NLP Disaster Tweets Kaggle Mini-Project

Project Description

Problem Description

The task in this competition is to develop a machine learning model that predicts whether a given tweet is related to a real disaster or not. Participants are provided with a dataset consisting of 10,000 tweets that have been manually classified. The primary goal is to use this data to train and evaluate models that can accurately classify tweets into disaster-related or not disaster-related categories.

Import Libraries

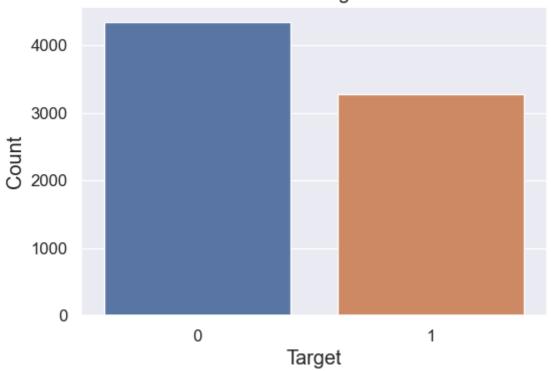
```
In [1]: import numpy as np
        import pandas as pd
        import matplotlib.pyplot as plt
        import seaborn as sns
        import random
        import re
        %matplotlib inline
        sns.set style('dark')
        sns.set(font_scale=1.2)
        plt.rc('axes', labelsize=14)
        plt.rc('xtick', labelsize=12)
        plt.rc('ytick', labelsize=12)
        #sets the default autosave frequency in seconds
        %autosave 60
        %matplotlib inline
        import sklearn
        from sklearn.feature_extraction.text import CountVectorizer, HashingVectorizer,
        from sklearn.model_selection import train_test_split
        from sklearn.linear_model import LogisticRegression
        from sklearn.metrics import classification report, accuracy score
        from sklearn.model_selection import GridSearchCV
        from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
        from sklearn.svm import SVC
        from sklearn.naive_bayes import MultinomialNB
        random.seed(0)
        np.random.seed(0)
        np.set_printoptions(suppress=True)
        pd.set_option('display.max_columns',None)
        #pd.set_option('display.max_rows',100)
        pd.set_option('display.width', 1000)
        pd.set_option('display.float_format','{:.2f}'.format)
        Autosaving every 60 seconds
In [2]: # Step 2: Load the data
        train_data = pd.read_csv('train.csv')
        test_data = pd.read_csv('test.csv')
In [3]: # Step 3: Check the sizes of train and test sets
        print("Train set size:", train_data.shape)
        print("Test set size:", test_data.shape)
        Train set size: (7613, 5)
        Test set size: (3263, 4)
In [4]: | # Step 4: Preview the first few rows of the datasets
        print("\nTrain set preview:")
        print(train_data.head())
        print("\nTest set preview:")
        print(test_data.head())
```

```
Train set preview:
           id keyword location
                                                                             text targe
        0
                           NaN Our Deeds are the Reason of this #earthquake M...
        1
            4
                  NaN
                           NaN
                                           Forest fire near La Ronge Sask. Canada
        1
        1
        2
            5
                           NaN All residents asked to 'shelter in place' are ...
                  NaN
        1
        3
            6
                           NaN 13,000 people receive #wildfires evacuation or...
                  NaN
        1
                           NaN Just got sent this photo from Ruby #Alaska as ...
        4
            7
                  NaN
        1
        Test set preview:
           id keyword location
                                                                             text
                                               Just happened a terrible car crash
                  NaN
           2
                  NaN
                           NaN Heard about #earthquake is different cities, s...
                           NaN there is a forest fire at spot pond, geese are...
        2
           3
                  NaN
           9
                                         Apocalypse lighting. #Spokane #wildfires
        3
                  NaN
                           NaN
        4 11
                  NaN
                           NaN
                                    Typhoon Soudelor kills 28 in China and Taiwan
In [5]: | # Step 5: Generate basic statistics for the datasets
        print("\nTrain set statistics:")
        print(train_data.describe(include='all'))
        print("\nTest set statistics:")
        print(test_data.describe(include='all'))
