

Natural Language Processing

HomeWork 3

Ilia Hashemi Rad

99102456

AmirMohammad Fakhimi

99170531

AmirMahdi Namjoo

402211467

Introduction

Explanations:

In this project, we have implemented a classification problem on sentiment analysis both on a document and text level on the Taaghche comments dataset. Taaghche is an Iranian online ebook marketplace featuring user comments on each book. Based on a dataset of scores and comments from users and other users' likes in the comments, we implemented a classifier for sentiment analysis of the text. We use both base models like LSTM or SVM and also Transformer-based models.

In addition, we also wrote a crawler Taaghche website to crawl all the book pages and get the book information like name and author names so that we can have a list of NER for them.

Crawlers

Explanations:

In the first part, we implemented a crawler on the Taaghche website to get all the book info. The crawler downloads the pages of all books on the website by iterating on the "id" in the URL and saving the result as HTML. We then feed the results into an extractor module written using BeautifulSoup to extract the main parts of book info, including name, publication, author, translator, etc., and save them into CSV files for further use.

crawler.py

```
1 base_url = "https://taaghche.com/book/"
2 def save_page(book_id, thread_exceptions):
3     url = f"{base_url}{book_id}/"
4     try:
5         response = requests.get(url)
6
7         if response.status_code == 404:
8             print(book_id, " : 404")
9             return
10
11         with open(os.path.join(output_dir, f"{book_id}.html"), 'w', encoding='utf-8') as f:
12             print(f'Saving book with id: {book_id}')
13             f.write(response.text)
14     except Exception as e:
15         print(e)
16     thread_exceptions.append(book_id)
```

Explanations:

This is the main part of the crawler that sends requests to Taaghche to get the book pages and save them into files. We also used Python threading features to make the whole process faster.

extractor.py

```
1 def extract_data_from_html(file_path):
2     with open(file_path, 'r', encoding='utf-8') as file:
3         soup = BeautifulSoup(file, 'html.parser')
4         script_tag = soup.find('script', type='application/ld+json')
5         if script_tag:
6             try:
7                 json_data = json.loads(script_tag.string)
8                 book_name = json_data.get('name', '')
9                 authors = ' $ '.join([author['name'] for author in json_data.get('author', [])])
```

```

10         translators = ' $ '.join(
11             [translator['name'] for translator in json_data.get('workExample', {}).get('
12                 translator', [])])
13         publisher = json_data.get('workExample', {}).get('publisher', {}).get('name', '')
14         data.append({
15             'name': book_name,
16             'author': authors,
17             'translator': translators,
18             'publisher': publisher
19         })
20     except json.JSONDecodeError:
21         pass
22
23     for x in os.listdir(input_dir):
24         file_name = x
25         file_path = os.path.join(input_dir, file_name)
26         if os.path.isfile(file_path) and file_path.endswith('.html'):
27             extract_data_from_html(file_path)
28         if len(data) >= 10000:
29             df = pd.DataFrame(data)
30             output_file = os.path.join(output_dir, f'books_data_part_{part_number}.csv')
31             df.to_csv(output_file, index=False, encoding='utf-8')
32             data = []
33             part_number += 1
34             print("index put into files: ", file_name)
35
36     if data:
37         df = pd.DataFrame(data)
38         output_file = os.path.join(output_dir, f'books_data_part_{part_number}.csv')
39         df.to_csv(output_file, index=False, encoding='utf-8')
40         data = []
41         part_number += 1
42         print("index put into files: ", file_name)
43         data = []

```

Explanations:

The main part of extracor.py is `extract_data_from_html` function. This function uses BS4 to find the JSON section that includes book data in the HTML and save data json data into a python dictionary. This data is then fed into a Pandas Dataframe and we save it into a csv file.

Document Classifier - Base Model

The base model is in `DocClassifier_Base.ipynb`. We investigate each segment in the following sections.

Note: Some parts of the explanations of code from here to the end of the document are written by partially ChatGPT to help us write clear and unambiguous descriptions.

Loading and Preparing Data

```

1 # Import the pandas library for data manipulation
2 import pandas as pd
3
4 # Load the CSV file into a pandas DataFrame
5 # The file 'taghche.csv' is located in the 'datasets/' directory
6 data = pd.read_csv('datasets/taghche.csv')
7
8 # Remove any duplicate rows in the DataFrame
9 data = data.drop_duplicates()
10
11 # Drop rows where the 'comment' or 'rate' columns have missing values (NaN)
12 data.dropna(subset=['comment', 'rate'], inplace=True)
13
14 # Print the first 5 rows of the cleaned DataFrame

```

```

15 print(data.head())
16
17
18 def label_sentiment(rate, positive_threshold, neutral_threshold):
19     """
20     Labels sentiment based on rating thresholds.
21
22     Args:
23     - rate (int or float): The numerical rating to evaluate.
24     - positive_threshold (int or float): The minimum rating value that qualifies as 'positive'.
25     - neutral_threshold (int or float): The minimum rating value that qualifies as 'neutral'; ratings
      below this are considered 'negative'.
26
27     Returns:
28     - str: The sentiment label ('positive', 'neutral', or 'negative') based on the rating.
29     """
30     # Check if the rating is greater than or equal to the positive threshold
31     if rate >= positive_threshold:
32         return 'positive'
33     # If the rating is not 'positive', check if it is greater than or equal to the neutral threshold
34     elif rate >= neutral_threshold:
35         return 'neutral'
36     # If the rating is neither 'positive' nor 'neutral', label the sentiment as 'negative'
37     else:
38         return 'negative'

```

Explanations:

At first, we load data using Pandas Library. We then define a function to label the sentiment of each comment based on its rating and threshold. Note that we have two thresholds, one for positive, and one for neutral comments. Everything below neutral is considered negative.

