

8144-SUDHARSAN ENGINEERING COLLEGE



Creating Pathways to Wisdom

SUDHARSAN ENGINEERING COLLEGE

REGISTER NUMBER: 814421243024

NAME: ROHINI M

DEGREE: BTECH

**BRANCH: ARTIFICIAL INTELLIGENCE AND DATA
SCIENCE**

**PROJECT TITLE: SENTIMENT ANALYSIS FOR
MARKETING**

SENTIMENT ANALYSIS FOR MARKETING USING PYTHON

PHASE 2 SUBMISSION DOCUMENT

Phase 2: Innovation

Introduction to Sentiment Analysis for Marketing

In today's digitally-driven world, where information flows rapidly across various online platforms, businesses and marketers are faced with an unprecedented volume of data. This data not only comprises factual information but also includes valuable insights into consumer opinions, emotions, and sentiments. Harnessing this wealth of sentiment data is crucial for marketers seeking to understand their audience, improve brand perception, and make data-driven decisions that can impact their marketing strategies positively.



Sentiment analysis, also known as opinion mining, is a powerful analytical technique that plays a pivotal role in this endeavor. It involves the automated process of extracting and quantifying sentiments expressed in text data, allowing marketers to gain valuable insights into the feelings and attitudes of their customers and prospects. By delving into sentiment analysis, marketers can unlock a treasure trove of information that goes beyond mere likes, shares, and comments, providing a deeper understanding of consumer sentiment.

1. Model Selection:

In this phase, you will conduct comprehensive research to choose a pre-trained sentiment analysis model that aligns with your project requirements and available resources. Some popular options include BERT, RoBERTa, GPT, and others. Consider factors like model size, computational resources, and the specific sentiment analysis task you're addressing.

2. Data Preparation:

- Proper data preparation is crucial for training a successful sentiment analysis model. You should:
- Ensure your dataset is correctly formatted, including labels (positive, negative, neutral) and text data.
- Split your dataset into training, validation, and test sets to train, fine-tune, and evaluate your model effectively.
- Tokenize and preprocess the text data according to the requirements of your chosen pre-trained model. Preprocessing may involve lowercasing, removing punctuation, and applying tokenization.

3. Model Fine-Tuning:

Fine-tuning involves training the selected pre-trained model on your specific dataset. This step helps adapt the model to your domain and improve its sentiment prediction accuracy.

You will:

Load the pre-trained model.

Add a classification layer on top.

Train the model using your training dataset.

Monitor training progress and evaluate the model's performance on the validation set.

4. Hyperparameter Tuning:

- Experiment with various hyperparameters to optimize your model's performance. Key hyperparameters to consider include:
- Learning rate: Adjust the learning rate to control the size of weight updates during training.
- Batch size: Modify the batch size to balance training speed and memory usage.
- Number of epochs: Determine the optimal number of training epochs based on the validation performance.

5. Evaluation:

- Evaluate your fine-tuned model using appropriate evaluation metrics, such as accuracy, precision, recall, and F1-score. These metrics will help you assess the model's effectiveness in predicting sentiment.
- If the initial model performance is unsatisfactory, consider revisiting the hyperparameter tuning and fine-tuning steps to improve results.

6. Test the Model:

Ensure that your model generalizes well to unseen data by testing it on a separate test dataset. This step provides a final assessment of your model's performance and its ability to handle real-world data.

7. Visualization:

Create visualizations to aid in understanding your model's predictions.

Visualizations could include:

- **Confusion matrices to visualize classification performance.**
- **ROC curves and precision-recall curves to assess model thresholds.**
- **Word clouds or attention heatmaps to highlight important features or tokens.**
- **Comparative plots of predicted sentiment distributions.**