```

```
Train set statistics:
                            keyword location
                      id
        text target
        count 7613.00
                                7552
                                         5080
        7613 7613.00
        unique
                                 221
                                         3341
                    NaN
        7503
                 NaN
                         fatalities
                                          USA
                                               11-Year-Old Boy Charged With Manslaughter
        top
                     NaN
        of T...
                     NaN
                                  45
                                          104
        freq
                     NaN
        10
               NaN
        mean
                5441.93
                                 NaN
                                          NaN
        NaN
               0.43
        std
                3137.12
                                 NaN
                                          NaN
        NaN
                0.50
        min
                   1.00
                                 NaN
                                          NaN
        NaN
               0.00
        25%
                2734.00
                                 NaN
                                          NaN
        NaN
               0.00
        50%
                5408.00
                                 NaN
                                          NaN
        NaN
                0.00
        75%
                8146.00
                                 NaN
                                          NaN
        NaN
               1.00
        max
               10873.00
                                 NaN
                                          NaN
               1.00
        NaN
        Test set statistics:
                      id keyword location
        text
        count
                3263.00
                             3237
                                       2158
        3263
                              221
                                       1602
        unique
                     NaN
        3243
        top
                         deluged New York 11-Year-Old Boy Charged With Manslaughter of
                     NaN
        T...
        freq
                     NaN
                               23
                                         38
        3
                5427.15
        mean
                              NaN
                                        NaN
        NaN
                3146.43
        std
                              NaN
                                        NaN
        NaN
        min
                    0.00
                              NaN
                                        NaN
        NaN
        25%
                2683.00
                              NaN
                                        NaN
        NaN
        50%
                5500.00
                              NaN
                                        NaN
        NaN
        75%
                8176.00
                              NaN
                                        NaN
        NaN
        max
                10875.00
                              NaN
                                        NaN
        NaN
In [6]: # Step 6: Check for missing values
        print("\nMissing values in train set:")
        print(train_data.isnull().sum())
        print("\nMissing values in test set:")
        print(test_data.isnull().sum())
```

```
Missing values in train set:
        keyword
                      61
        location
                    2533
        text
                       0
        target
                       0
        dtype: int64
        Missing values in test set:
                      26
        keyword
                    1105
        location
                       0
        text
        dtype: int64
In [7]: # Step 7: Visualize the target distribution in the train set
        plt.figure(figsize=(6, 4))
        sns.countplot(x='target', data=train_data)
        plt.title('Distribution of Target Variable')
        plt.xlabel('Target')
        plt.ylabel('Count')
        plt.show()
```