Balancing Data

```

1  # Function to prepare data and labels based on given thresholds
2  def prepare_data(positive_threshold, neutral_threshold):
3      """
4      Prepares the data and labels based on given thresholds for sentiment classification.
5
6      Args:
7      - positive_threshold (int or float): The minimum rating value that qualifies as 'positive'.
8      - neutral_threshold (int or float): The minimum rating value that qualifies as 'neutral';
        ratings below this are considered 'negative'.
9
10     Returns:
11     - tuple: A tuple containing:
12         - pandas.Series: The comments from the balanced dataset.
13         - pandas.Series: The corresponding sentiment labels from the balanced dataset.
14     """
15     # Create a copy of the original data to avoid modifying it
16     labeled_data = data.copy()
17
18     # Apply the label_sentiment function to the 'rate' column to create a new 'sentiment' column
19     labeled_data['sentiment'] = labeled_data['rate'].apply(lambda x: label_sentiment(x,
20         positive_threshold, neutral_threshold))
21
22     # Combine the 'comment' and 'sentiment' columns into a single DataFrame
23     df = pd.concat([labeled_data['comment'], labeled_data['sentiment']], axis=1)
24
25     # Separate the DataFrame into three classes based on sentiment
26     positive = df[df['sentiment'] == 'positive']
27     neutral = df[df['sentiment'] == 'neutral']
28     negative = df[df['sentiment'] == 'negative']
29
30     # Determine the size of the smallest class to balance the dataset

```

```

30 min_class_size = min(len(positive), len(neutral), len(negative))
31
32 # Downsample each class to the size of the smallest class to ensure balance
33 positive_downsampled = resample(positive, replace=False, n_samples=min_class_size, random_state
    =42)
34 neutral_downsampled = resample(neutral, replace=False, n_samples=min_class_size, random_state
    =42)
35 negative_downsampled = resample(negative, replace=False, n_samples=min_class_size, random_state
    =42)
36
37 # Combine the downsampled classes into a single DataFrame
38 df_balanced = pd.concat([positive_downsampled, neutral_downsampled, negative_downsampled])
39
40 # Shuffle the balanced DataFrame to mix the rows
41 df_balanced = df_balanced.sample(frac=1, random_state=42).reset_index(drop=True)
42
43 # Return the 'comment' and 'sentiment' columns as separate pandas Series
44 return df_balanced['comment'], df_balanced['sentiment']

```

Explanations:

- The function starts by creating a copy of the original dataset to avoid modifying it:

```

1 labeled_data = data.copy()

```

- It applies the `label_sentiment` function to the `rate` column to create a new `sentiment` column:

```

1 labeled_data['sentiment'] = labeled_data['rate'].apply(lambda x: label_sentiment(
    x, positive_threshold, neutral_threshold))

```

- The function then combines the `comment` and `sentiment` columns into a single DataFrame:

```

1 df = pd.concat([labeled_data['comment'], labeled_data['sentiment']], axis=1)

```

- It separates the DataFrame into three classes based on sentiment:

```

1 positive = df[df['sentiment'] == 'positive']
2 neutral = df[df['sentiment'] == 'neutral']
3 negative = df[df['sentiment'] == 'negative']

```

- The function determines the size of the smallest class to balance the dataset:

```

1 min_class_size = min(len(positive), len(neutral), len(negative))

```

- It down-samples each class to the size of the smallest class to ensure balance:

```

1 positive_downsampled = resample(positive, replace=False, n_samples=min_class_size
    , random_state=42)
2 neutral_downsampled = resample(neutral, replace=False, n_samples=min_class_size,
    random_state=42)
3 negative_downsampled = resample(negative, replace=False, n_samples=min_class_size
    , random_state=42)

```

- The function combines the down-sampled classes into a single DataFrame:

```

1 df_balanced = pd.concat([positive_downsampled, neutral_downsampled,
    negative_downsampled])

```

- It shuffles the balanced DataFrame to mix the rows:

```

1 df_balanced = df_balanced.sample(frac=1, random_state=42).reset_index(drop=True)

```

- Finally, the function returns the `comment` and `sentiment` columns as separate pandas Series:

```

1 return df_balanced['comment'], df_balanced['sentiment']

```

Preprocess

```
1 # Function to preprocess and normalize the text
2 def preprocess(text):
3     """
4     Preprocesses and normalizes text data by removing special characters,
5     non-Persian characters, digits, and multiple spaces.
6
7     Args:
8     - text (str): Input text to be processed.
9
10    Returns:
11    - str: Processed text with normalized format.
12    """
13    # Replace one or more newline characters with a single newline
14    pattern = re.compile(r"\n+")
15    text = pattern.sub("\n", text)
16
17    # Replace '\n' and '\n' with a single space
18    text = re.sub(r'\\n|\\n', ' ', text)
19
20    # Remove non-Persian characters and digits
21    text = re.sub(r'[^-\s]', ' ', text)
22
23    # Replace one or more spaces with a single space
24    pattern = re.compile(r" +")
25    text = pattern.sub(" ", text)
26
27    return text
28
29 # Apply the preprocess function to the 'comment' column in the DataFrame data
30 data['comment'] = data['comment'].apply(preprocess)
31 # Remove any duplicate rows in the DataFrame
32 data = data.drop_duplicates()
33
34 # Drop rows where the 'comment' or 'rate' columns have missing values (NaN)
35 data.dropna(subset=['comment', 'rate'], inplace=True)
```

Explanations:

- The function starts by replacing one or more newline characters with a single newline:

```
1 pattern = re.compile(r"\n+")
2 text = pattern.sub("\n", text)
```

- Next, it replaces occurrences of the newline character (`\n`) and escaped newline (`\\n`) with a single space:

```
1 text = re.sub(r'\\n|\\n', ' ', text)
```

- The function then removes all non-Persian characters and digits. This is done using a regular expression that matches any character not in the Persian alphabet (`-`) or whitespace:

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

- Finally, it replaces one or more spaces with a single space to normalize the spacing in the text:

```
1 pattern = re.compile(r" +")
2 text = pattern.sub(" ", text)
```

- The processed text is then returned by the function.

The following lines apply the `preprocess` function to the `comment` column of the DataFrame `data`:

```

1 # Apply the preprocess function to the 'comment' column in the DataFrame data
2 data['comment'] = data['comment'].apply(preprocess)
3
4 # Remove any duplicate rows in the DataFrame
5 data = data.drop_duplicates()
6
7 # Drop rows where the 'comment' or 'rate' columns have missing values (NaN)
8 data.dropna(subset=['comment', 'rate'], inplace=True)

```

- The `preprocess` function is applied to each entry in the `comment` column to clean and normalize the text.
- After preprocessing, any duplicate rows in the DataFrame are removed using:

```

1 data = data.drop_duplicates()

```

- Finally, rows where the `comment` or `rate` columns have missing values (NaN) are dropped:

```

1 data.dropna(subset=['comment', 'rate'], inplace=True)