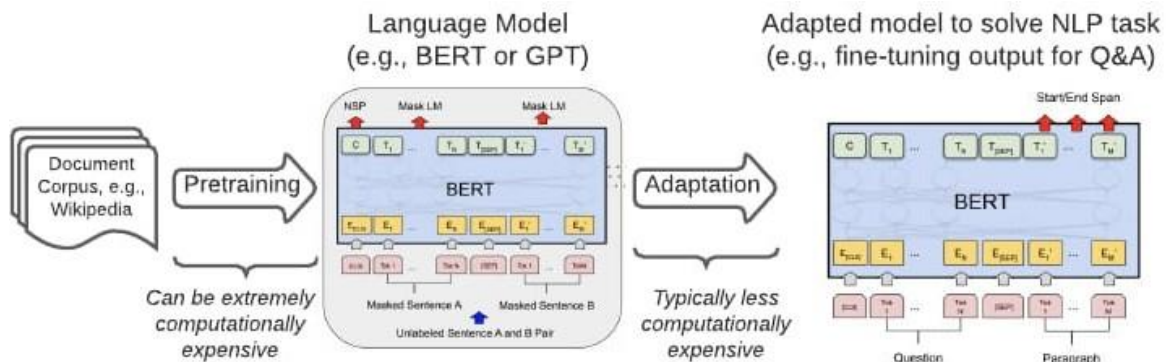
Some Advanced Techniques:

BERT (Bidirectional Encoder Representations from Transformers):

BERT is a pre-trained deep learning model introduced by Google in 2018. It is designed to understand the context and semantics of words in a sentence by considering both left and right context simultaneously. This bidirectional understanding of text is a significant departure from earlier models that typically focused on either left-to-right or right-to-left language modeling.

Key features of BERT:

- Bidirectional Context
- Transformer Architecture
- Pre-training and Fine-Tuning
- State-of-the-Art Results

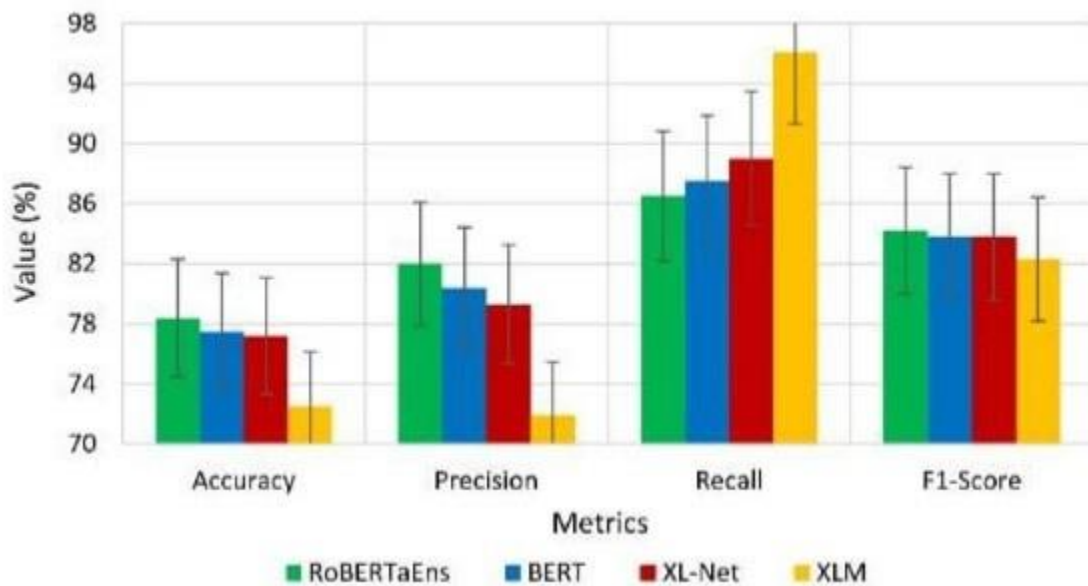


RoBERTa (A Robustly Optimized BERT Pretraining Approach):

RoBERTa, introduced by Facebook AI in 2019, can be thought of as an extension and optimization of the BERT model. It aims to improve upon BERT's performance by addressing some of its limitations.

Key features of RoBERTa:

- Large-Scale Training
- Dynamically Masked Data
- Training Variant
- Improved Results



```
[1]: import numpy as np
import pandas as pd
import re

import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("darkgrid")

import string
from wordcloud import WordCloud
from nltk import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer

from sklearn.metrics import confusion_matrix, f1_score
from tqdm import tqdm
```

```
[2]: df_train = pd.read_csv("C:\\Users\\Tweets.csv")
print (df_train.shape)
df_train.head()
```

