

optimizing spam filtering - Colab

colab.research.google.com/drive/1fMLtGNapuPS1MpogBeQmn20F7t7t6rQ#scrollTo=5i0dLa2TTn

Comment

Share

optimizing spam filtering

File Edit View Insert Runtime Tools Help All changes saved

Files

sample_data

spam.csv

+ Code + Text

[40] import numpy as np # scientific computation
import pandas as pd # loading dataset file
import matplotlib.pyplot as plt # visualization
import nltk # preprocessing our text
from nltk.corpus import stopwords # removing all the stop words
from nltk.stem.porter import PorterStemmer # stemming of words

df=pd.read_csv("/content/spam.csv")
df.head()

[41] from google.colab import drive
drive.mount("/content/drive")

[42]
df.info()

[43]
df.isna().sum()

[26] df.rename({'v1':'label','v2':'text'},inplace=True,axis=1)

[27] # bottom 5 rows of the dataframe
df.tail()

[28] from sklearn.preprocessing import LabelEncoder

0s completed at 1:41 PM

33°C Partly sunny

ENG IN

1:40 PM 3/25/2023

optimizing spam filtering - Colab

colab.research.google.com/drive/1fMx5GhspuP51Mpag3eQnn2047h76rQ#scrollTo=i50dLs2TIn

optimizing spam filtering

File Edit View Insert Runtime Tools Help All changes saved

Files

sample_data

spam.csv

+ Code + Text

Connecting

```
[ ] nltk.download("stopwords")
[nltk_data]Downloading package stopwords to
[nltk_data] c:\users\smart\AppData\Roaming\nltk_data...
[nltk_data] package stopwords is already up-to-date!

import nltk
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer

[ ] import re
corpus = []
length = len(df)

for i in range(0,length):
    text = re.sub("[a-zA-Z0-9]", "",df['text'][i])
    text = text.lower()
    text = text.split()
    ps = PorterStemmer()
    stopword = stopwords.words("english")
    text = [ps.stem(word) for word in text if not word in set(stopword)]
    text = " ".join(text)
    corpus.append(text)

[ ] from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=35000)
X = cv.fit_transform(corpus).toarray()
```

0s completed at 1:41 PM

33°C Partly sunny

ENG IN 1:51 PM 3/23/2023

optimizing spam filtering - Colab

colab.research.google.com/drive/19MkxGNgpuF51MpeqbeQnn2047H76rQ#scrollTo=i50dLa2TIn

optimizing spam filtering

File Edit View Insert Runtime Tools Help All changes saved

Files

sample_data
spam.csv

Code

Test

```
[28] from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['label'] = le.fit_transform(df['label'])

[29] #Splitting data into train and validation sets using train_test_split

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20, random_state = 0)

#train size 80% and test size 20%

[30] print("Before OverSampling, counts of label '1': {}".format(sum(y_train==1)))
print("Before OverSampling, counts of label '0': {}".format(sum(y_train==0)))
#import SMOTE module from imblearn library
#pip install imblearn (if you don't have imblearn in your system)
from imblearn.over_sampling import SMOTE
sm=SMOTE(random_state = 0)
X_train_res, y_train_res = sm.fit_resample(X_train,y_train.ravel())

print('After OverSampling, the shape of train_X:{}'.format(X_train_res.shape))
print('After OverSampling, the shape of train_Y: {}'.format(y_train_res.shape))

print("After OverSampling, counts of label '1': {}".format(sum(y_train_res == 1)))
print("After OverSampling, counts of label '0': {}".format(sum(y_train_res == 0)))

[ ] nltk.download("stopwords")
[nltk_data]Downloading package stopwords to
```

0s completed at 1:41 PM

optimizing spam filtering - Colab

colab.research.google.com/drive/19Mx2GhspuP51MpqgbeQnn20F7H76rQ

optimizing spam filtering

File Edit View Insert Runtime Tools Help Last saved at 1:42 PM

Files

Connecting to a runtime to enable file browsing

(x)

+ Code + Test

Allocating

```
from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=15000)
X = cv.fit_transform(corpus).toarray()

[ ] df.describe()

[ ] df.shape

[ ] df["label"].value_counts().plot(kind="bar",figsize=(12,6))
plt.xticks(np.arange(2),('Non spam','spam'),rotation=0);

[ ] # performing feature scaling operation using standard scaler on x part of the dataset because]
# there different type of values in the column
sc=StandardScaler()
x_bal=sc.fit_transform(x_bal)
x_bal = pd.DataFrame(x_bal,columns=names)

[ ] #splitting data into train and validation sets using train_test_split

from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size = 0.20, random_state = 0)

##train size 80% and test size 20%

[ ] from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
model.fit(X_train,y_train,X_test,y_test)
```

Waiting for client@google.com...

33°C Partly sunny

ENG 1:56 PM