```

TF IDF - Logistic Regression

Next, we implement a TF-IDF vectorizer and use logistic regression for the task. Based on our study and the papers we investigated, logistic regression is the best base model for this task.

```

1 # Create a pipeline with TF-IDF and logistic regression
2 logReg_PL = Pipeline([
3     ('tfidf', TfidfVectorizer()),
4     ('logreg', LogisticRegression(max_iter=500, solver='newton-cg'))
5 ])
6
7 # Define the parameter grid for GridSearchCV
8 param_grid = {
9     'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
10    'tfidf__max_features': [5000, 10000],
11    'logreg__C': [0.01, 0.1, 1, 10]
12 }
13
14 # Custom GridSearchCV implementation to iterate over parameter grid
15 best_score = 0
16 best_params = None
17
18 # Thresholds to evaluate
19 rate_thresholds = [(1, 2), (1, 3), (1, 4), (2, 3), (2, 4), (3, 4)]
20
21 # Iterate over each pair of thresholds and perform GridSearchCV
22 for neutral_threshold, positive_threshold in tqdm(rate_thresholds):
23     # Prepare the data using the specified thresholds
24     X_prepared, y_prepared = prepare_data(positive_threshold, neutral_threshold)
25
26     # Split the data into training and testing sets
27     X_train, X_test, y_train, y_test = train_test_split(X_prepared, y_prepared, test_size=0.1,
28                                                         random_state=42)
29
30     # Initialize GridSearchCV with the pipeline and parameter grid
31     grid_search = GridSearchCV(logReg_PL, param_grid, cv=5, scoring='accuracy')
32
33     # Fit GridSearchCV on the training data
34     grid_search.fit(X_train, y_train)
35
36     # Get the best score and parameters from GridSearchCV
37     score = grid_search.best_score_
38
39     # Update the best score and best parameters if the current score is better

```

```

39     if score > best_score:
40         best_score = score
41         best_params = grid_search.best_params_
42         best_params['positive_threshold'] = positive_threshold
43         best_params['neutral_threshold'] = neutral_threshold
44
45 # Print the best parameters found by GridSearchCV
46 print("Best parameters for TF-IDF model are:", best_params)
47
48 # Prepare the data using the best parameters found from GridSearchCV
49 X_prepared, y_prepared = prepare_data(best_params['positive_threshold'], best_params['
    neutral_threshold'])
50
51 # Split the prepared data into training and testing sets
52 X_train, X_test, y_train, y_test = train_test_split(X_prepared, y_prepared, test_size=0.1,
    random_state=42)
53
54 # Create a pipeline for the best Logistic Regression model with the best parameters
55 best_logReg_model = Pipeline([
56     ("tfidf", TfidfVectorizer(ngram_range=best_params['tfidf__ngram_range'], max_features=
    best_params['tfidf__max_features'])),
57     ("logreg", LogisticRegression(C=best_params['logreg__C'], max_iter=500, solver='newton-cg'))
58 ])
59
60 # Fit the best Logistic Regression model on the training data
61 best_logReg_model.fit(X_train, y_train)
62
63 # Predict the labels on the test set using the best model
64 y_test_pred = best_logReg_model.predict(X_test)
65
66 # Calculate the accuracy score of the best model on the test set
67 test_accuracy = accuracy_score(y_test, y_test_pred)
68
69 # Print the test accuracy score of the best Logistic Regression model
70 print("Test accuracy of Logistic Regression model:", test_accuracy)

```

Explanations:

```

1     logReg_PL = Pipeline([
2         ("tfidf", TfidfVectorizer()),
3         ("logreg", LogisticRegression(max_iter=500, solver='newton-cg'))
4     ])

```

- The pipeline logReg_PL is created with two steps:

1. TfidfVectorizer(): Converts text data into TF-IDF features.
2. LogisticRegression(): Applies logistic regression for classification, with a maximum of 500 iterations and the 'newton-cg' solver.

```

1     # Define the parameter grid for GridSearchCV
2     param_grid = {
3         'tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
4         'tfidf__max_features': [5000, 10000],
5         'logreg__C': [0.01, 0.1, 1, 10]
6     }

```

- The param_grid defines the hyperparameters for GridSearchCV to search over:

- tfidf__ngram_range: N-gram ranges (unigrams, bigrams, trigrams).
- tfidf__max_features: Maximum number of features (5000 or 10000).
- logreg__C: Inverse of regularization strength (0.01, 0.1, 1, 10).


```

1 # Custom GridSearchCV implementation to iterate over parameter grid
2 best_score = 0
3 best_params = None
4
5 # Thresholds to evaluate
6 rate_thresholds = [(1, 2), (1, 3), (1, 4), (2, 3), (2, 4), (3, 4)]

```

- `best_score` and `best_params` are initialized to store the best score and corresponding parameters.
- `rate_thresholds` contains pairs of thresholds to evaluate for neutral and positive sentiment classification.

```

1 # Iterate over each pair of thresholds and perform GridSearchCV
2 for neutral_threshold, positive_threshold in tqdm(rate_thresholds):
3     # Prepare the data using the specified thresholds
4     X_prepared, y_prepared = prepare_data(positive_threshold, neutral_threshold)
5
6     # Split the data into training and testing sets
7     X_train, X_test, y_train, y_test = train_test_split(X_prepared, y_prepared,
8                                                         test_size=0.1, random_state=42)
9
10    # Initialize GridSearchCV with the pipeline and parameter grid
11    grid_search = GridSearchCV(logReg_PL, param_grid, cv=5, scoring='accuracy')
12
13    # Fit GridSearchCV on the training data
14    grid_search.fit(X_train, y_train)
15
16    # Get the best score and parameters from GridSearchCV
17    score = grid_search.best_score_
18
19    # Update the best score and best parameters if the current score is better
20    if score > best_score:
21        best_score = score
22        best_params = grid_search.best_params_
23        best_params['positive_threshold'] = positive_threshold
24        best_params['neutral_threshold'] = neutral_threshold
25
26    # Print the best parameters found by GridSearchCV
27    print("Best parameters for TF-IDF model are:", best_params)

```

- The code iterates over each pair of thresholds in `rate_thresholds`.
- For each pair:
 1. `prepare_data` is called to prepare the dataset with the current thresholds.
 2. The data is split into training and testing sets using `train_test_split`.
 3. `GridSearchCV` is initialized with the pipeline and parameter grid, and fitted to the training data.
 4. The best score and parameters are retrieved from the grid search results.
 5. If the current score is better than the best score, update `best_score` and `best_params`.

```

1 # Prepare the data using the best parameters found from GridSearchCV
2 X_prepared, y_prepared = prepare_data(best_params['positive_threshold'], best_params['
3     neutral_threshold'])
4
5 # Split the prepared data into training and testing sets
6 X_train, X_test, y_train, y_test = train_test_split(X_prepared, y_prepared, test_size
7     =0.1, random_state=42)
8
9 # Create a pipeline for the best Logistic Regression model with the best parameters
10 best_logReg_model = Pipeline([
11     ("tfidf", TfidfVectorizer(ngram_range=best_params['tfidf__ngram_range'],
12                             max_features=best_params['tfidf__max_features'])),

