INTERNSHIP: PROJECT REPORT

Internship Project Title	Automate emotion anlaysis of texual comments and feedback		
Project Title Automate emotion anlaysis of texual comments and			
Name of the Company	TATA CONSULTANCY SERVICES		
Name of the Industry Mentor	Debashis Roy		
Name of the Institute	PREPINSTA TECHNOLOGIES PRIVATE LIMITED		

	1	1		i
Start Date	End Date	Total Effort (hrs.)	Project Environment	Tools used
15-06-2023	15-07-2023	125	Programming	Integrated
			Language: Python	Development
			Operating System:	Environment (IDE):
			Windows 10	Google Colab
			Hardware	Machine Learning
			Requirements:	Libraries: scikit-
			CPU: Intel Core i5 or	learn.
			equivalent	Text Visualization
			RAM: 8GB or higher	Tools: Matplotlib,
			Storage: Sufficient	seaborn, word
			disk space for datasets	clouds.
			and libraries	Additional Tools:
			Development Tools:	Pandas (data
			Google Colab	manipulation),
				NumPy (numerical
				computations)

Project Synopsis:

The aim of this project is to automate emotion analysis for textual comments and feedback. Emotion analysis, also known as sentiment analysis, involves identifying and categorizing the emotions expressed in text data, such as comments and feedback. The project seeks to develop a system that can accurately analyze the emotions conveyed in text and provide insights to understand the sentiment of the users.

Solution Approach:

The proposed solution approach for automating emotion analysis involves the following steps:

- 1. Data Collection: Gather a dataset of textual comments and feedback that cover a wide range of emotions. This dataset will be used to train and evaluate the emotion analysis model.
- 2. Data Preprocessing: Clean and preprocess the collected data by removing noise, special characters, and stopwords. Convert the text into a suitable format for further analysis.
- 3. Feature Extraction: Extract relevant features from the preprocessed text data. This step may involve techniques like bag-of-words, word embeddings, or TF-IDF to represent the text.

- 4. Emotion Classification Model: Train a machine learning or deep learning model using the preprocessed data and extracted features. This model should be capable of accurately classifying the emotions expressed in the text data.
- 5. Evaluation: Assess the performance of the emotion analysis model using appropriate evaluation metrics such as accuracy, precision, recall, and F1-score. This step will help determine the effectiveness of the model in identifying different emotions.
- 6. Integration: Integrate the trained model into a user-friendly interface or API, allowing users to input their text comments or feedback and receive the emotion analysis results.

Assumptions:

The project makes the following assumptions:

- 1. The collected dataset is representative of the emotions expressed in textual comments and feedback.
- 2. The input text provided by users for emotion analysis is in the same language as the trained model.
- 3. The emotions analyzed in the project are limited to a predefined set, such as positive, negative, neutral, happiness, sadness, anger, etc.

Project Diagrams:

The system takes in text as input and outputs the emotions expressed in the text. The system is divided into the following components:

Data Collection: This component collects text from a variety of sources, such as social media, customer reviews, or product feedback surveys.

Data Preprocessing: This component cleans the data and removes any noise or irrelevant information. It may also need to tokenize the data, which means breaking it down into individual words or phrases.

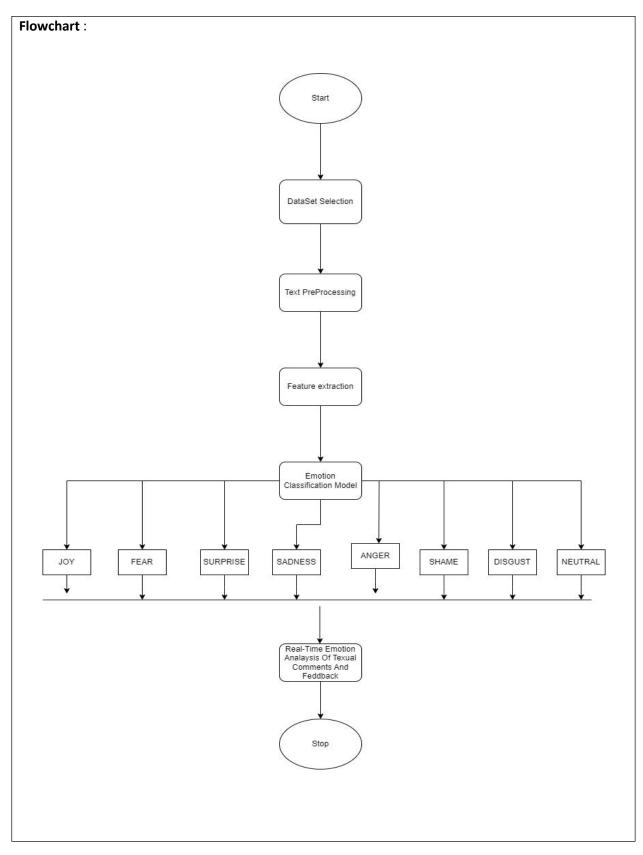
Emotion Keyword Identification: This component identifies emotion keywords. These are words or phrases that are associated with specific emotions. For example, the word "happy" is associated with the emotion of happiness, while the word "sad" is associated with the emotion of sadness.

