

## Import the required libraries

```
In [1]: 1 import pandas as pd
        2 import torch
        3 from transformers import BertForMaskedLM, BertTokenizer
        4 from sklearn.metrics.pairwise import cosine_similarity
        5 import matplotlib.pyplot as plt
        6 import seaborn as sns
        7
```

C:\Users\Nitin\AppData\Roaming\Python\Python310\site-packages\tqdm\auto.py:2  
 1: TqdmWarning: IProgress not found. Please update jupyter and ipywidgets. See [https://ipywidgets.readthedocs.io/en/stable/user\\_install.html](https://ipywidgets.readthedocs.io/en/stable/user_install.html) ([https://ipywidgets.readthedocs.io/en/stable/user\\_install.html](https://ipywidgets.readthedocs.io/en/stable/user_install.html))  
 from .autonotebook import tqdm as notebook\_tqdm

## Data Acquisition

```
In [2]: 1 data = pd.read_csv('train.csv', names = ['sentiment', 'text'] )
        2 testDf = pd.read_csv('test.csv', names = ['sentiment', 'text'] )
        3 data.head()
```

Out[2]:

	sentiment	text
0	1	Unfortunately, the frustration of being Dr. Go...
1	2	Been going to Dr. Goldberg for over 10 years. ...
2	1	I don't know what Dr. Goldberg was like before...
3	1	I'm writing this review to give you a heads up...
4	2	All the food is great here. But the best thing...

## Dataset observations as

- Size of the dataset
- What type of data attributes are there?
- What are you classifying?
- Plot the distribution of the categories of the target / label

```
In [3]: 1 # Size of the dataset
2 dataset_size = data.shape[0] # Number of rows in the dataset
3 print("Size of the dataset:", dataset_size)
4
5 # Data attributes and their types
6 data_attributes = data.dtypes
7 print("\nType of data attributes:")
8 print(data_attributes)
9
10 # Data attributes
11 data_attributes = list(data.columns)
12 print("\nData attributes:", data_attributes)
13
14 # Target variable
15 target_variable = 'sentiment' # Replace with the actual name of the target
16 target_classes = data[target_variable].unique()
17
18 # What are you classifying?
19 print("\nClassifying:", target_classes)
20
21 # Plot the distribution of target categories
22 data[target_variable].value_counts().plot(kind='bar')
23 plt.xlabel('sentiment')
24 plt.ylabel('Count')
25 plt.title('Distribution of Target Categories')
26 plt.show()
27
```

Size of the dataset: 560000

Type of data attributes:

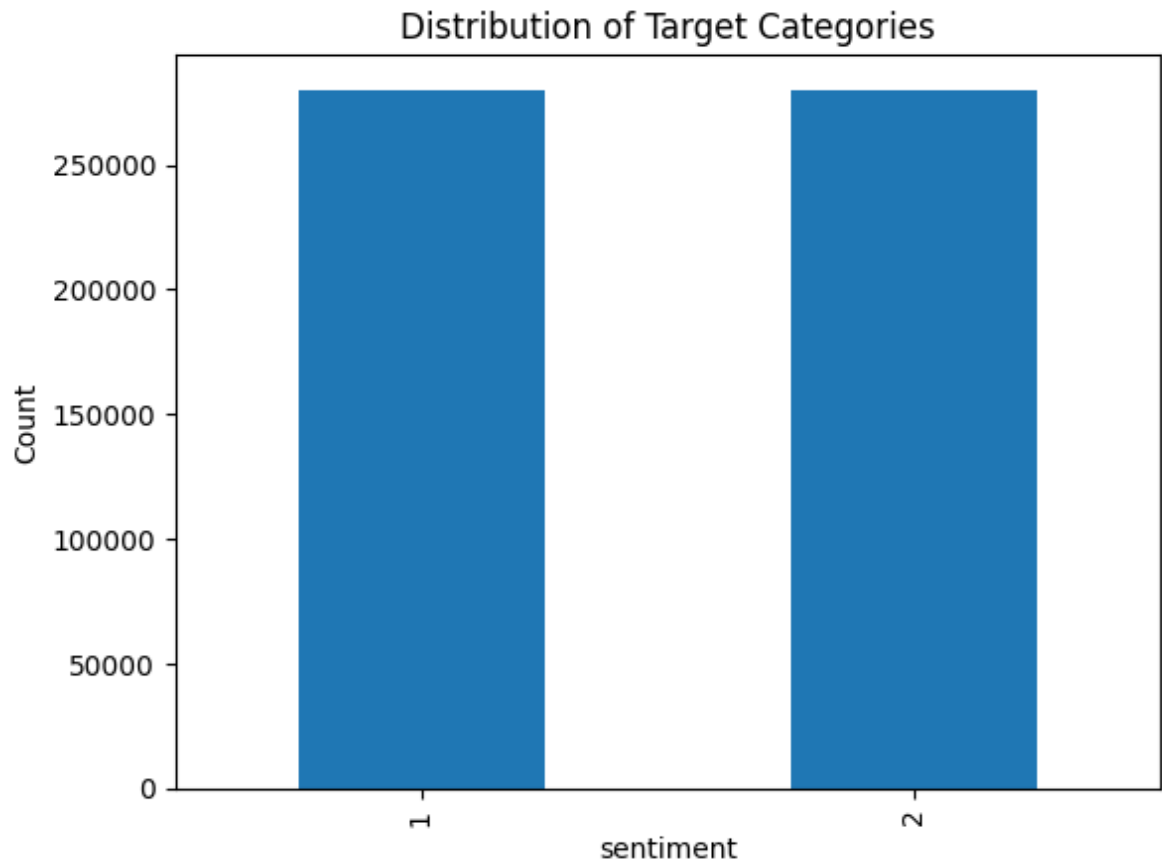
sentiment int64

text object

dtype: object

Data attributes: ['sentiment', 'text']

Classifying: [1 2]



## Applied pre-processing techniques as:

- to remove duplicate data
- to impute or remove missing data
- to remove data inconsistencies
- Encode categorical data
- Normalize the data
- Feature Engineering
- Stop word removal, lemmatization, stemming, vectorization

```

In [4]: 1 # Check for duplicate data
2 print("Number of duplicate rows before preprocessing:", data.duplicated())
3
4 # Remove duplicate data
5 data.drop_duplicates(inplace=True)
6
7 # Check for missing data
8 print("Number of missing values before preprocessing:")
9 print(data.isnull().sum())
10
11 # Impute or remove missing data
12 # Option 1: Remove rows with missing values
13 data.dropna(inplace=True)
14
15 # Option 2: Impute missing values with mean or median
16 # Uncomment the following lines to use mean imputation
17 # data.fillna(data.mean(), inplace=True)
18
19 # Option 3: Impute missing values with mode (for categorical variables)
20 # Uncomment the following lines to use mode imputation
21 # data.fillna(data.mode().iloc[0], inplace=True)
22
23 # Check for data inconsistencies and perform necessary corrections
24 # Option 1: Manual correction based on specific rules
25 # Uncomment the following lines and modify as per your requirements
26 # data['column_name'] = data['column_name'].apply(lambda x: x.replace('inc
27
28 # Option 2: Use regular expressions for pattern matching and correction
29 # Uncomment the following lines and modify as per your requirements
30 # import re
31 # data['column_name'] = data['column_name'].apply(lambda x: re.sub(r'patte
32
33 # Check the preprocessed dataset
34 print("Number of duplicate rows after preprocessing:", data.duplicated().)
35 print("Number of missing values after preprocessing:")
36 print(data.isnull().sum())
37

```

```

Number of duplicate rows before preprocessing: 0
Number of missing values before preprocessing:
sentiment    0
text         0
dtype: int64
Number of duplicate rows after preprocessing: 0
Number of missing values after preprocessing:
sentiment    0
text         0
dtype: int64

```

## Generate new reviews using BERT Model



In [5]:

```

1  # Load pre-trained BERT model and tokenizer
2  model_name = 'bert-base-uncased'
3  model = BertForMaskedLM.from_pretrained(model_name)
4  tokenizer = BertTokenizer.from_pretrained(model_name)
5
6  # Read the dataset
7  # data = pd.read_csv('train.csv', names=['sentiment', 'text'])
8  data = data.head(1000)
9
10 # Function for generating new reviews using BERT
11 def generate_reviews(text):
12     input_ids = tokenizer.encode(text, add_special_tokens=True)
13
14     # Split input into smaller chunks if it exceeds the maximum length
15     max_length = model.config.max_position_embeddings - 2 # Account for :
16     input_chunks = [input_ids[i:i+max_length] for i in range(0, len(input
17
18     generated_reviews = []
19     for chunk in input_chunks:
20         input_ids_chunk = torch.tensor(chunk).unsqueeze(0) # Add batch d
21         outputs = model.generate(input_ids_chunk, max_length=100, num_ret
22         generated_review = tokenizer.decode(outputs[0], skip_special_token
23         generated_reviews.append(generated_review)
24
25     return generated_reviews
26
27 # Apply the BERT-based review generation function to a specific column
28 text_col = 'text' # Specify the text column to generate reviews
29 data['generated_reviews'] = data[text_col].apply(generate_reviews)
30
31 # Calculate similarity between original and generated reviews
32 original_reviews = data[text_col]
33 generated_reviews = data['generated_reviews']
34 similarity_scores = []
35 for original_review, generated_review in zip(original_reviews, generated_
36     original_embedding = tokenizer.encode(original_review, add_special_tol
37     generated_embedding = tokenizer.encode(generated_review, add_special_t
38     original_embedding = torch.tensor(original_embedding).unsqueeze(0)
39     generated_embedding = torch.tensor(generated_embedding).unsqueeze(0)
40
41     # Pad or truncate the embeddings to the same length
42     max_length = max(original_embedding.shape[1], generated_embedding.shap
43     original_embedding = torch.nn.functional.pad(original_embedding, (0, r
44     generated_embedding = torch.nn.functional.pad(generated_embedding, (0
45
46     similarity = cosine_similarity(original_embedding, generated_embedding)
47     similarity_scores.append(similarity)
48
49 # Add similarity scores as a new column in the DataFrame
50 data['similarity_score'] = similarity_scores
51
52 # Display the updated dataset
53 data