```

```

0         ("logreg", LogisticRegression(C=best_params['logreg__C'], max_iter=500, solver='
1         newton-cg'))
2
3     # Fit the best Logistic Regression model on the training data
4     best_logReg_model.fit(X_train, y_train)
5
6     # Predict the labels on the test set using the best model
7     y_test_pred = best_logReg_model.predict(X_test)
8
9     # Calculate the accuracy score of the best model on the test set
10    test_accuracy = accuracy_score(y_test, y_test_pred)
11
12    # Print the test accuracy score of the best Logistic Regression model
13    print("Test accuracy of Logistic Regression model:", test_accuracy)

```

- The data is prepared using the best parameters found by GridSearchCV.
- The prepared data is split into training and testing sets.
- A pipeline is created for the best logistic regression model with the best parameters.
- The model is fitted to the training data.
- Predictions are made on the test set.
- The accuracy of the model is calculated on the test set.
- The test accuracy is printed.

Evaluation Metrics

```

1  def evaluate_model(y_true, y_pred, class_names):
2      """
3      Evaluates the performance of a classification model using various metrics and visualizations.
4
5      Args:
6      - y_true (array-like): True labels of the data.
7      - y_pred (array-like): Predicted labels of the data.
8      - class_names (list): List of class names in the same order as the confusion matrix.
9
10     Returns:
11     - pd.DataFrame: DataFrame containing the classification report.
12     """
13     # Generate and print the classification report
14     report = classification_report(y_true, y_pred, target_names=class_names, output_dict=True)
15     report_df = pd.DataFrame(report).transpose()
16     print("Classification Report:\n", report_df)
17
18     # Generate and display the confusion matrix as a heatmap
19     cm = confusion_matrix(y_true, y_pred)
20     plt.figure(figsize=(10, 7))
21     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=class_names, yticklabels=
22         class_names)
23     plt.xlabel('Predicted')
24     plt.ylabel('True')
25     plt.title('Confusion Matrix')
26     plt.show()
27
28     # Calculate and print overall metrics: accuracy, precision, recall, and F1 score
29     accuracy = accuracy_score(y_true, y_pred)
30     precision = precision_score(y_true, y_pred, average='weighted')
31     recall = recall_score(y_true, y_pred, average='weighted')
32     f1 = f1_score(y_true, y_pred, average='weighted')

```

```

32
33 metrics = {
34     "Accuracy": accuracy,
35     "Precision": precision,
36     "Recall": recall,
37     "F1 Score": f1
38 }
39
40 print("\nOverall Metrics:")
41 for metric, value in metrics.items():
42     print(f"{metric}: {value:.4f}")
43
44 return report_df
45
46 # Evaluate the model
47 report_df = evaluate_model(y_test, y_test_pred, ["Negative", "Neutral", "Positive"])

```

Explanations:

The `evaluate_model` function evaluates the performance of a classification model using various metrics and visualizations.

- The function takes three arguments:
 - `y_true`: The true labels of the data.
 - `y_pred`: The predicted labels of the data.
 - `class_names`: A list of class names in the same order as the confusion matrix.
- The function returns a pandas DataFrame containing the classification report.

Classification Report

- The classification report is generated using `classification_report` from scikit-learn and printed:

```

1 report = classification_report(y_true, y_pred, target_names=class_names, output_dict=
  True)
2 report_df = pd.DataFrame(report).transpose()
3 print("Classification Report:\n", report_df)

```

Confusion Matrix

- The confusion matrix is generated and displayed as a heatmap using seaborn:

```

1 cm = confusion_matrix(y_true, y_pred)
2 plt.figure(figsize=(10, 7))
3 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=class_names,
4             yticklabels=class_names)
5 plt.xlabel('Predicted')
6 plt.ylabel('True')
7 plt.title('Confusion Matrix')
8 plt.show()

```

Overall Metrics

- The function calculates and prints overall metrics including accuracy, precision, recall, and F1 score:

```

1 accuracy = accuracy_score(y_true, y_pred)
2 precision = precision_score(y_true, y_pred, average='weighted')
3 recall = recall_score(y_true, y_pred, average='weighted')
4 f1 = f1_score(y_true, y_pred, average='weighted')
5
6 metrics = {
7     "Accuracy": accuracy,

```

```

8         "Precision": precision,
9         "Recall": recall,
10        "F1 Score": f1
11    }
12
13    print("\nOverall Metrics:")
14    for metric, value in metrics.items():
15        print(f"{metric}: {value:.4f}")

```

- These metrics are printed in a readable format.

Function Return

- The function returns the classification report DataFrame:

```

1    return report_df

```

Model Evaluation

- The `evaluate_model` function is called to evaluate the model:

```

1    report_df = evaluate_model(y_test, y_test_pred, ["Negative", "Neutral", "Positive"])

```

0.0.1 Results and Interesting Notes

For the results section, we analyzed different approaches to the data to check whether we could get a better result. At first, we tested it without any special technique and got an accuracy of about 0.61 on the whole data. Another approach we tested is removing stopwords, but it seems stopwords are important for our task, and after removing them, our accuracy dropped.

Another approach was to use a formalizer to change the text into a formal Persian text. We used a formalizer based on T5 using models from previous semesters, but we got worse results. It is predictable, though, because some informal Persian slang can significantly change the sentiment of a sentence, and they lose their meaning even from a human perspective when we formalize them.

We also used another approach for formalization: giving the Persian text to Google Translate to translate it into English and then back to Persian. It didn't work fine either.

One of the things we find in papers that seem good for this task is a lemmatizer. We tried to use it on our data, but unfortunately, Hazm was very slow on our 70K database, so we could not use this technique for our task.

It seems our result is good enough, though. As it is a multiclass classification and the best results are for ParsBERT on a significantly larger database than ours that is near 70%, it seems our accuracy of 61% is good enough.