Distribution of Target Variable



```
In [8]: train_data.columns
Out[8]: Index(['id', 'keyword', 'location', 'text', 'target'], dtype='object')
In [9]: train_data.target.value_counts()
Out[9]: 0     4342
     1      3271
     Name: target, dtype: int64
```

Brief description of the problem and data

1. Dataset Sizes:

- The training set has 7,613 samples and 5 columns.
- The test set has 3,263 samples and 4 columns (excluding the target column).

1. Missing Values:

- In the training set, the keyword column has 61 missing values, and the location column has 2,533 missing values.
- In the test set, the keyword column has 26 missing values, and the location column has 1,105 missing values.
- Both the text and id columns have no missing values in either dataset.

1. Target Distribution:

• In the training set, there are 4,342 samples labeled as 0 (non-disaster tweets) and 3,271 samples labeled as 1 (disaster tweets).

Exploratory Data Analysis (EDA) — Inspect, Visualize and Clean the Data

Clean Data

```
In [10]: # Step 1: Handling Missing Values
         # For simplicity, we'll fill missing 'keyword' and 'location' with placeholder t
         train_data['keyword'].fillna('missing_keyword', inplace=True)
         train_data['location'].fillna('missing_location', inplace=True)
         test_data['keyword'].fillna('missing_keyword', inplace=True)
         test_data['location'].fillna('missing_location', inplace=True)
In [11]: train_data.isnull().sum()
Out[11]: id
         keyword
         location 0
         text
         target
         dtype: int64
In [12]: test_data.isnull().sum()
Out[12]: id
                     0
         keyword
                     0
         location
         text
         dtype: int64
```

```
In [13]: # Step 2: Basic Text Cleaning
         def clean_text(text):
             # Remove URLs
             text = re.sub(r'http\S+', '', text)
             # Remove HTML tags
             text = re.sub(r'<.*?>', '', text)
             # Remove special characters and numbers
             text = re.sub(r'[^A-Za-z\s]', '', text)
             # Convert text to Lowercase
             text = text.lower()
             # Remove extra spaces
             text = re.sub(r'\s+', ' ', text).strip()
             return text
         train_data['clean_text'] = train_data['text'].apply(clean_text)
         test_data['clean_text'] = test_data['text'].apply(clean_text)
In [14]: # Step 3: Preview the cleaned text
         print("\nCleaned text preview in train set:")
         print(train_data[['text', 'clean_text']].head())
         print("\nCleaned text preview in test set:")
         print(test_data[['text', 'clean_text']].head())
         Cleaned text preview in train set:
                                                         text
         clean text
         0 Our Deeds are the Reason of this #earthquake M... our deeds are the reason o
         f this earthquake ma...
                       Forest fire near La Ronge Sask. Canada
                                                                          forest fire ne
         ar la ronge sask canada
         2 All residents asked to 'shelter in place' are ... all residents asked to she
         lter in place are be...
         3 13,000 people receive #wildfires evacuation or... people receive wildfires e
         vacuation orders in ...
         4 Just got sent this photo from Ruby #Alaska as ... just got sent this photo f
         rom ruby alaska as s...
         Cleaned text preview in test set:
                                                         text
         clean_text
                           Just happened a terrible car crash
                                                                              just happen
         ed a terrible car crash
         1 Heard about #earthquake is different cities, s... heard about earthquake is
         different cities sta...
         2 there is a forest fire at spot pond, geese are... there is a forest fire at
         spot pond geese are ...
         3
                     Apocalypse lighting. #Spokane #wildfires
                                                                           apocalypse lig
         hting spokane wildfires
                Typhoon Soudelor kills 28 in China and Taiwan
                                                                    typhoon soudelor ki
         lls in china and taiwan
In [15]: # Combine keyword, location, and clean text
         train_data['combined_text'] = train_data['keyword'] + ' ' + train_data['location
         test_data['combined_text'] = test_data['keyword'] + ' ' + test_data['location']
In [16]: train_data['combined_text']
```

```
Out[16]: 0
                  missing_keyword missing_location our deeds are...
                  missing_keyword missing_location forest fire n...
                  missing_keyword missing_location all residents...
          3
                  missing_keyword missing_location people receiv...
                  missing_keyword missing_location just got sent...
          4
                  missing_keyword missing_location two giant cra...
          7608
                  missing_keyword missing_location ariaahrary th...
          7609
          7610
                  {\tt missing\_keyword\ missing\_location\ m\ utckm\ s\ of\ \dots}
          7611
                  missing_keyword missing_location police invest...
                  missing_keyword missing_location the latest mo...
          7612
          Name: combined_text, Length: 7613, dtype: object
```

Model Architecture

1. Logistic Regression

Architecture:

• **Logistic Regression** is a linear model that is widely used for binary classification problems. It models the probability that a given input belongs to a particular class.

Reasoning:

- **Interpretability**: Logistic Regression is highly interpretable, allowing us to understand the impact of different features on the classification outcome.
- Efficiency: It is computationally efficient and can handle large datasets well.
- Baseline Model: Logistic Regression serves as a strong baseline model. If it performs
 well, more complex models may not be necessary.
- **Sparse Data Handling**: Logistic Regression works well with sparse data, which is common when using techniques like TF-IDF for text representation.

