

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import re
import nltk
nltk.download('omw-1.4')
from nltk.corpus import stopwords
nltk.download('stopwords')
from nltk.stem import PorterStemmer
from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')
from nltk.tokenize import word_tokenize
from nltk.tokenize import sent_tokenize
nltk.download('punkt')
from sklearn.feature_extraction.text import CountVectorizer
from wordcloud import WordCloud
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import make_scorer, roc_curve, roc_auc_score
from sklearn.metrics import precision_recall_fscore_support as score
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.multiclass import OneVsRestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC, LinearSVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
dataset = pd.read_csv("sample_data/BBC News.csv")
print(dataset.head())
dataset.shape
dataset.info()
target_category = dataset['Category'].unique()
print(target_category)
dataset['CategoryId'] = dataset['Category'].factorize()[0]
dataset.head()
category = dataset[['Category', 'CategoryId']].drop_duplicates().sort_values('CategoryId')
category
dataset.groupby('Category').CategoryId.value_counts().plot(kind = "bar", color = ["pink", "orange", "red", "yellow", "blue"])
plt.xlabel("Category of data")
plt.title("Visulaize numbers of Category of data")
plt.show()
fig = plt.figure(figsize = (5,5))
colors = ["skyblue"]
business = dataset[dataset['CategoryId'] == 0 ]
tech = dataset[dataset['CategoryId'] == 1 ]
politics = dataset[dataset['CategoryId'] == 2]
sport = dataset[dataset['CategoryId'] == 3]
entertainment = dataset[dataset['CategoryId'] == 4]
count = [business['CategoryId'].count(), tech['CategoryId'].count(), politics['CategoryId'].count(), sport['CategoryId'].count(), entertainment['CategoryId'].count()]
pie = plt.pie(count, labels = ['business', 'tech', 'politics', 'sport', 'entertainment'],
              autopct = "%1.1f%%",
              shadow = True,
              colors = colors,
              startangle = 45,
              explode = (0.05, 0.05, 0.05, 0.05, 0.05))

stop = set(stopwords.words('english'))

business = dataset[dataset['CategoryId'] == 0]

business = business['Text']

tech = dataset[dataset['CategoryId'] == 1]

tech = tech['Text']

politics = dataset[dataset['CategoryId'] == 2]

politics = politics['Text']

```

```

sport = dataset[dataset['CategoryId'] == 3]

sport = sport['Text']

entertainment = dataset[dataset['CategoryId'] == 4]

entertainment = entertainment['Text']

def wordcloud_draw(dataset, color = 'white'):

    words = ' '.join(dataset)

    cleaned_word = ' '.join([word for word in words.split()

if (word != 'news' and word != 'text')])

    wordcloud = WordCloud(stopwords = stop,background_color = color,

width = 2500, height = 2500).generate(cleaned_word)

    plt.figure(1, figsize = (10,7))

    plt.imshow(wordcloud)

    plt.axis("off")

    plt.show()

print("business related words:")

#wordcloud_draw(business, 'white')

print("tech related words:")

#wordcloud_draw(tech, 'white')

print("politics related words:")

#wordcloud_draw(politics, 'white')

print("sport related words:")

#wordcloud_draw(sport, 'white')

print("entertainment related words:")

#wordcloud_draw(entertainment, 'white')

def remove_tags(text):
    remove = re.compile(r'')
    return re.sub(remove, '', text)
dataset['Text'] = dataset['Text'].apply(remove_tags)
def special_char(text):
    reviews = ''
    for x in text:
        if x.isalnum():
            reviews = reviews + x
        else:
            reviews = reviews + ' '
    return reviews
dataset['Text'] = dataset['Text'].apply(special_char)

def convert_lower(text):
    return text.lower()
dataset['Text'] = dataset['Text'].apply(convert_lower)
dataset['Text'][1]

def remove_stopwords(text):
    stop_words = set(stopwords.words('english'))
    words = word_tokenize(text)
    return [x for x in words if x not in stop_words]
dataset['Text'] = dataset['Text'].apply(remove_stopwords)
dataset['Text'][1]
def lemmatize_word(text):
    wordnet = WordNetLemmatizer()
    return " ".join([wordnet.lemmatize(word) for word in text])
dataset['Text'] = dataset['Text'].apply(lemmatize_word)
dataset['Text'][1]

```

```
x = dataset['Text']
y = dataset['CategoryId']

x = np.array(dataset.iloc[:,0].values)
y = np.array(dataset.CategoryId.values)
cv = CountVectorizer(max_features = 5000)
x = cv.fit_transform(dataset.Text).toarray()
print("X.shape = ",x.shape)
print("y.shape = ",y.shape)
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.3, random_state = 0, shuffle = True)
print(len(x_train))
print(len(x_test))
perform_list = [ ]
def run_model(model_name, est_c, est_pnlty):

    mdl=''

    if model_name == 'Logistic Regression':

        mdl = LogisticRegression()
        oneVsRest = OneVsRestClassifier(mdl)
        oneVsRest.fit(x_train, y_train)

        y_pred = oneVsRest.predict(x_test)

# Performance metrics

    accuracy = round(accuracy_score(y_test, y_pred) * 100, 2)

# Get precision, recall, f1 scores

    precision, recall, f1score, support = score(y_test, y_pred, average='micro')

    print(f'Test Accuracy Score of Basic {model_name}: % {accuracy}')
```

```
print(f'Precision : {precision}')
```

```
print(f'Recall : {recall}')
```

```
print(f'F1-score : {f1score}')
```

```
# Add performance parameters to list

    perform_list.append(dict([

        ('Model', model_name),

        ('Test Accuracy', round(accuracy, 2)),

        ('Precision', round(precision, 2)),

        ('Recall', round(recall, 2)),

        ('F1', round(f1score, 2))

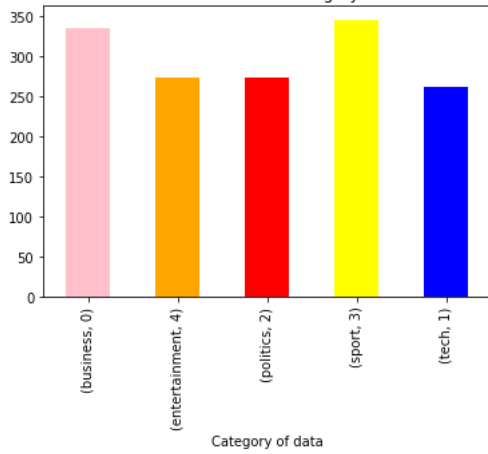
    ]))

run_model('Logistic Regression', est_c=None, est_pnlty=None)
```



```
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
ArticleId      Text      Category
0      1833  worldcom ex-boss launches defence lawyers defe...  business
1       154  german business confidence slides german busin...  business
2      1101  bbc poll indicates economic gloom citizens in ...  business
3      1976  lifestyle governs mobile choice faster bett...   tech
4       917  enron bosses in $168m payout eighteen former e...  business
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1490 entries, 0 to 1489
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   ArticleId  1490 non-null   int64
1   Text       1490 non-null   object
2   Category   1490 non-null   object
dtypes: int64(1), object(2)
memory usage: 35.0+ KB
['business' 'tech' 'politics' 'sport' 'entertainment']
```

Visulaize numbers of Category of data



```
business related words:
tech related words:
politics related words:
sport related words:
entertainment related words:
X.shape = (1490, 5000)
y.shape = (1490,)
1043
447
Test Accuracy Score of Basic Logistic Regression: % 97.09
Precision : 0.970917225950783
Recall : 0.970917225950783
F1-score : 0.9709172259507831
```

