.word_tokenize(doc)
139     if(maxlen<len(tokens)):
140         maxlen = len(tokens)
141 print("The maximum number of words in a title is =", maxlen)
142 fig = px.histogram(x = [len(nltk.word_tokenize(x)) for x in df.clean_joined_title], nbins =
143 fig.show()
144
145 # %% [markdown]
146 # *Creating Prediction Model*
147
148 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:57.606819Z","iopub.execute_
149 X_train, X_test, y_train, y_test = train_test_split(df.clean_joined_title, df.target, test_
150 vec_train = CountVectorizer().fit(X_train)
151 X_vec_train = vec_train.transform(X_train)
152 X_vec_test = vec_train.transform(X_test)
153
154 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:20:59.134487Z","iopub.execute_
155 #model
156 model = LogisticRegression(C=2)
157
158 #fit the model
159 model.fit(X_vec_train, y_train)
160 predicted_value = model.predict(X_vec_test)
161
162 #accuracy & predicted value
163 accuracy_value = roc_auc_score(y_test, predicted_value)
164 print(accuracy_value)
165
166 # %% [markdown]
167 # *Create the confusion matrix*
168
169 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:21:01.621433Z","iopub.execute_
170 cm = confusion_matrix(list(y_test), predicted_value)
171 plt.figure(figsize = (7, 7))
172 sns.heatmap(cm, annot = True,fmt='g',cmap='viridis')
173
174 # %% [markdown]
175 # - 4465 Fake News have been Classified as Fake
176 # - 4045 Real News have been classified as Real
177
178 # %% [markdown]
179 # *Checking the content of news*
180
```

```

181 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:21:02.055182Z","iopub.execute_
182 df['clean_text'] = df['text'].apply(preprocess)
183 df['clean_joined_text']=df['clean_text'].apply(lambda x:" ".join(x))
184
185 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:22:28.723648Z","iopub.execute_
186 plt.figure(figsize = (20,20))
187 wc = WordCloud(max_words = 2000 , width = 1600 , height = 800 , stopwords = stop_words).gen
188 plt.imshow(wc, interpolation = 'bilinear')
189
190 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:23:30.622781Z","iopub.execute_
191 maxlen = -1
192 for doc in df.clean_joined_text:
193     tokens = nltk.word_tokenize(doc)
194     if(maxlen<len(tokens)):
195         maxlen = len(tokens)
196 print("The maximum number of words in a News Content is =", maxlen)
197 fig = px.histogram(x = [len(nltk.word_tokenize(x)) for x in df.clean_joined_text], nbins =
198 fig.show()
199
200
201 # %% [markdown]
202 # *Predicting the Model*
203
204 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:27:23.232691Z","iopub.execute_
205 X_train, X_test, y_train, y_test = train_test_split(df.clean_joined_text, df.target, test_s
206 vec_train = CountVectorizer().fit(X_train)
207 X_vec_train = vec_train.transform(X_train)
208 X_vec_test = vec_train.transform(X_test)
209 model = LogisticRegression(C=2.5)
210 model.fit(X_vec_train, y_train)
211 predicted_value = model.predict(X_vec_test)
212 accuracy_value = roc_auc_score(y_test, predicted_value)
213 print(accuracy_value)
214
215 # %% [code] {"execution":{"iopub.status.busy":"2023-10-24T21:30:15.444444Z","iopub.execute_
216 prediction = []
217 for i in range(len(predicted_value)):
218     if predicted_value[i].item() > 0.5:
219         prediction.append(1)
220     else:
221         prediction.append(0)
222 cm = confusion_matrix(list(y_test), prediction)
223 plt.figure(figsize = (6, 6))
224 sns.heatmap(cm, annot = True,fmt='g')
225
226 # %% [code]

```