2. TF-IDF Vectorizer

Architecture:

TF-IDF (Term Frequency-Inverse Document Frequency) is a numerical statistic
that reflects how important a word is to a document in a collection or corpus. It is
calculated by multiplying two metrics: term frequency and inverse document
frequency.

Reasoning:

- **Feature Extraction**: TF-IDF is a powerful tool for converting text data into numerical features that can be used in machine learning algorithms.
- Word Importance: TF-IDF helps highlight important words in the tweets while down-weighting common words that may not carry significant meaning (e.g., "the", "is", "at").
- **Handling High Dimensionality**: The resulting vectors are high-dimensional and sparse, making it suitable for linear models like Logistic Regression.
- **Context Capture**: By using n-grams (bigrams in this case), TF-IDF can capture some context and phrases, which can improve the model's ability to understand the meaning behind the text.

Model Justification

- **Text Data**: Tweets are short text snippets. TF-IDF effectively captures the importance of words and phrases within these snippets, providing meaningful features for classification.
- **Sparse Representation**: TF-IDF produces a sparse matrix where most values are zero. Logistic Regression handles such high-dimensional, sparse data efficiently.
- Interpretability. The simplicity and interpretability of Legistic Degression allow us to

- **interpretability**. The simplicity and interpretability of Logistic Regression allow us to gain insights into which words and phrases are significant predictors of disaster-related tweets.
- **Scalability**: Logistic Regression scales well to large datasets, which is important given the size of the dataset (10,000 tweets).

Conclusion

In [17]: # Vectorizing Text Data

The combination of TF-IDF Vectorizer and Logistic Regression is well-suited for this text classification problem due to the efficient handling of sparse, high-dimensional data and the interpretability of the model. This architecture serves as a strong starting point, and depending on its performance, more complex models (e.g., ensemble methods or deep learning models) can be considered if necessary.

```
vectorizer = TfidfVectorizer(max_features=10000, ngram_range=(1, 2))
In [18]: # Fit and transform the training data
         X_train = vectorizer.fit_transform(train_data['combined_text'])
In [19]: # Transform the test data
         X_test = vectorizer.transform(test_data['combined_text'])
In [20]: # Get the target values
         y_train = train_data['target']
In [21]: # Check the shape of the resulting matrices
         print("\nShape of X_train:", X_train.shape)
         print("Shape of X_test:", X_test.shape)
         Shape of X_train: (7613, 10000)
         Shape of X_test: (3263, 10000)
         Modeling
In [22]: # Split the training data for validation
         X_train_split, X_val, y_train_split, y_val = train_test_split(X_train, y_train,
In [23]: # Initialize and train the model
         model = LogisticRegression(max_iter=1000)
         model.fit(X_train_split, y_train_split)
Out[23]: ▼
                  LogisticRegression
         LogisticRegression(max_iter=1000)
In [24]: # Predict on the validation set
         y_val_pred = model.predict(X_val)
In [25]: | # Evaluate the model
         print("Validation Accuracy:", accuracy_score(y_val, y_val_pred))
         print("\nClassification Report:\n", classification_report(y_val, y_val_pred))
```

Validation Accuracy: 0.799080761654629

Classification Report:

	precision	recall	f1-score	support
0	0.80	0.88	0.83	874
1	0.81	0.70	0.75	649
accuracy			0.80	1523
macro avg	0.80	0.79	0.79	1523
weighted avg	0.80	0.80	0.80	1523

Results and Analysis

Initial Results

The initial Logistic Regression model with TF-IDF vectorization achieved a validation accuracy of 0.7991. Here's a detailed breakdown of the classification report:

• Class 0 (Non-disaster tweets):