The code used for formalization based on T5 is as follows:

```

1    from transformers import (T5ForConditionalGeneration, AutoTokenizer, pipeline)
2    import torch
3
4    model = T5ForConditionalGeneration.from_pretrained('parsi-ai-nlpclass/PersianTextFormalizer')
5    tokenizer = AutoTokenizer.from_pretrained('parsi-ai-nlpclass/PersianTextFormalizer')
6
7    pipe = pipeline(task='text2text-generation', model=model, tokenizer=tokenizer)
8    def formalizer(text):
9        device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
10       model.to(device)
11       formalized = ""
12       for line in text.splitlines():
13
14           inputs = tokenizer.encode("informal: " + line, return_tensors='pt', max_length=128,
15                                   truncation=True, padding='max_length')
16           inputs = inputs.to(device)
17
18           outputs = model.generate(inputs, max_length=128, num_beams=4)

```

```

18         formalized = formalized + tokenizer.decode(outputs[0], skip_special_tokens=True) + " "
19     return formalized
20
21 data['comment'] = data['comment'].apply(formalizer)

```

The code used for formalization based on Google Translate is as follows:

```

1 import googletrans as gt
2 translator = gt.Translator()
3
4 def translation(text):
5     en_translated = translator.translate(text, 'en', 'fa')
6     fa_translated = translator.translate(en_translated.text, 'fa', 'en')
7     return fa_translated.text
8
9 data['comment'] = data['comment'].apply(translation)

```

Figure 1: Results of running Logistic Regression on 10000 rows of dataset - With Formalization

Best parameters for TF-IDF model are: {'logreg_C': 1, 'tfidf_max_features': 10000, 'tfidf_ngram_range': (1, 2), 'positive_threshold': 4, 'neutral_threshold': 2}

Test accuracy of Logistic Regression model: 0.5208333333333334

	precision	recall	f1-score	support
negative	0.48	0.49	0.49	117
neutral	0.56	0.55	0.55	137
positive	0.52	0.52	0.52	130
accuracy			0.52	384
macro avg	0.52	0.52	0.52	384
weighted avg	0.52	0.52	0.52	384

Figure 2: Results of running Logistic Regression on 10000 rows of dataset - Without Formalization

Best parameters for TF-IDF model are: {'logreg_C': 0.1, 'tfidf_max_features': 5000, 'tfidf_ngram_range': (1, 3), 'positive_threshold': 4, 'neutral_threshold': 2}

Test accuracy of Logistic Regression model: 0.5416666666666667

	precision	recall	f1-score	support
negative	0.47	0.56	0.52	117
neutral	0.54	0.50	0.52	137
positive	0.62	0.57	0.59	130
accuracy			0.54	384
macro avg	0.55	0.54	0.54	384
weighted avg	0.55	0.54	0.54	384

Figure 3: Results of running Logistic Regression on 10000 rows of dataset - without stopwords

Best parameters for TF-IDF model are: {'logreg_C': 0.1, 'tfidf_max_features': 10000, 'tfidf_ngram_range': (1, 1), 'positive_threshold': 4, 'neutral_threshold': 2}

Test accuracy of Logistic Regression model: 0.4719541666666667

	precision	recall	f1-score	support
negative	0.40	0.60	0.48	117
neutral	0.47	0.42	0.44	137
positive	0.63	0.41	0.50	130
accuracy			0.47	384
macro avg	0.50	0.48	0.47	384
weighted avg	0.50	0.47	0.47	384

+ Code + Markdown

Figure 4: Results of running Logistic Regression on 70000 rows of dataset - without stopwords

Best parameters for TF-IDF model are: {'logreg_C': 0.1, 'tfidf_max_features': 10000, 'tfidf_ngram_range': (1, 2), 'positive_threshold': 4, 'neutral_threshold': 2}

Test accuracy of Logistic Regression model: 0.540631329203545

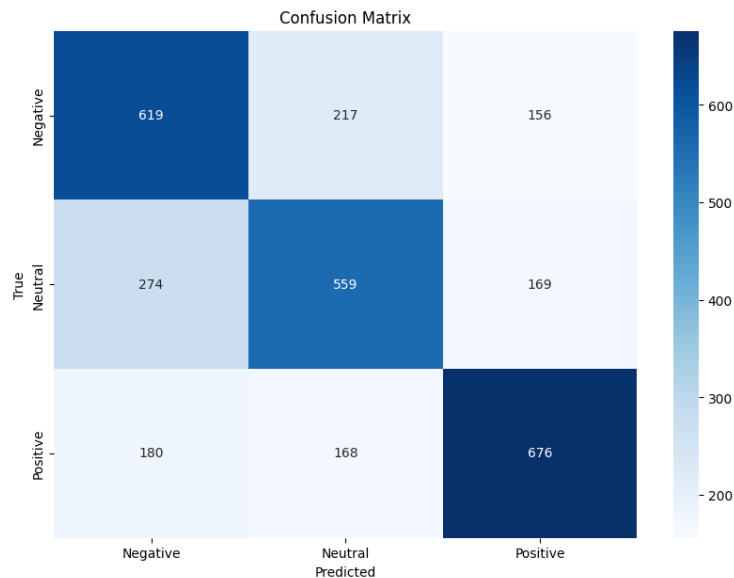
	precision	recall	f1-score	support
negative	0.53	0.59	0.56	1020
neutral	0.52	0.52	0.52	897
positive	0.65	0.58	0.61	1011
accuracy			0.56	3018
macro avg	0.56	0.56	0.56	3018
weighted avg	0.57	0.56	0.56	3018

And the Final Result is as follows:

Classification Report:

	precision	recall	f1-score	support
Negative	0.576887	0.623992	0.599516	992.000000
Neutral	0.592161	0.557884	0.574512	1002.000000
Positive	0.675325	0.660156	0.667654	1024.000000
accuracy	0.614314	0.614314	0.614314	0.614314
macro avg	0.614791	0.614011	0.613894	3018.000000
weighted avg	0.615358	0.614314	0.614333	3018.000000

Figure 5: Confusion Matrix



Overall Metrics:

Accuracy: 0.6143

Precision: 0.6154

Recall: 0.6143

F1 Score: 0.6143

Document Classifier - Transformer Model

The Transformer model is in `DocClassifier_transformer.ipynb`. We investigate each segment in the following sections.