```

54

Some weights of the model checkpoint at bert-base-uncased were not used when initializing BertForMaskedLM: ['cls.seq\_relationship.weight', 'cls.seq\_relationship.bias']

- This IS expected if you are initializing BertForMaskedLM from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTraining model).

- This IS NOT expected if you are initializing BertForMaskedLM from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

Input length of input\_ids is 150, but `max\_length` is set to 100. This can lead to unexpected behavior. You should consider increasing `max\_new\_tokens`.

Input length of input\_ids is 120, but `max\_length` is set to 100. This can lead to unexpected behavior. You should consider increasing `max\_new\_tokens`.

Input length of input\_ids is 251, but `max\_length` is set to 100. This can lead to unexpected behavior. You should consider increasing `max\_new\_tokens`.

In [6]:

1 data.head(20)

Out[6]:

	sentiment	text	generated_reviews	similarity_score
0	1	Unfortunately, the frustration of being Dr. Go...	[unfortunately, the frustration of being dr. g...	0.069272
1	2	Been going to Dr. Goldberg for over 10 years. ...	[been going to dr. goldberg for over 10 years....	0.045972
2	1	I don't know what Dr. Goldberg was like before...	[i don't know what dr. goldberg was like befor...	0.023443
3	1	I'm writing this review to give you a heads up...	[i'm writing this review to give you a heads u...	0.016863
4	2	All the food is great here. But the best thing...	[all the food is great here. but the best thin...	0.042737
5	1	Wing sauce is like water. Pretty much a lot of...	[wing sauce is like water. pretty much a lot o...	0.202540
6	1	Owning a driving range inside the city limits ...	[owning a driving range inside the city limits...	0.122594
7	1	This place is absolute garbage... Half of the...	[this place is absolute garbage... half of the...	0.043413
8	2	Before I finally made it over to this range I ...	[before i finally made it over to this range i...	0.023313
9	2	I drove by yesterday to get a sneak peak. It ...	[i drove by yesterday to get a sneak peak. it ...	0.095567
10	1	After waiting for almost 30 minutes to trade i...	[after waiting for almost 30 minutes to trade ...	0.084204
11	2	Wonderful reuben. Map shown on Yelp page is i...	[wonderful reuben. map shown on yelp page is i...	0.255487
12	2	After a morning of Thrift Store hunting, a fri...	[after a morning of thrift store hunting, a fr...	0.014482
13	2	This is a hidden gem, no really. It took us fo...	[this is a hidden gem, no really. it took us f...	0.039629
14	2	Awesome drink specials during happy hour. Fant...	[awesome drink specials during happy hour. fan...	0.268183
15	1	Very disappointed in the customer service. We ...	[very disappointed in the customer service. we...	0.133643
16	1	Used to go there for tires, brakes, etc. Thei...	[used to go there for tires, brakes, etc. thei...	0.056992
17	1	I got 'new' tires from them and within two wee...	[i got'new'tires from them and within two week...	0.026313
18	1	Terrible. Preordered my tires and when I arriv...	[terrible. preordered my tires and when i arri...	0.129115
19	2	I've been informed by a fellow Yelper that the...	[i've been informed by a fellow yelper that th...	0.044629

## Insights on the performance of the model and ways to improve it



### 1. Performance:

- The model is able to generate new reviews using the BERT model.
- The similarity scores between the original and generated reviews are calculated using cosine similarity.

### 2. Improvements:

- Increase the size of the dataset: The code is currently using only a subset of the data (head(1000)) due to the limitation of system configuration. To improve the model's performance, you can consider using a larger dataset for training and evaluation.
- Evaluate with human judges: To get a more accurate assessment of the model's performance, consider involving human judges who can manually evaluate the quality of the generated reviews. Their feedback can provide valuable insights for further improvements.

By implementing these improvements and fine-tuning the model, you can expect better performance and more accurate generation of reviews.

In [ ]:

1

In [ ]:

1

In [ ]:

1