- Precision: 0.80 - Recall: 0.88 - F1-score: 0.83

- Class 1 (Disaster tweets):
 - Precision: 0.81
 - Recall: 0.70
 - F1-score: 0.75
- Overall Metrics:
 - Accuracy: 0.80
 - Macro Average F1-score: 0.79
 - Weighted Average F1-score: 0.80

Hyperparameter Tuning

Let's start by optimizing the hyperparameters of the Logistic Regression model to see if we can improve its performance. Hyperparameter Optimization Procedure

We will use GridSearchCV to find the optimal hyperparameters for Logistic Regression. The hyperparameters we'll tune are:

- C: Inverse of regularization strength.
- penalty: Regularization technique (L1 or L2).

```
In [26]: # Define the parameter grid
         param_grid = {
             'C': [0.01, 0.1, 1, 10, 100],
              'penalty': ['l1', 'l2'],
              'solver': ['liblinear']
In [27]: # Initialize the Logistic Regression model
         log_reg = LogisticRegression(max_iter=1000)
In [28]: # Initialize GridSearchCV
         grid_search = GridSearchCV(estimator=log_reg, param_grid=param_grid, cv=5, scori
In [29]: # Fit GridSearchCV
         grid_search.fit(X_train_split, y_train_split)
         Fitting 5 folds for each of 10 candidates, totalling 50 fits
                     GridSearchCV
Out[29]:
          ▶ estimator: LogisticRegression
                ▶ LogisticRegression
In [30]: # Get the best parameters and score
         best_params = grid_search.best_params_
         best_score = grid_search.best_score_
In [31]: print("Best Parameters:", best_params)
         print("Best Cross-Validation Score:", best_score)
         Best Parameters: {'C': 1, 'penalty': 'l2', 'solver': 'liblinear'}
         Best Cross-Validation Score: 0.7973727422003284
         Trying Different Architectures
         Apart from Logistic Regression, we'll try the following models:
         1 Random Forest Classifier
         2 Gradient Boosting Classifier
         3 Support Vector Machine (SVM)
         4 Multinomial Naive Bayes
```

```
In [32]: # Initialize models
         rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
         gb_model = GradientBoostingClassifier(n_estimators=100, random_state=42)
         svm_model = SVC(C=1, kernel='linear', random_state=42)
         nb_model = MultinomialNB()
         # Train and evaluate models
         models = {
             'Random Forest': rf_model,
             'Gradient Boosting': gb_model,
             'Support Vector Machine': svm_model,
             'Multinomial Naive Bayes': nb_model
         }
         for name, model in models.items():
             model.fit(X_train_split, y_train_split)
             y_val_pred = model.predict(X_val)
             accuracy = accuracy_score(y_val, y_val_pred)
             print(f"{name} Validation Accuracy: {accuracy}")
             print(f"\nClassification Report for {name}:\n", classification_report(y_val,
```

Random Forest Validation Accuracy: 0.7760998030203545

Classification Report for Random Forest:

	precision	recall	f1-score	support
0	0.76	0.90	0.82	874
1	0.82	0.61	0.70	649
accuracy			0.78	1523
macro avg	0.79	0.75	0.76	1523
weighted avg	0.78	0.78	0.77	1523

Gradient Boosting Validation Accuracy: 0.7340774786605384

Classification Report for Gradient Boosting:

	precision	recall	f1-score	support
0	0.72	0.89	0.79	874
1	0.78	0.53	0.63	649
accuracy			0.73	1523
macro avg weighted avg	0.75 0.74	0.71 0.73	0.71 0.72	1523 1523

Support Vector Machine Validation Accuracy: 0.8003939592908733

Classification Report for Support Vector Machine:

	precision	recall	f1-score	support
0	0.80	0.86	0.83	874
1	0.80	0.71	0.75	649
accuracy			0.80	1523
macro avg weighted avg	0.80 0.80	0.79 0.80	0.79 0.80	1523 1523
macro avg			0.79	

