# Theory Activity No. 1

Name: Vinay Bhure

PRN:202401070031

Roll no: ET2-12

Batch: ET-21

## Dataset: Text Classification

Dataset Link: https://www.kaggle.com/datasets/vstepanenko/disaster-tweets

### 20 Problem statements:

1. Find the total number of tweets in the dataset.

2. Find the number of real disaster and non-disaster tweets.

3. Calculate the percentage of real disaster tweets using NumPy.

4. Calculate the percentage of non-disaster tweets using NumPy.

5. Find the number of missing values in keyword and location.

6. Replace missing keywords with "unknown" using NumPy.

7. Replace missing locations with "unknown" using NumPy.

8. Add a new column text_length representing the number of characters in the tweet.

9. Add a new column word_count representing the number of words in the tweet.

10. Calculate the average word count of disaster tweets using NumPy.

11. Calculate the average word count of non-disaster tweets using NumPy.

12. Find the most common keyword among real disaster tweets.

13. Find the most common keyword among non-disaster tweets.

14. Find the tweet with the maximum number of characters.

15. Find the tweet with the minimum number of characters.

16. Remove duplicate tweets and display the new shape of the dataset.

17. Calculate the correlation between text length and target using NumPy.

18. Group tweets by location and count disaster tweets.

19. Find the top 5 locations with the highest number of disaster tweets.

20. Find how many disaster tweets contain the word "help" using Numpy and String Operations.

## Code

```python
1    import pandas as pd
2    import numpy as np
3
4    df = pd.read_csv('tweets.csv')
5
6    # 1. Find the total number of tweets
7    print("1. Find the total number of tweets:")
8    total_tweets = df.shape[0]
9    print(f"Total number of tweets: {total_tweets}\n")
10
11   # 2. Find the number of real disaster and non-disaster tweets
12   print("2. Find the number of real disaster and non-disaster tweets:")
13   target_counts = df['target'].value_counts()
14   print("\nDisaster (1) and Non-disaster (0) counts:\n", target_counts, "\n")
15
16   # 3. Calculate the percentage of real disaster tweets using np
17   print("3. Calculate the percentage of real disaster tweets using NumPy:")
18   disaster_percentage = np.round((np.sum(df['target'] == 1) / total_tweets) * 100, 2)
19   print(f"Percentage of real disaster tweets: {disaster_percentage}%\n")
20
21   # 4. Calculate the percentage of non-disaster tweets using np
22   print("4. Calculate the percentage of non-disaster tweets using NumPy:")
23   non_disaster_percentage = np.round((np.sum(df['target'] == 0) / total_tweets) * 100, 2)
24   print(f"Percentage of non-disaster tweets: {non_disaster_percentage}%\n")
25
26   # 5. Find missing values
27   print("5. Find the number of missing values in 'keyword' and 'location':")
28   missing_keyword = df['keyword'].isnull().sum()
29   missing_location = df['location'].isnull().sum()
30   print(f"Missing keywords: {missing_keyword}")
31   print(f"Missing locations: {missing_location}\n")
32
33   # 6. Fill missing keywords using np.where
34   print("6. Replace missing keywords with 'unknown' using NumPy:")
35   df['keyword'] = np.where(df['keyword'].isnull(), 'unknown', df['keyword'])
36   print("Missing keywords replaced with 'unknown'.\n")
37
```

```python
 38    # 7. Fill missing locations using np.where
 39    print("7. Replace missing locations with 'unknown' using NumPy:")
 40    df['location'] = np.where(df['location'].isnull(), 'unknown', df['location'])
 41    print("Missing locations replaced with 'unknown'.\n")
 42
 43    # 8. Add a column 'text_length'
 44    print("8. Add a new column 'text_length' representing the number of characters in the tweet:")
 45    df['text_length'] = df['text'].apply(len)
 46    print("Column 'text_length' added.\n")
 47
 48    # 9. Add a column 'word_count'
 49    print("9. Add a new column 'word_count' representing the number of words in the tweet:")
 50    df['word_count'] = df['text'].apply(lambda x: len(str(x).split()))
 51    print("Column 'word_count' added.\n")
 52
 53    # 10. Average word count of disaster tweets using np
 54    print("10. Calculate the average word count of disaster tweets using NumPy:")
 55    avg_disaster_words = np.round(df[df['target'] == 1]['word_count'].mean(), 2)
 56    print(f"Average word count in disaster tweets: {avg_disaster_words}\n")
 57
 58    # 11. Average word count of non-disaster tweets using np
 59    print("11. Calculate the average word count of non-disaster tweets using NumPy:")
 60    avg_nondisaster_words = np.round(df[df['target'] == 0]['word_count'].mean(), 2)
 61    print(f"Average word count in non-disaster tweets: {avg_nondisaster_words}\n")
 62
 63    # 12. Most common keyword among disaster tweets
 64    print("12. Find the most common keyword among real disaster tweets:")
 65    common_disaster_keyword = df[df['target'] == 1]['keyword'].mode()[0]
 66    print(f"Most common keyword among disaster tweets: {common_disaster_keyword}\n")
 67
 68    # 13. Most common keyword among non-disaster tweets
 69    print("13. Find the most common keyword among non-disaster tweets:")
 70    common_nondisaster_keyword = df[df['target'] == 0]['keyword'].mode()[0]