The first part of code is simply loading data and doing simple stuff like dropping missing values.

```

1  # Import required packages
2  import numpy as np
3  import pandas as pd
4  from sklearn.model_selection import train_test_split
5  from sklearn.metrics import classification_report
6  from sklearn.metrics import f1_score
7  from sklearn.utils import shuffle
8  import hazm
9  from cleantext import clean
10 import plotly.express as px
11 import plotly.graph_objects as go
12 from tqdm.notebook import tqdm
13 import os
14 import re
15 import json
16 import copy
17 import collections
18 import seaborn as sns
19 import matplotlib.pyplot as plt
20 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score,
    precision_score, recall_score, f1_score
21 from transformers import BertConfig, BertTokenizer
22 from transformers import TFBertModel, TFBertForSequenceClassification
23 from transformers import glue_convert_examples_to_features
24 import tensorflow as tf
25
26 # Import the pandas library for data manipulation and analysis
27 import pandas as pd

```

```

28
29 # Load the CSV file into a DataFrame
30 data = pd.read_csv('/kaggle/input/taghchemain/taghche.csv')
31
32 # Select only the 'comment' and 'rate' columns for further analysis
33 data = data[['comment', 'rate']]
34
35 # Display the first 10 rows of the DataFrame to get an overview of the data
36 data.head(10)
37
38 # Print the general information about the DataFrame
39 print('data information')
40 print(data.info(), '\n')
41
42 # Print the statistics of missing values in the DataFrame
43 print('missing values stats')
44 print(data.isnull().sum(), '\n')
45
46 # Print the first 5 rows where the 'rate' column has missing values
47 print('some missing values')
48 print(data[data['rate'].isnull()].iloc[:5], '\n')
49
50
51 # Handle conflicts in the dataset structure
52 # For simplicity, remove invalid data combinations
53
54 # Replace ratings that are 6 or higher with None
55 # This assumes that valid ratings should be between 0 and 5
56 data['rate'] = data['rate'].apply(lambda r: r if r < 6 else None)
57
58 # Drop rows where the 'rate' column has missing values
59 data = data.dropna(subset=['rate'])
60
61 # Drop rows where the 'comment' column has missing values
62 data = data.dropna(subset=['comment'])
63
64 # Remove duplicate comments, keeping only the first occurrence
65 data = data.drop_duplicates(subset=['comment'], keep='first')
66
67 # Reset the index of the DataFrame
68 data = data.reset_index(drop=True)
69
70 # Print the general information about the DataFrame
71 print('data information')
72 print(data.info(), '\n')
73
74 # Print the statistics of missing values in the DataFrame
75 print('missing values stats')
76 print(data.isnull().sum(), '\n')
77
78 # Print the first 5 rows where the 'rate' column has missing values
79 print('some missing values')
80 print(data[data['rate'].isnull()].iloc[:5], '\n')

```

Preprocess

Explanations:

Preprocessing

The comments have different lengths based on words. Detecting the most normal range could help us find the maximum length of the sequences for the preprocessing step. On the other hand, we suppose that the minimum word combination for having a meaningful phrase for our learning process is 1.

```

1 # Calculate the length of comments based on the number of words
2 # This uses hazm's word_tokenize function to split comments into words and then counts

```

```

    the number of words
data['comment_len_by_words'] = data['comment'].apply(lambda t: len(hazm.word_tokenize(t)
))

# Determine the minimum and maximum comment lengths
# This provides insights into the range of comment lengths in the dataset
min_max_len = data["comment_len_by_words"].min(), data["comment_len_by_words"].max()

# Print the minimum and maximum comment lengths
# This helps understand the variation in comment lengths
print(f'Min: {min_max_len[0]} \tMax: {min_max_len[1]}')

```

First, the code calculates the length of each comment based on the number of words. It uses the `word_tokenize` function from the `hazm` library to split comments into words. The minimum and maximum lengths of the comments are then determined and printed to understand the variation in comment lengths.

```

def data_gl_than(data, less_than=100.0, greater_than=0.0, col='comment_len_by_words'):
    """
    Calculate the percentage of comments with lengths greater than 'greater_than' and
    less than or equal to 'less_than'.

    Parameters:
    data (DataFrame): The DataFrame containing the data.
    less_than (float): The upper bound for the comment length.
    greater_than (float): The lower bound for the comment length.
    col (str): The column name that contains the comment lengths.

    Returns:
    None
    """
    # Extract the lengths of the comments from the specified column
    data_length = data[col].values

    # Count the number of comments that have a length greater than 'greater_than' and
    # less than or equal to 'less_than'
    data_glt = sum([1 for length in data_length if greater_than < length <= less_than])

    # Calculate the percentage of such comments relative to the total number of comments
    data_glt_rate = (data_glt / len(data_length)) * 100

    # Print the result
    print(f'Texts with word length of greater than {greater_than} and less than {
        less_than} includes {data_glt_rate:.2f}% of the whole!')

# Call the function with specific bounds
data_gl_than(data, 256, 0)

```

The function `data_gl_than` calculates the percentage of comments with lengths greater than a specified lower bound and less than or equal to an upper bound. This is used to analyze the distribution of comment lengths within a specified range.

```

# Define minimum and maximum limits for comment length
minlim, maxlim = 0, 256

# Remove comments with a length of fewer than three words or more than 256 words
data['comment_len_by_words'] = data['comment_len_by_words'].apply(lambda len_t: len_t if
    minlim < len_t <= maxlim else None)

# Drop rows where the 'comment_len_by_words' column has missing values
data = data.dropna(subset=['comment_len_by_words'])

# Reset the index of the DataFrame
data = data.reset_index(drop=True)

```

Comments with a length of fewer than three words or more than 256 words are removed. The DataFrame is then updated to drop rows with missing values in the `comment_len_by_words` column, and the index is reset.


```

1 # This initializes an empty figure to which we will add traces
2 fig = go.Figure()
3
4 # Add a histogram trace to the figure
5 # The x-axis data is taken from the 'comment_len_by_words' column
6 fig.add_trace(go.Histogram(
7     x=data['comment_len_by_words'],
8     marker=dict(
9         color='rgb(0, 123, 255)', # Set bar color to a blue shade
10        line=dict(
11            color='rgb(8, 48, 107)', # Set bar border color to a darker blue
12            width=1.5, # Set the width of the bar borders
13        ),
14    ),
15    opacity=0.75, # Set the opacity of the bars for a slight transparency effect
16))
17
18 # Update the layout of the figure
19 # This includes setting titles for the plot and axes, and customizing the appearance
20 fig.update_layout(
21     title=dict(
22         text='Distribution of Word Counts Within Comments', # Set the title of the plot
23         font=dict(size=24), # Set the font size of the title
24         x=0.5, # Center the title on the plot
25     ),
26     xaxis=dict(
27         title=dict(
28             text='Word Count', # Set the x-axis title
29             font=dict(size=18), # Set the font size of the x-axis title
30         ),
31         tickfont=dict(size=14), # Set the font size of the x-axis tick labels
32     ),
33     yaxis=dict(
34         title=dict(
35             text='Frequency', # Set the y-axis title
36             font=dict(size=18), # Set the font size of the y-axis title
37         ),
38         tickfont=dict(size=14), # Set the font size of the y-axis tick labels
39     ),
40     bargap=0.2, # Set the gap between individual bars
41     bargroupgap=0.2, # Set the gap between groups of bars
42     template='plotly_white' # Use the 'plotly_white' template for a cleaner look
43 )
44
45 # Show the figure
46 fig.show()

