

Theory Activity No. 1

NAME: PUSHKAR HARSHKUMAR PARAKH

PRN: 202401040313

CLASS:CS4

ROLL NO.:54

Text Classification Assignment - IMDb Dataset

CO

Text_Classification_IMDb_Assignment.ipynb

☆

🔗

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

RAM 100% Disk 100%

↑ ↓ ↺ ↻ 📄 📁 🗑️ ⋮

Text Classification Assignment - IMDb Dataset

Import Libraries and Dataset

[7]

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from google.colab import files
uploaded = files.upload()

IMDb Dataset.csv

IMDb Dataset.csv(text/csv) - 66212309 bytes, last modified: 19/10/2019 - 100% done
Saving IMDb Dataset.csv to IMDb Dataset.csv

[8]

Load Dataset
df = pd.read_csv('IMDb Dataset.csv')
df.head()

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 Matte's "Love in the Time of Money" is...

positive

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

1. View the shape of the dataset

Text_Classification_IMDb_Assignment.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

1. View the shape of the dataset

```
[9] df.shape
```

(50000, 2)

2. Display the first 10 rows

```
df.head(10)
```

	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 Matte's "Love in the Time of Money" is...	positive
5	Probably my all-time favorite movie, a story o...	positive
6	I sure would like to see a resurrection of a u...	positive
7	This show was an amazing, fresh & innovative i...	negative
8	Encouraged by the positive comments about this...	negative
9	If you like original gut wrenching laughter yo...	positive

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

Text_Classification_IMDb_Assignment.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

3. Check for missing values

```
[11] df.isnull().sum()
```

	review	sentiment
	0	0

dtype: int64

4. Distribution of Sentiment Labels

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

	count
positive	25000
negative	25000

dtype: int64

5. Percentage of Sentiments

```
[13] df['sentiment'].value_counts(normalize=True) * 100
```

	proportion
positive	50.0
negative	50.0

dtype: float64

colab.research.google.com/drive/1g26C61evaX6L_HFHBeKOlUPotAEKb40#scrollTo=WbFhfxjhcB5D

Text_Classification_IMDb_Assignment.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

RAM Disk

6. Add Review Length Column

```
df['review_length'] = df['review'].apply(len)
df[['review', 'review_length']].head()
```

	review	review_length
0	One of the other reviewers has mentioned that ...	1761
1	A wonderful little production. The...	998
2	I thought this was a wonderful way to spend ti...	926
3	Basically there's a family where a little boy ...	748
4	Petter Mattei's "Love in the Time of Money" is...	1317

7. Average Review Length

```
df['review_length'].mean()
np.float64(1389.43182)
```

8. Longest Review

```
df.loc[df['review_length'].idxmax()]
```

	review	sentiment	review_length
31481	Match 1: Tag Team Table Match Bubba Ray and Sp...	positive	13704

dtype: object

colab.research.google.com/drive/1g26C61evaX6L_HFHBeKOlUPotAEKb40#scrollTo=WbFhfxjhcB5D

Text_Classification_IMDb_Assignment.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

RAM Disk

9. Shortest Review

```
df.loc[df['review_length'].idxmin()]
```

	review	sentiment	review_length
27521	Read the book, forget the movie!	negative	32

dtype: object

10. Number of Reviews Containing 'good'

```
df['review'].str.contains('good', case=False).sum()
np.int64(19472)
```

11. Number of Reviews Containing 'bad'

```
df['review'].str.contains('bad', case=False).sum()
np.int64(12784)
```

colab.research.google.com/drive/1g26C61evaX6L_HFHBekOIUPotAEKb40#scrollTo=WbFhFjhcBSD

Text_Classification_IMDb_Assignment.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

RAM Disk

12. Add Word Count Column

```
[20] df['word_count'] = df['review'].apply(lambda x: len(x.split()))
df[['review', 'word_count']].head()
```

	review	word_count
0	One of the other reviewers has mentioned that ...	307
1	A wonderful little production. The...	162
2	I thought this was a wonderful way to spend ti...	166
3	Basically there's a family where a little boy ...	138
4	Petter Mattei's "Love in the Time of Money" is...	230

13. Average Word Count

```
[21] df['word_count'].mean()
np.float64(231.15694)
```

14. Correlation between Review Length and Word Count

```
df[['review_length', 'word_count']].corr()
```

	review_length	word_count
review_length	1.00000	0.99683
word_count	0.99683	1.00000

colab.research.google.com/drive/1g26C61evaX6L_HFHBekOIUPotAEKb40#scrollTo=WbFhFjhcBSD

Text_Classification_IMDb_Assignment.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

RAM Disk

15. Reviews Containing Numbers

```
[23] df['review'].str.contains(r'\d').sum()
np.int64(28005)
```

16. Reviews Fully Uppercase

```
[24] df['review'].apply(lambda x: x.isupper()).sum()
np.int64(1)
```