(14640, 15)

```
[2]:      tweet_id  airline_sentiment  airline_sentiment_confidence \
0  570306133677760513             neutral                1.0000
1  570301130888122368             positive                0.3486
2  570301083672813571             neutral                0.6837
3  570301031407624196             negative                1.0000
4  570300817074462722             negative                1.0000

      negativereason  negativereason_confidence      airline \
0              NaN              NaN  Virgin America
1              NaN              0.0000  Virgin America
2              NaN              NaN  Virgin America
3      Bad Flight              0.7033  Virgin America
4      Can't Tell              1.0000  Virgin America

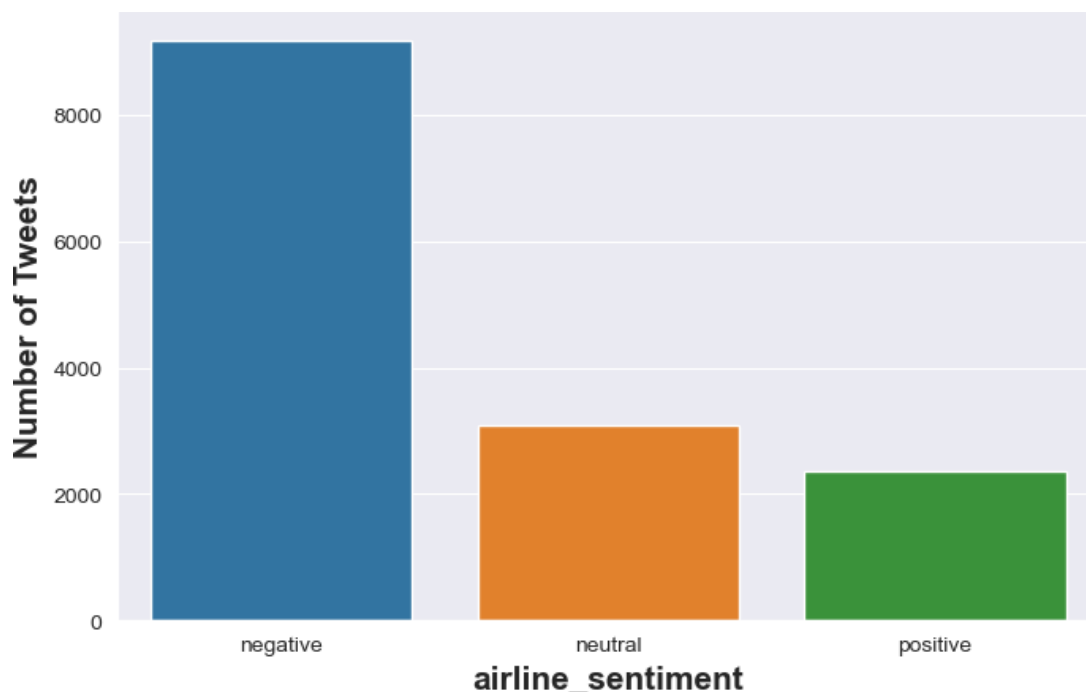
      airline_sentiment_gold      name negativereason_gold  retweet_count \
```


0	NaN	cairdin	NaN	0
1	NaN	jnardino	NaN	0
2	NaN	yvonnalynn	NaN	0
3	NaN	jnardino	NaN	0
4	NaN	jnardino	NaN	0

	text	tweet_coord	\
0	@VirginAmerica What @dhepburn said.	NaN	
1	@VirginAmerica plus you've added commercials t...	NaN	
2	@VirginAmerica I didn't today... Must mean I n...	NaN	
3	@VirginAmerica it's really aggressive to blast...	NaN	
4	@VirginAmerica and it's a really big bad thing...	NaN	