```

A histogram is created to visualize the distribution of word counts within comments. The layout of the figure is customized, including setting titles for the plot and axes, and using the `plotly_white` template for a cleaner look. The figure is then displayed.

```

1 # Extract unique 'rate' values from the 'data' DataFrame, sort them, and convert them
   into a list
2 unique_rates = list(sorted(data['rate'].unique()))
3
4 # Print the number of unique rates and the list of unique rates
5 print(f'We have {len(unique_rates)} unique rates: {unique_rates}')
6
7
8 # Create a figure for the plot
9 fig = go.Figure()
10
11 # Group the data by 'rate' and count the occurrences of each rate
12 groupby_rate = data.groupby('rate')['rate'].count()
13
14 # Add a bar trace to the figure

```

```

5 fig.add_trace(go.Bar(
6     x=list(sorted(groupby_rate.index)), # Set the x-axis data to the sorted unique
       rates
7     y=groupby_rate.tolist(), # Set the y-axis data to the frequency counts of each rate
8     text=groupby_rate.tolist(), # Display the frequency counts as text on the bars
9     textposition='auto', # Automatically position the text within the bars
10    marker=dict(
11        color='rgb(255, 123, 0)', # Set bar color to an orange shade
12        line=dict(
13            color='rgb(107, 48, 8)', # Set bar border color to a darker orange/brown
14            width=1.5, # Set the width of the bar borders
15        ),
16    ),
17    opacity=0.75 # Set the opacity of the bars for a slight transparency effect
18))
19
20 # Update the layout of the figure
21 # This includes titles for the plot and axes, as well as setting bar gaps
22 fig.update_layout(
23     title=dict(
24         text='Distribution of Rates Within Comments', # Set the title of the plot
25         font=dict(size=24), # Set title font size
26         x=0.5, # Center the title
27     ),
28     xaxis=dict(
29         title=dict(
30             text='Rate', # Set the x-axis title
31             font=dict(size=18), # Set x-axis title font size
32         ),
33         tickfont=dict(size=14), # Set x-axis tick font size
34     ),
35     yaxis=dict(
36         title=dict(
37             text='Frequency', # Set the y-axis title
38             font=dict(size=18), # Set y-axis title font size
39         ),
40         tickfont=dict(size=14), # Set y-axis tick font size
41     ),
42     bargap=0.2, # Set the gap between individual bars
43     bargroupgap=0.2, # Set the gap between groups of bars
44     template='plotly_white' # Set the plot template for a cleaner look
45 )
46
47 # Show the figure
48 fig.show()

```

A bar chart is created to visualize the distribution of rates within comments. The layout of the figure is customized, including setting titles for the plot and axes, and using the `plotly_white` template for a cleaner look. The figure is then displayed.

```

1 # Function to label the sentiment based on rating thresholds
2 def label_sentiment(rate, positive_threshold, neutral_threshold):
3     """
4     Labels sentiment based on rating thresholds.
5
6     Args:
7     - rate (int or float): The numerical rating to evaluate.
8     - positive_threshold (int or float): The minimum rating value that qualifies as '
9       positive'.
10    - neutral_threshold (int or float): The minimum rating value that qualifies as '
11      neutral'; ratings below this are considered 'negative'.
12
13    Returns:
14    - str: The sentiment label ('positive', 'neutral', or 'negative') based on the
15      rating.
16    """
17    # Check if the rating is greater than or equal to the positive threshold

```

```

5  if rate >= positive_threshold:
6      return 'positive'
7  # If the rating is not 'positive', check if it is greater than or equal to the
8  neutral threshold
9  elif rate >= neutral_threshold:
10     return 'neutral'
11  # If the rating is neither 'positive' nor 'neutral', label the sentiment as '
12  negative'
13  else:
14     return 'negative'
15
16  # Apply the label_sentiment function to the 'rate' column to create a new 'label' column
17  data['label'] = data['rate'].apply(lambda t: label_sentiment(t, 4.0, 2.0))
18  # Print unique labels
19  labels = list(sorted(data['label'].unique()))
20  # Display the first 5 rows of the DataFrame to get an overview of the data
21  data.head()

```

The function `label_sentiment` is defined to label the sentiment based on rating thresholds. Ratings are classified as *positive*, *neutral*, or *negative* based on specified thresholds. The function is applied to the `rate` column to create a new label column.

```

1  def preprocess(text):
2      """
3      Preprocess and normalize the given text string by performing several transformations
4      .
5
6      Args:
7      text (str): The text string to be preprocessed.
8
9      Returns:
10     str: The preprocessed and normalized text string.
11     """
12     # Define a regular expression pattern that matches one or more newline characters
13     pattern = re.compile(r"\n+")
14
15     # Replace multiple newlines with a single newline
16     text = pattern.sub("\n", text)
17
18     # Remove special characters and replace newline characters with spaces
19     text = re.sub(r'\\n|\\n', ' ', text)
20
21     # Remove non-Persian characters and digits
22     text = re.sub(r'[^\\s\\u200c]', ' ', text)
23
24     # Define a regular expression pattern that matches one or more spaces
25     pattern = re.compile(r" +")
26
27     # Replace multiple spaces with a single space
28     text = pattern.sub(" ", text)
29
30     return text
31
32  # Apply the preprocess function to the 'comment' column of the DataFrame
33  data['cleaned_comment'] = data['comment'].apply(preprocess)
34
35  # Remove duplicate rows from the DataFrame
36  data = data.drop_duplicates()
37
38  # Drop rows with empty or NaN values in the 'comment' or 'rate' columns
39  data.dropna(subset=['comment', 'rate'], inplace=True)
40
41  # Display the first 10 rows of the cleaned DataFrame
42  data.head(10)