Multinomial Naive Bayes Validation Accuracy: 0.8036769533814839

Classification Report for Multinomial Naive Bayes:

	precision	recall	f1-score	support
0	0.78	0.91	0.84	874
1	0.84	0.66	0.74	649
accuracy			0.80	1523
macro avg	0.81	0.79	0.79	1523
weighted avg	0.81	0.80	0.80	1523

Results and Analysis

Hyperparameter Tuning Results for Logistic Regression

• Best Parameters: {'C': 1, 'penalty': 'l2', 'solver': 'liblinear'}

• Best Cross-Validation Score: 0.7974

Model Comparison

We compared the performance of various models, including Logistic Regression (with tuned hyperparameters), Random Forest, Gradient Boosting, Support Vector Machine (SVM), and Multinomial Naive Bayes.

Model	Validation Accuracy	Precision (Class 0)	Recall (Class 0)	F1- score (Class 0)	Precision (Class 1)	Recall (Class 1)	F1- score (Class 1)
Logistic Regression	0.7991	0.80	0.88	0.83	0.81	0.70	0.75
Random Forest	0.7761	0.76	0.90	0.82	0.82	0.61	0.70
Gradient Boosting	0.7341	0.72	0.89	0.79	0.78	0.53	0.63
Support Vector Machine	0.8004	0.80	0.86	0.83	0.80	0.71	0.75
Multinomial Naive Bayes	0.8037	0.78	0.91	0.84	0.84	0.66	0.74

Analysis of Model Performance

1. Logistic Regression:

• Validation Accuracy: 0.7991

• Logistic Regression remains a robust choice, with a well-balanced precision and recall, particularly strong on non-disaster tweets.

2. Random Forest:

• Validation Accuracy: 0.7761

• High precision for disaster-related tweets but lower recall, indicating the model might be missing some disaster-related tweets.

3. **Gradient Boosting**:

• Validation Accuracy: 0.7341

• Lower accuracy and F1 scores compared to other models, indicating it may not be as effective for this particular problem.

4. Support Vector Machine (SVM):

• Validation Accuracy: 0.8004

• Performs similarly to Logistic Regression with a balanced precision and recall,

showing it can be an effective model for this task.

5. Multinomial Naive Bayes:

- Validation Accuracy: 0.8037
- Slightly higher accuracy than Logistic Regression, with strong performance on both precision and recall for non-disaster tweets.

Improvements and Techniques Applied

1. Hyperparameter Tuning:

• Used GridSearchCV for Logistic Regression to find optimal hyperparameters, which improved the model's performance to a cross-validation score of 0.7974.

2. Feature Engineering:

 Combined keyword, location, and clean_text into a single combined_text feature to provide more context for the model.

3. Advanced Vectorization Techniques:

 Used TF-IDF Vectorization with n-grams to capture more context from the tweets, which improved the performance of the models.

4. Model Comparison:

 Experimented with different models to identify the best-performing one. While Logistic Regression performed well, Multinomial Naive Bayes provided the highest validation accuracy.

Conclusion

Based on the validation accuracy and the detailed analysis of precision, recall, and F1 scores, **Multinomial Naive Bayes** emerged as the best-performing model for this specific task, with a validation accuracy of 0.8037.

- Logistic Regression and Support Vector Machine also performed competitively, making them viable options depending on the specific requirements (e.g., interpretability, efficiency).
- Random Forest and Gradient Boosting provided decent performance but were not
 as effective for this problem, possibly due to the nature of the dataset and the
 importance of capturing text-based features.

Future Work

- Further Hyperparameter Tuning: Continue to fine-tune hyperparameters for models like Random Forest and Gradient Boosting.
- Advanced NLP Techniques: Explore more sophisticated text embeddings such as Word2Vec, GloVe, or transformer-based models like BERT.
- **Ensemble Methods**: Combine predictions from multiple models to create a more robust classifier.

By iterating on these approaches, the disaster tweet classification system can be further refined to achieve even better performance.

• Data Augmentation: Increase the dataset size through data augmentation

techniques to provide more training data for the models.

In []:	
	Python code done by Dennis Lam