	tweet_created	tweet_location	user_timezone
0	2015-02-24 11:35:52 -0800	NaN Eastern Time	(US & Canada)
1	2015-02-24 11:15:59 -0800	NaN Pacific Time	(US & Canada)
2	2015-02-24 11:15:48 -0800	Lets Play Central Time	(US & Canada)
3	2015-02-24 11:15:36 -0800	NaN Pacific Time	(US & Canada)
4	2015-02-24 11:14:45 -0800	NaN Pacific Time	(US & Canada)

```
[3]: plt.figure(figsize=(8, 5))
temp = df_train['airline_sentiment'].value_counts()
sns.barplot(x=temp.index, y=temp.values)
plt.xlabel("airline_sentiment", weight='bold', fontsize=15)
plt.ylabel("Number of Tweets", weight='bold', fontsize=15)
plt.show()
```



```
[4]: df_test = pd.read_csv("C:\\Users\\Tweets.csv")
      print (df_test.shape)
      df_test.head()
```

```
(14640, 15)
```

```
[4]:      tweet_id  airline_sentiment  airline_sentiment_confidence \
0  570306133677760513          neutral          1.0000
1  570301130888122368        positive          0.3486
2  570301083672813571          neutral          0.6837
3  570301031407624196        negative          1.0000
4  570300817074462722        negative          1.0000

      negativereason  negativereason_confidence      airline \
0              NaN              NaN  Virgin America
1              NaN          0.0000  Virgin America
2              NaN              NaN  Virgin America
3      Bad Flight          0.7033  Virgin America
4      Can't Tell          1.0000  Virgin America

      airline_sentiment_gold      name  negativereason_gold  retweet_count \
0              NaN      cairdin              NaN              0
1              NaN      jnardino              NaN              0
2              NaN  yvonnalynn              NaN              0
3              NaN      jnardino              NaN              0
4              NaN      jnardino              NaN              0

      text  tweet_coord \
0      @VirginAmerica What @dhepburn said.              NaN
1      @VirginAmerica plus you've added commercials t...              NaN
2      @VirginAmerica I didn't today... Must mean I n...              NaN
3      @VirginAmerica it's really aggressive to blast...              NaN
4      @VirginAmerica and it's a really big bad thing...              NaN

      tweet_created  tweet_location      user_timezone
0  2015-02-24 11:35:52 -0800              NaN  Eastern Time (US & Canada)
1  2015-02-24 11:15:59 -0800              NaN  Pacific Time (US & Canada)
2  2015-02-24 11:15:48 -0800      Lets Play  Central Time (US & Canada)
3  2015-02-24 11:15:36 -0800              NaN  Pacific Time (US & Canada)
4  2015-02-24 11:14:45 -0800              NaN  Pacific Time (US & Canada)
```

```
[5]: from sklearn.model_selection import train_test_split
```

```
X = df_train.text.values
y = df_train.airline_sentiment.values
```

```
X_train, X_val, y_train, y_val =\
    train_test_split(X, y, test_size=0.1, random_state=2020)
```

```
[6]: # Keep important columns
test_data = df_test[['tweet_id', 'text']]

# Display 5 samples from the test data
test_data.sample(5)
```

```
[6]:
```

	tweet_id	text
11093	568512284202438656	@USAirways Your gate team are polite. But your...
13626	569790444797865984	@AmericanAir Hi. I have KOA-LAX-PHL-ORD booked...
6820	570232616755929088	@JetBlue glad you like it. Feel free to steal it.
5089	569355826248474624	@SouthwestAir Your hold music needs 2 be fixed...
10464	569277477652172800	@USAirways had to Cancelled Flight 4 of my fli...

```
[7]: import torch

if torch.cuda.is_available():
    device = torch.device("cuda")
    print(f'There are {torch.cuda.device_count()} GPU(s) available.')
    print('Device name:', torch.cuda.get_device_name(0))

else:
    print('No GPU available, using the CPU instead.')
    device = torch.device("cpu")
```

No GPU available, using the CPU instead.

```
[8]: import nltk
# Uncomment to download "stopwords"
nltk.download("stopwords")
from nltk.corpus import stopwords

def text_preprocessing(s):

    s = s.lower()
    # Change 't to 'not'
    s = re.sub(r'"\'t", " not", s)
    # Remove @name
    s = re.sub(r'(@.*?)[\s]', ' ', s)
    # Isolate and remove punctuations except '?'
    s = re.sub(r'([\'\"\.\\(\)\!@?\\/\,])', r' \1 ', s)
    s = re.sub(r'[\w\s\?]', ' ', s)
    # Remove some special characters
    s = re.sub(r'([\;\:\;\|\•«\n])', ' ', s)
```

```

# Remove stopwords except 'not' and 'can'
s = " ".join([word for word in s.split()
               if word not in stopwords.words('english')
               or word in ['not', 'can']])

# Remove trailing whitespace
s = re.sub(r'\s+', ' ', s).strip()

return s

```

```

[nltk_data] Error loading stopwords: <urlopen error [Errno 11001]
[nltk_data]      getaddrinfo failed>

