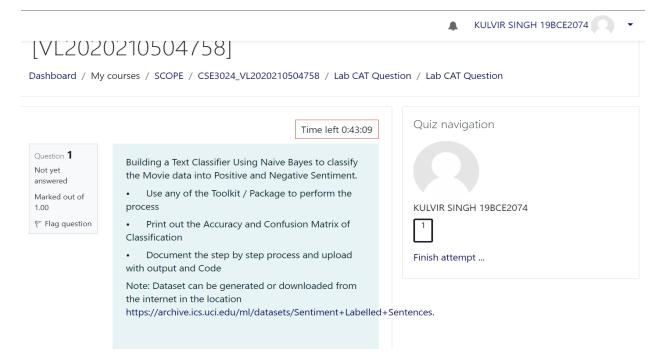
Web Mining Lab CAT

Name: Kulvir Singh

Register Number: 19BCE2074

Question:



Dataset Link:

https://archive.ics.uci.edu/ml/machine-learning-databases/00331/sentiment%20labelled%20sentences.zip lmbd_labelled.txt is used

Procedure

- 1) Download data set from the given link
- 2) Import the necessary packages into the code using the import statement
- 3) Using nltk download the stop words
- 4) Train the txt file as a data frame with sentiment and review as headers using pandas
- 5) Train each of the headers as two separate entities
- 6) Create a tokenizer class with appropriate functions used to assist the tokenizing of the reviews of the document
- 7) Use Naive Bayes classifier to classify the movie data into positive and negative sentiments
- 8) Display the accuracy of the classification
- 9) Display the confusion matrix of the classification
- 10) Display the classification report

Code

```
import numpy as np
import tensorflow as tf
from tensorflow import keras
import pandas as pd
import seaborn as sns
from pylab import rcParams
import string
import re
import matplotlib.pyplot as plt
import math
from matplotlib import rc
from google.colab import drive
from sklearn.model selection import train test split
from collections import Counter, defaultdict
from bs4 import BeautifulSoup
from sklearn.metrics import accuracy score
from sklearn.metrics import classification report, confusion matrix
import nltk
from nltk.corpus import stopwords
from wordcloud import WordCloud
%matplotlib inline
sns.set(style='whitegrid', palette='muted', font scale=1.5)
rcParams['figure.figsize'] = 14, 8
RANDOM SEED = 42
np.random.seed(RANDOM SEED)
nltk.download('stopwords')
train = pd.read csv("imdb labelled.txt", delimiter="\t", names=['review','
sentiment'])
train.head()
X = train['review'].values
y = train['sentiment'].values
class Tokenizer:
  def clean(self, text):
      no html = BeautifulSoup(text).get text()
      clean = re.sub("[^a-z\s]+", " ", no html, flags=re.IGNORECASE)
      return re.sub("(\s+)", " ", clean)
```

```
def tokenize(self, text):
      clean = self.clean(text).lower()
      stopwords en = stopwords.words("english")
      return [w for w in re.split("\W+", clean) if not w in stopwords en]
class MultinomialNaiveBayes:
    def init (self, classes, tokenizer):
      self.tokenizer = tokenizer
      self.classes = classes
    def group by class(self, X, y):
      data = dict()
      for c in self.classes:
        data[c] = X[np.where(y == c)]
      return data
    def fit(self, X, y):
        self.n class items = {}
        self.log class priors = {}
        self.word counts = {}
        self.vocab = set()
       n = len(X)
        grouped data = self.group by class(X, y)
        for c, data in grouped data.items():
          self.n class items[c] = len(data)
          self.log class priors[c] = math.log(self.n class items[c] / n)
          self.word counts[c] = defaultdict(lambda: 0)
          for text in data:
            counts = Counter(self.tokenizer.tokenize(text))
            for word, count in counts.items():
                if word not in self.vocab:
                    self.vocab.add(word)
                self.word counts[c][word] += count
        return self
    def laplace smoothing(self, word, text class):
      num = self.word counts[text class][word] + 1
      denom = self.n class items[text class] + len(self.vocab)
```

```
return math.log(num / denom)
    def predict(self, X):
        result = []
        for text in X:
          class scores = {c: self.log class priors[c] for c in self.classe
s}
          words = set(self.tokenizer.tokenize(text))
          for word in words:
              if word not in self.vocab: continue
              for c in self.classes:
                log w given c = self.laplace smoothing(word, c)
                class scores[c] += log w given c
          result.append(max(class scores, key=class scores.get))
        return result
X train, X test, y train, y test = train test split(X, y, test size=0.2, r
andom state=RANDOM SEED)
MNB = MultinomialNaiveBayes(
    classes=np.unique(y),
    tokenizer=Tokenizer()
).fit(X train, y train)
y hat = MNB.predict(X test)
print("Accuracy = ",accuracy score(y test, y hat))
cnf matrix = confusion matrix(y test, y hat)
print("Confusion Matrix =\n", cnf matrix)
print("Classification Report = \n", classification report(y test, y hat))
```

Output

macro avg

0.81

weighted avg 0.81 0.81 0.81

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
           Package stopwords is already up-to-date!
[nltk_data]
KULVIR SINGH
19BCE2074
____OUTPUT____
Confusion Matrix =
[[59 17]
[12 62]]
Classification Report =
             precision recall f1-score support
                       0.78
         0
                0.83
                                 0.80
                                            76
         1
                0.78
                        0.84
                                 0.81
                                           74
                                 0.81
   accuracy
                                           150
```

0.81

0.81

150

150