Sentiment Lexicon Creation: This component creates a sentiment lexicon. A sentiment lexicon is a list of words and phrases that have been manually labeled with their corresponding emotions. This lexicon can be used to train a machine learning model to detect emotions in text.

Machine Learning: This component trains a machine learning model to detect emotions in text. There are a variety of machine learning algorithms that can be used for this task, such as support vector machines, naive Bayes classifiers, and deep learning models.

Emotion Analysis: This component analyzes the text and identifies the emotions expressed in the text. The emotions are then classified into one of the following categories: happiness, sadness, anger, fear, or neutral.

.....



Algorithms:

The project may utilize various algorithms, they are:

1. Text preprocessing algorithms

NeatText is a Python library that provides various text preprocessing functions for cleaning and preparing textual data. It offers a range of functionalities to handle common text preprocessing tasks effectively. Here are some of the preprocessing algorithms provided by the neattext.functions module:

remove_emails: Removes email addresses from the text.

remove_urls: Removes URLs or web links from the text.

remove_html_tags: Removes HTML tags from the text.

remove_punctuation: Removes punctuation marks from the text.

remove_special_characters: Removes special characters from the text.

remove_stopwords: Removes common stopwords from the text (e.g., "the", "is", "and"). remove whitespace: Removes leading, trailing, and excess whitespaces from the text.

normalize_whitespace: Replaces multiple consecutive whitespaces with a single space.

remove digits: Removes digits or numbers from the text.

remove_multiple_spaces: Replaces multiple consecutive spaces with a single space.

remove_accents: Removes accents and diacritical marks from the text.

expand contractions: Expands contractions in the text (e.g., "don't" to "do not").

replace_currency_symbols: Replaces currency symbols with their respective word representation.

replace emojis: Replaces emojis with their textual representation.

replace_urls: Replaces URLs or web links with a placeholder or specific text.

These functions can be applied to the text data as preprocessing steps before further analysis or modeling. NeatText simplifies the text cleaning process and provides a convenient way to handle common text preprocessing tasks.

2. Feature extraction algorithms

"Feature Extraction Using CountVectorizer: A Powerful Algorithm for Text Analysis"

In the context of text analysis and natural language processing (NLP), feature extraction plays a crucial role in converting textual data into numerical representations that machine learning algorithms can process. CountVectorizer is one such powerful algorithm commonly used for feature extraction in NLP tasks.

The algorithm's purpose is to transform a collection of text documents into a matrix of token counts, where each row represents a document, and each column represents a unique token (word or n-gram). CountVectorizer performs the following steps:

- 1. Tokenization: It breaks down the text into individual words or tokens. For example, the sentence "The cat is on the mat" would be tokenized into ['the', 'cat', 'is', 'on', 'the', 'mat'].
- 2. Vocabulary Creation: CountVectorizer builds a vocabulary of all unique tokens present in the corpus

(the collection of text documents). It assigns a unique index to each token and stores it as a dictionary-like structure.

- 3. Counting: The algorithm then counts the frequency of each token in each document and represents it as a numerical value. For instance, the sentence "The cat is on the mat" would be transformed into a vector [1, 1, 1, 1, 1, 1] because each token appears once.
- 4. Document-Term Matrix: Finally, CountVectorizer constructs a document-term matrix, where each row represents a document, and each column represents the count of a token in that document. This matrix is a numerical representation of the text data suitable for machine learning algorithms.

CountVectorizer offers various configurable parameters, such as removing stop words, applying lowercase normalization, adjusting n-gram ranges, or limiting the vocabulary size. These parameters allow customization based on the specific needs of the text analysis task.

By using CountVectorizer for feature extraction, you can effectively represent text documents as numerical vectors, capturing the frequency of each token. This enables the application of various machine learning algorithms for tasks such as sentiment analysis, text classification, topic modeling, and more.

3. Effective Text Classification Using Logistic Regression and Naive Bayes Algorithms

Machine learning algorithms are essential components in the field of natural language processing (NLP) for tasks such as sentiment analysis, document classification, and text categorization. Logistic Regression and Naive Bayes are two popular algorithms that excel in text classification tasks.

1. Logistic Regression:

Logistic Regression is a widely used algorithm for binary and multiclass classification. It is particularly effective when dealing with linearly separable data, making it a suitable choice for text classification. The algorithm learns a logistic function to model the relationship between the features (textual data) and the target class labels. It predicts the probability of each class and assigns the class with the highest probability as the final prediction.

