**1.** Outline the key steps involved in developing a subject extraction algorithm using Python.

**A. Key Steps in Developing a Subject Extraction Algorithm Using Python:**

- Data Loading: Load the provided dataset into a pandas DataFrame.
- Data Pre-processing: Clean and pre-process the text data.
- Model Selection: Choose a machine learning or deep learning model for subject classification.
- Model Training: Train the model with the pre-processed dataset.
- Model Evaluation: Evaluate the model's performance using metrics like accuracy, precision, recall, and F1-score.
- Model Deployment: Deploy the model for subject extraction on new data.

**2.** Describe the structure and format of the sample dataset required for subject extraction.

**A. Structure and Format of the Sample Dataset**

The sample dataset should be structured in a CSV format with the following columns:

- text: The textual data or comments containing the sentences or paragraphs.
- subject: The subject label associated with each comment, categorized into "News", "politics", "Government News", "US-News", "left-news", and "Middle east".

**3.** Implement the Python code to read and pre-process the sample dataset for subject analysis. Ensure that the code correctly handles text data and labels.

**A. Python Code to Read and Pre-process the Sample Dataset:**

```
import pandas as pd

import re

from sklearn.model_selection import train_test_split

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.preprocessing import LabelEncoder

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import classification_report


url = "https://drive.google.com/uc?id=1EdI_HyUeI_Fi2nld7rQnnGEpQqn_BwM-"

df = pd.read_csv(url)

print(df.head())


def preprocess_text(text):

    text = text.lower()

    text = re.sub(r'[^\w\s]', '', text)
```

```
    text = re.sub(r'\d+', '', text)

    text = re.sub(r'\s+', ' ', text)

    return text


df['text'] = df['text'].apply(preprocess_text)


label_encoder = LabelEncoder()

df['subject'] = label_encoder.fit_transform(df['subject'])
```

**4.** Discuss the process of classifying subjects into the specified categories: "News", "politics", "Government News", "US-News", "left-news", and "Middle east". Explain any techniques or algorithms employed for this classification task.

## A. Classifying Subjects into Specified Categories:

# splitting the dataset into training and testing sets

X_train, X_test, y_train, y_test = train_test_split(df['text'], df['subject'], test_size=0.2, random_state=42)

# converting text data to numerical data using TF-IDF

tfidf_vectorizer = TfidfVectorizer(max_features=5000)

X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)

X_test_tfidf = tfidf_vectorizer.transform(X_test)

# training logistic regression model

model = LogisticRegression(max_iter=1000)

model.fit(X_train_tfidf, y_train)


**5.** Evaluate the effectiveness of the subject extraction algorithm on the provided sample dataset. Consider metrics such as accuracy, precision, recall, and F1-score.

## A. Evaluate the Effectiveness of the Subject Extraction Algorithm:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Government News | 0.28 | 0.07 | 0.12 | 316 |
| Middle-east | 0.14 | 0.13 | 0.13 | 159 |
| News | 0.91 | 0.97 | 0.94 | 1821 |
| US_News | 0.11 | 0.09 | 0.10 | 160 |
| left-news | 0.26 | 0.16 | 0.20 | 897 |
| politics | 0.47 | 0.65 | 0.55 | 1344 |

| | | | | |
|---|---|---|---|---|
| accuracy | | | 0.61 | 4697 |
| macro avg | 0.36 | 0.35 | 0.34 | 4697 |
| weighted avg | 0.57 | 0.61 | 0.58 | 4697 |

**6.** Propose potential enhancements or modifications to improve the performance of the sentiment extraction algorithm. Justify your recommendations.

**A. Potential Enhancements:**

To improve the performance, consider the following:

- Hyperparameter Tuning: Optimize model parameters using techniques like grid search or random search.
- Advanced Models: Use deep learning models like LSTM, BERT, or transformers.
- Data Augmentation: Increase the dataset size using data augmentation techniques.

**7.** Reflect on the ethical considerations associated with sentiment analysis, particularly regarding privacy, bias, and potential misuse of extracted sentiments.

**A. Ethical Considerations:**

- Privacy: Ensure that data collection respects user privacy.
- Bias: Address potential biases in the dataset and model.
- Misuse: Be aware of how analysis results are used to avoid misuse.

**8.** Write a complete code for this assignment.

**A. Complete code:**

```
import pandas as pd

import re

from sklearn.model_selection import train_test_split

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.preprocessing import LabelEncoder

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import classification_report


url = "https://drive.google.com/uc?id=1EdI_HyUeI_Fi2nld7rQnnGEpQqn_BwM-"

df = pd.read_csv(url)


print(df.head())
```

```python
def preprocess_text(text):
    text = text.lower()
    text = re.sub(r'[^\w\s]', '', text)
    text = re.sub(r'\d+', '', text)
    text = re.sub(r'\s+', ' ', text)
    return text


df['text'] = df['text'].apply(preprocess_text)


label_encoder = LabelEncoder()
df['subject'] = label_encoder.fit_transform(df['subject'])


X_train, X_test, y_train, y_test = train_test_split(df['text'], df['subject'], test_size=0.2, random_state=42)


tfidf_vectorizer = TfidfVectorizer(max_features=5000)
X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
X_test_tfidf = tfidf_vectorizer.transform(X_test)


model = LogisticRegression(max_iter=1000)
model.fit(X_train_tfidf, y_train)


y_pred = model.predict(X_test_tfidf)


print(classification_report(y_test, y_pred, target_names=label_encoder.classes_))
```