 71    print(f"Most common keyword among non-disaster tweets: {common_nondisaster_keyword}\n")
 72
```

```python
 73    # 14. Tweet with maximum characters
 74    print("14. Find the tweet with the maximum number of characters:")
 75    max_char_tweet = df.loc[df['text_length'].idxmax()]
 76    print("\nTweet with maximum characters:\n", max_char_tweet['text'], "\n")
 77
 78    # 15. Tweet with minimum characters
 79    print("15. Find the tweet with the minimum number of characters:")
 80    min_char_tweet = df.loc[df['text_length'].idxmin()]
 81    print("\nTweet with minimum characters:\n", min_char_tweet['text'], "\n")
 82
 83    # 16. Remove duplicate tweets
 84    print("16. Remove duplicate tweets and display the new shape of the dataset:")
 85    df = df.drop_duplicates(subset='text')
 86    print("Shape after removing duplicate tweets:", df.shape, "\n")
 87
 88    # 17. Calculate the correlation between text length and target using np
 89    print("17. Calculate the correlation between text length and target using NumPy:")
 90    correlation = np.corrcoef(df['text_length'], df['target'])[0, 1]
 91    print(f"Correlation between text length and target: {correlation:.4f}\n")
 92
 93    # 18. Group tweets by location and count disaster tweets
 94    print("18. Group tweets by location and count disaster tweets:")
 95    disaster_by_location = df[df['target'] == 1].groupby('location').size().sort_values(ascending=False)
 96    print("\nDisaster tweets by location:\n", disaster_by_location.head(), "\n")
 97
 98    # 19. Top 5 locations with the highest number of disaster tweets
 99    print("19. Find the top 5 locations with the highest number of disaster tweets:")
100    top5_locations = disaster_by_location.head(5)
101    print("\nTop 5 locations with most disaster tweets:\n", top5_locations, "\n")
102
103    # 20. Find how many disaster tweets contain the word 'help'
104    print("20. Find how many disaster tweets contain the word 'help' using NumPy and String Operations:")
105    help_tweets = df[(df['target'] == 1) & (np.char.find(df['text'].values.astype(str), 'help') >= 0)]
106    print(f"\nNumber of disaster tweets containing 'help': {help_tweets.shape[0]}")
107    print("\nSample tweets containing 'help':\n", help_tweets['text'].head())
108
```

```
(EDS) PS D:\EDS> python .\EDS_Yashraj221B.py
1. Find the total number of tweets:
Total number of tweets: 11370

2. Find the number of real disaster and non-disaster tweets:

Disaster (1) and Non-disaster (0) counts:
 target
0    9256
1    2114
Name: count, dtype: int64

3. Calculate the percentage of real disaster tweets using NumPy:
Percentage of real disaster tweets: 18.59%

4. Calculate the percentage of non-disaster tweets using NumPy:
Percentage of non-disaster tweets: 81.41%

5. Find the number of missing values in 'keyword' and 'location':
Missing keywords: 0
Missing locations: 3418

6. Replace missing keywords with 'unknown' using NumPy:
Missing keywords replaced with 'unknown'.

7. Replace missing locations with 'unknown' using NumPy:
Missing locations replaced with 'unknown'.

8. Add a new column 'text_length' representing the number of characters in the tweet:
Column 'text_length' added.

9. Add a new column 'word_count' representing the number of words in the tweet:
Column 'word_count' added.

10. Calculate the average word count of disaster tweets using NumPy:
Average word count in disaster tweets: 17.36

11. Calculate the average word count of non-disaster tweets using NumPy:
Average word count in non-disaster tweets: 17.12

12. Find the most common keyword among real disaster tweets:
Most common keyword among disaster tweets: thunderstorm
```

```
13. Find the most common keyword among non-disaster tweets:
Most common keyword among non-disaster tweets: drowning

14. Find the tweet with the maximum number of characters:

Tweet with maximum characters:
 &gt; Get new bicycle saddle &gt; Manual entirely in Chinese &gt; I've got engineering qualifications I'm sure I can figure o… https://t.co/mL94RxUiyx

15. Find the tweet with the minimum number of characters:

Tweet with minimum characters:
 Hello

16. Remove duplicate tweets and display the new shape of the dataset:
Shape after removing duplicate tweets: (11223, 7)

17. Calculate the correlation between text length and target using NumPy:
Correlation between text length and target: 0.1148

18. Group tweets by location and count disaster tweets:

Disaster tweets by location:
 location
unknown          573
UK                29
Ireland           21
United States     20
London, England   19
dtype: int64

19. Find the top 5 locations with the highest number of disaster tweets:

Top 5 locations with most disaster tweets:
 location
unknown          573
UK                29
Ireland           21
United States     20
London, England   19
dtype: int64
```

```
20. Find how many disaster tweets contain the word 'help' using NumPy and String Operations:

Number of disaster tweets containing 'help': 23

Sample tweets containing 'help':
 1938    Australian guy has been helping koalas after t...
3085    Dozen of people reportedly dead in iceberg in ...
3953    Shot and killed in front of his wife, sister a...
4078    Morocco allegedly offered GWB 3,000 monkeys to...
4086    Our country has been devastated by bushfires &...
Name: text, dtype: object
(EDS) PS D:\EDS> |
```