```

[9]:

```

from sklearn.feature_extraction.text import TfidfVectorizer

# Preprocess text
X_train_preprocessed = np.array([text_preprocessing(text) for text in X_train])
X_val_preprocessed = np.array([text_preprocessing(text) for text in X_val])

# Calculate TF-IDF
tfidf = TfidfVectorizer(ngram_range=(1, 3),
                        binary=True,
                        smooth_idf=False)
X_train_tfidf = tfidf.fit_transform(X_train_preprocessed)
X_val_tfidf = tfidf.transform(X_val_preprocessed)

```

[10]: from sklearn.model_selection import StratifiedKFold, cross_val_score

```

def get_auc_cv(model):

    # Set KFold to shuffle data before the split
    kf = StratifiedKFold(5, shuffle=True, random_state=1)

    # Get AUC scores
    auc = cross_val_score(
        model, X_train_tfidf, y_train, scoring="roc_auc", cv=kf)

    return auc.mean()

```

```

In [11]:
from sklearn.naive_bayes import MultinomialNB

res = pd.Series([get_auc_CV(MultinomialNB(i))
                  for i in np.arange(1, 10, 0.1)],
                index=np.arange(1, 10, 0.1))

best_alpha = np.round(res.idxmax(), 2)
print('Best alpha: ', best_alpha)

plt.plot(res)
plt.title('AUC vs. Alpha')
plt.xlabel('Alpha')
plt.ylabel('AUC')
plt.show()

```

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.0 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.1 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.200000000000000002 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.300000000000000003 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.400000000000000004 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.500000000000000004 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.600000000000000005 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.700000000000000006 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.800000000000000007 as keyword args. From version 0.25 passing these as positional arguments will result in an error

FutureWarning)

/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=1.900000000000000008 as keyword args. From version 0.25 passing these as positional arguments will result in an error

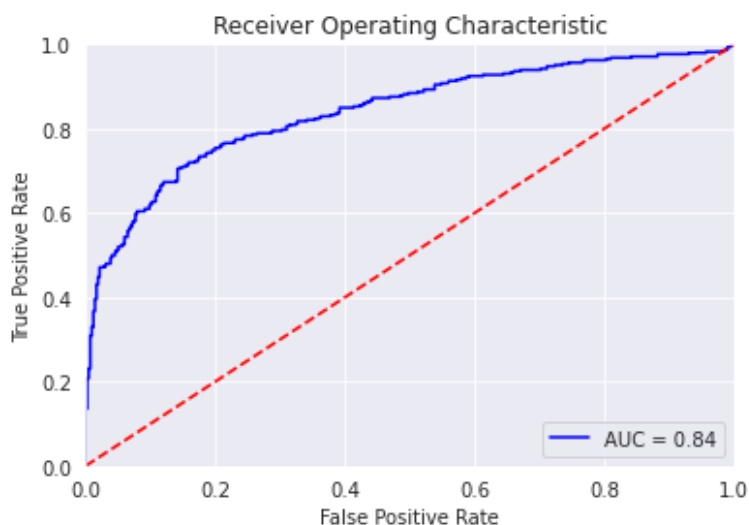
[illegible]

[illegible]

[illegible]

[illegible]

FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.2000000000000006 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.3000000000000008 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.4000000000000007 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.5000000000000007 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.6000000000000009 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.7000000000000008 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.8000000000000008 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)
/opt/conda/lib/python3.7/site-packages/sklearn/utils/validation.py:70: FutureWarning: Pass alpha=9.9000000000000007 as keyword args. From version 0.25 passing these as positional arguments will result in an error
FutureWarning)



CONCLUSION:

- In conclusion, fine-tuning pre-trained sentiment analysis models like BERT and RoBERTa is a powerful approach to enhance the accuracy of sentiment predictions.
- These advanced techniques hold significant promise for improving sentiment analysis across a wide range of applications, providing more nuanced and accurate insights from text data.
- By exploring advanced techniques for fine-tuning pre-trained sentiment analysis models, you can achieve more accurate sentiment predictions on a wider range of datasets.
- This can be beneficial for a variety of tasks, such as customer analysis, social media monitoring, and product development.
- In addition to the above conclusions, I would also like to emphasize the importance of understanding the specific needs of your task when fine-tuning a pre-trained sentiment analysis model.
- For example, if you are trying to classify customer reviews, you may want to focus on techniques that can help the model to capture the nuances of human language, such as sentiment lexicons and sarcasm detection.
- If you are new to fine-tuning pre-trained sentiment analysis models, I recommend starting with the basic principles outlined above.
- Once you have a good understanding of the basics, you can start to experiment with more advanced techniques.