2. Naive Bayes:

Naive Bayes is a probabilistic algorithm that relies on Bayes' theorem and the assumption of feature independence, hence the "naive" assumption. Despite its simplifying assumption, Naive Bayes performs remarkably well in many NLP tasks. It calculates the conditional probability of each class given the observed features (words or tokens) and selects the class with the highest probability as the prediction. Naive Bayes is known for its efficiency, scalability, and ability to handle high-dimensional data like text.

Combining Logistic Regression and Naive Bayes for text classification allows for complementary strengths. Logistic Regression can capture more complex relationships between features and class labels, while Naive Bayes excels in handling high-dimensional feature spaces and can be faster in training and prediction.

By leveraging these algorithms in text classification tasks, you can effectively categorize textual data into predefined classes or labels. These algorithms can be applied to various domains such as sentiment analysis, spam detection, topic classification, and more.

4. Evaluation metrics algorithms

"Robust Evaluation Metrics for Text Classification: Accuracy, Classification Report, and Confusion Matrix"

When evaluating the performance of text classification models, it is essential to employ appropriate evaluation metrics to assess their accuracy and effectiveness. Three commonly used evaluation metrics for text classification tasks are accuracy_score, classification_report, and confusion_matrix.

1. Accuracy Score:

Accuracy Score is a widely used metric that measures the overall correctness of the model's predictions. It calculates the proportion of correctly predicted instances (both true positives and true negatives) out of the total number of instances in the dataset. Accuracy Score provides a simple and intuitive measure of the model's performance, but it can be misleading if the classes are imbalanced.

2. Classification Report:

The Classification Report is a comprehensive evaluation metric that provides insights into the model's performance for each individual class in a multiclass classification problem. It includes several key metrics such as precision, recall, F1-score, and support. Precision measures the proportion of correctly predicted instances for a given class out of the total predicted instances for that class. Recall (also known as sensitivity) measures the proportion of correctly predicted instances for a given class out of the total actual instances of that class. F1-score is the harmonic mean of precision and recall, providing a balanced measure of the model's performance. Support represents the number of instances for each class in the dataset.

3. Confusion Matrix:

A Confusion Matrix is a tabular representation that visually summarizes the performance of a classification model. It provides a detailed breakdown of the model's predictions, showing the number of instances that were correctly or incorrectly classified for each class. The confusion matrix enables further analysis of the model's performance, such as identifying specific types of errors (e.g., false positives, false negatives) and assessing the performance across different classes.

By utilizing these evaluation metrics, you can effectively assess the performance and accuracy of your text classification models. Accuracy Score provides an overall measure of correctness, the Classification Report offers detailed insights for individual classes, and the Confusion Matrix allows for a more granular analysis of the model's predictions. These metrics collectively help in understanding the strengths and weaknesses of the text classification model and guide improvements in its performance.

Outcome:

The expected outcome of this project is an automated system capable of accurately analyzing the emotions expressed in textual comments and feedback. Users will be able to input their text, and the system will provide insights into the sentiment and emotions conveyed in the text. The system's output may include the identified emotion categories, confidence scores, and visualizations if applicable.

```
[ ] # Make A Prediction
    ex1 = "This book was so interesting it made me happy"

pipe_lr.predict([ex1])

array(['joy'], dtype=object)
```

Output 1

```
# Make A Prediction
ex1 = "I get afraid from side effects of this cream"

pipe_lr.predict([ex1])
array(['fear'], dtype=object)
```

Output 2

```
# Make A Prediction
ex1 = "Customer scervice is bad"

pipe_lr.predict([ex1])
array(['sadness'], dtype=object)
```

Output 3

Exceptions considered:

During the development of the project, the following exceptions and challenges should be considered:

1. Handling out-of-vocabulary words or unknown emotions not present in the training data.

INTERNSHIP: PROJECT REPORT

- 2. Dealing with sarcasm, irony, or other linguistic nuances that may affect the accuracy of emotion analysis.
- 3. Addressing potential biases in the training data that could lead to skewed results.
- 4. Ensuring the system's scalability and performance, especially when handling a large volume of text data.

Enhancement Scope:

Some potential enhancements for this project include:

- 1. Multilingual Emotion Analysis: Extend the system to support multiple languages, enabling emotion analysis for text in different languages.
- 2. Real-Time Analysis: Develop a mechanism to perform emotion analysis in real-time, enabling instant feedback on the emotional content of text data.
- 3. Fine-grained Emotion Analysis: Expand the emotion categories to include more nuanced emotions or sub-categories for a more detailed analysis.
- 4. User Feedback Incorporation: Implement a feedback loop to continually improve the emotion analysis model based on user feedback and corrections.

Link to Code and executable file:

Code File Link:

https://colab.research.google.com/drive/19Bcs3e-bfQaluGCno0Rq_uBCHTDkl2DK?usp=sharing

Dataset File Link:

https://drive.google.com/file/d/1pKrP68WchLhTVlydTP9ILpdmtQzMOhW6/view?usp=sharing