

# Assignment 8

## NLP Sentiment Analysis (IMDB Dataset)

**Name :** Atharva Ramgirkar

**Registration Number:** 19BCE0114

**Submission Date :** 15 July, 2021

**Program :** VIT-AI Industry Certification

**Email :** atharva.ramgirkar2019@vitstudent.ac.in

Other Assignments can be found in the link:

[https://drive.google.com/drive/folders/1QGOLHyZykoj\\_CroTJu6-YkZWf32JZ-QH?usp=sharing](https://drive.google.com/drive/folders/1QGOLHyZykoj_CroTJu6-YkZWf32JZ-QH?usp=sharing)

## Table of Content

- [Importing Libraries](#)
  - [Initailizing Objects](#)
- [Reading Data](#)
- [Understanding Data](#)
  - [Checking For Null Values](#)
  - [Checking For Duplicate Values](#)
    - [Dropping Duplicate Rows](#)
  - [Checking if Dataset is Balanced](#)
- [Encoding Target Coulmn](#)
- [Re-ordering the Index](#)
- [Separating the Features and the Target Variable](#)
  - [Storing the Target Variable](#)
  - [NLP on the "review" Column](#)
  - [Vectorizing the Features](#)
- [Train Test Split](#)
- [Model Building](#)
  - [Initialize the Model](#)
  - [Add Input Layer](#)
  - [Add Hidden Layers](#)
  - [Add Output Layer](#)
- [Compile the Model](#)
- [Train the Model](#)
- [Make prediction on Test Data](#)
- [Single Predictions](#)

In [ ]:

In [ ]:

### 1. Importing Libraries

[Back to Top](#)

In [1]:

```
import pandas as pd
import numpy as np

# For NLP
import re
import nltk
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
from sklearn.feature_extraction.text import CountVectorizer

# For Train Test Split
from sklearn.model_selection import train_test_split

# For Neural Network
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense

# For handling Missing Values
import missingno as ms
import matplotlib.pyplot as plt

# Model Evaluation
from sklearn.metrics import accuracy_score
```

#### 1.1 Initailizing Objects

In [2]:

```
ps = PorterStemmer()
cv = CountVectorizer(max_features=4000)
```

In [ ]:

### 2. Reading Data

[Back to Top](#)

In [3]:

```
df = pd.read_csv("IMDB Dataset.csv")
```

In [ ]:

### 3. Understanding Data

[Back to Top](#)

In [4]:

```
df.head()
```

Out[4]:

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production.   The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive

#### 3.1 Checking For Null Values

In [5]:

```
df.isnull().sum()
```

Out[5]:

```
review      0
sentiment   0
dtype: int64
```

#### 3.2 Checking For Duplicate Values

In [6]:

```
df.shape
```

Out[6]:

```
(50000, 2)
```

In [7]:

```
len(df['review'].unique())
```

Out[7]:

```
49582
```

##### 3.2.1 Dropping Duplicate Rows

In [4]:

```
df.drop_duplicates(inplace=True)
```

In [9]:

```
df.shape
```

Out[9]:

```
(49582, 2)
```

In [10]:

```
len(df['review'].unique())
```

Out[10]:

```
49582
```

In [11]:

```
df.head()
```

Out[11]:

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production.   The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive

#### 3.3 Checking if Dataset is Balanced

In [12]:

```
df['sentiment'].value_counts()
```

Out[12]:

```
positive      24884
negative      24698
Name: sentiment, dtype: int64
```

### 4. Encoding Target Coulmn

[Back to Top](#)

In [5]:

```
df['sentiment'] = np.where(df['sentiment']=="positive",1,0)
```

In [14]:

```
df.head()
```

Out[14]:

	review	sentiment
0	One of the other reviewers has mentioned that ...	1
1	A wonderful little production.   The...	1
2	I thought this was a wonderful way to spend ti...	1
3	Basically there's a family where a little boy ...	0
4	Petter Mattei's "Love in the Time of Money" is...	1

### 5. Re-ordering the Index

[Back to Top](#)

In [6]:

```
df.reset_index(drop=True,inplace=True)
```

In [ ]:

### 6. Separating the Features and the Target Variable

[Back to Top](#)

In [7]:

```
y = df['sentiment'].values
```

#### 6.1 Storing the Target Variable

In [16]:

```
y
```

Out[16]:

```
array([1, 1, 1, ..., 0, 0, 0])
```

#### 6.2 NLP on the "review" Column

In [8]:

```
data = []
for i in range(0,49582):
    rev = df['review'][i]

    # Removing Special Characters
    rev = re.sub('[^a-zA-Z]'," ",rev)

    # Converting to Lower Case
    rev = rev.lower()

    # Spliting Sentences to List of Words
    rev = rev.split()

    # Stemming and Stop Word Removal
    rev = [ps.stem(word) for word in rev if not word in set(stopwords.words('english'))]

    # Re-Forming Sentence
    rev = " ".join(rev)

    # Appending to Corpus
    data.append(rev)
```

#### 6.3 Vectorizing the Features

In [9]:

```
X = cv.fit_transform(data).toarray()
```

In [ ]:

### 7. Train Test Split

[Back to Top](#)

In [10]:

```
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.3,random_state=114,shuf
```

In [ ]:

### 8. Model Building

[Back to Top](#)

#### 8.1 Initialize the Model

In [11]:

```
model = Sequential()
```

#### 8.2 Add Input Layer

In [12]:

```
model.add(Dense(units = 4000,
                 kernel_initializer="random_uniform",
                 activation="relu"))
```

#### 8.3 Add Hidden Layers

In [13]:

```
model.add(Dense(units = 2000,
                 kernel_initializer="random_uniform",
                 activation="relu"))

model.add(Dense(units = 1000,
                 kernel_initializer="normal",
                 activation="relu"))
```

#### 8.4 Add Output Layer

In [14]:

```
model.add(Dense(units = 1,
                 kernel_initializer="random_uniform",
                 activation="sigmoid"))
```

### 9. Compile the Model

[Back to Top](#)

In [15]:

```
model.compile(optimizer="adam",
              loss="binary_crossentropy",
              metrics=['accuracy'])
```

In [ ]:

### 10. Train the Model

[Back to Top](#)

In [16]:

```
model.fit(X_train,y_train,epochs=2)
```

Epoch 1/2

```
1085/1085 [=====] - 199s 184ms/step - loss: 0.3416 - accurac
y: 0.8583
Epoch 2/2
```

```
1085/1085 [=====] - 189s 174ms/step - loss: 0.1636 - accurac
y: 0.9359
```

Out[16]:

```
<tensorflow.python.keras.callbacks.History at 0x16a8b5dba90>
```

In [ ]:

### 11. Make prediction on Test Data

[Back to Top](#)

In [17]:

```
pred = model.predict(X_test)
pred = pred>0.64
```

In [18]:

```
accuracy_score(y_test, pred)
```

Out[18]:

```
0.864672268907563
```

In [ ]:

### 12. Single Predictions

[Back to Top](#)

In [19]:

```
model.predict(cv.transform(["amazing good awesome movie"]))
```

Out[19]:

```
array([[0.6684723]], dtype=float32)
```

In [20]:

```
model.predict(cv.transform(["bad no good boring long movie"]))
```

Out[20]:

```
array([[0.42939442]], dtype=float32)
```

In [21]:

```
model.predict(cv.transform(["movie waste of time"]))
```

Out[21]:

```
array([[0.65461886]], dtype=float32)
```