```

The function `preprocess` is defined to preprocess and normalize the text. It performs several transforma-

tions, including replacing multiple newlines with a single newline, removing special characters, and removing non-Persian characters and digits. The function is applied to the `comment` column. Duplicate rows are removed from the DataFrame, and rows with empty or NaN values in the `comment` or `rate` columns are dropped.

```
1 # Create a figure for the plot
2 fig = go.Figure()
3
4 # Group the data by 'label' and count the occurrences of each label
5 groupby_label = data.groupby('label')['label'].count()
6
7 # Add a bar trace to the figure
8 fig.add_trace(go.Bar(
9     x=list(sorted(groupby_label.index)), # Set the x-axis data to the sorted unique
10         labels
11     y=groupby_label.tolist(), # Set the y-axis data to the frequency counts of each
12         label
13     text=groupby_label.tolist(), # Display the frequency counts as text on the bars
14     textposition='auto', # Automatically position the text within the bars
15     marker=dict(
16         color='rgb(0, 123, 255)', # Set bar color to a blue shade
17         line=dict(
18             color='rgb(8, 48, 107)', # Set bar border color to a darker blue
19             width=1.5, # Set the width of the bar borders
20         ),
21     ),
22     opacity=0.75 # Set the opacity of the bars for a slight transparency effect
23 ))
24
25 # Update the layout of the figure
26 fig.update_layout(
27     title=dict(
28         text='Distribution of Labels Within Comments', # Set the title of the plot
29         font=dict(size=24), # Set title font size
30         x=0.5, # Center the title
31     ),
32     xaxis=dict(
33         title=dict(
34             text='Label', # Set the x-axis title
35             font=dict(size=18), # Set x-axis title font size
36         ),
37         tickfont=dict(size=14), # Set x-axis tick font size
38     ),
39     yaxis=dict(
40         title=dict(
41             text='Frequency', # Set the y-axis title
42             font=dict(size=18), # Set y-axis title font size
43         ),
44         tickfont=dict(size=14), # Set y-axis tick font size
45     ),
46     bargap=0.2, # Set the gap between individual bars
47     bargroupgap=0.2, # Set the gap between groups of bars
48     template='plotly_white' # Set the plot template for a cleaner look
49 )
50
51 # Show the figure
52 fig.show()
```

The code creates a bar chart to visualize the distribution of sentiment labels within the comments. The layout is customized, and the plot is displayed.

Balancing the Labels via Resampling

```
1 # Filter data into separate DataFrames based on 'label' values
2 negative_data = data[data['label'] == 'negative']
3 positive_data = data[data['label'] == 'positive']
4 neutral_data = data[data['label'] == 'neutral']
5
```

```

6 # Determine the smallest length among the three filtered DataFrames
7 cutting_point = min(len(negative_data), len(positive_data), len(neutral_data))
8
9 # If the cutting_point is less than or equal to the length of negative_data, sample
10 # down negative_data
11 if cutting_point <= len(negative_data):
12     negative_data = negative_data.sample(n=cutting_point).reset_index(drop=True)
13
14 # If the cutting_point is less than or equal to the length of positive_data, sample
15 # down positive_data
16 if cutting_point <= len(positive_data):
17     positive_data = positive_data.sample(n=cutting_point).reset_index(drop=True)
18
19 # If the cutting_point is less than or equal to the length of neutral_data, sample
20 # down neutral_data
21 if cutting_point <= len(neutral_data):
22     neutral_data = neutral_data.sample(n=cutting_point).reset_index(drop=True)
23
24 # Concatenate the sampled DataFrames back into a single DataFrame
25 new_data = pd.concat([negative_data, positive_data, neutral_data])
26
27 # Shuffle the rows of the new DataFrame
28 new_data = new_data.sample(frac=1).reset_index(drop=True)
29
30 # Display summary information about the new DataFrame
31 new_data.info()

```

The code balances the dataset by resampling each sentiment label to have an equal number of samples. This is done by determining the smallest length among the sentiment labels and sampling down the other labels to match this length. The DataFrames are concatenated, and the rows are shuffled.

Distribution of Labels Within Comments After Balancing

```

1 # Create a figure for the plot
2 fig = go.Figure()
3
4 # Group the data by 'label' and count the occurrences of each label
5 groupby_label = new_data.groupby('label')['label'].count()
6
7 # Add a bar trace to the figure
8 fig.add_trace(go.Bar(
9     x=list(sorted(groupby_label.index)), # Set the x-axis data to the sorted unique
10     y=groupby_label.tolist(), # Set the y-axis data to the frequency counts of each
11     text=groupby_label.tolist(), # Display the frequency counts as text on the bars
12     textposition='auto', # Automatically position the text within the bars
13     marker=dict(
14         color='rgb(0, 123, 255)', # Set bar color to a blue shade
15         line=dict(
16             color='rgb(8, 48, 107)', # Set bar border color to a darker blue
17             width=1.5, # Set the width of the bar borders
18         ),
19     ),
20     opacity=0.75 # Set the opacity of the bars for a slight transparency effect
21 ))
22
23 # Update the layout of the figure
24 fig.update_layout(
25     title=dict(
26         text='Distribution of Labels Within Comments After Balancing Labels via
27         Resampling', # Set the title of the plot
28         font=dict(size=24), # Set title font size
29         x=0.5, # Center the title
30     ),
31     xaxis=dict(
32         title=dict(

```

```

32         text='Label', # Set the x-axis title
33         font=dict(size=18), # Set x-axis title font size
34     ),
35     tickfont=dict(size=14), # Set x-axis tick font size
36 ),
37 yaxis=dict(
38     title=dict(
39         text='Frequency', # Set the y-axis title
40         font=dict(size=18), # Set y-axis title font size
41     ),
42     tickfont=dict(size=14), # Set y-axis tick font size
43 ),
44 bargap=0.2, # Set the gap between individual bars
45 bargroupgap=0.2, # Set the gap between groups of bars
46 template='plotly_white' # Set the plot template for a cleaner look
47 )
48
49 # Show the figure
50 fig.show()

```

A second bar chart is created to visualize the distribution of sentiment labels within the comments after balancing the dataset via resampling. The layout is customized, and the plot is displayed.

Training Configuration

Explanations:

Mapping Labels to Numerical IDs and Data Splitting

```

1 # Map labels to numerical ids and add a new column 'label_id' to new_data
2 new_data['label_id'] = new_data['label'].apply(lambda t: labels.index(t))
3
4 # Split new_data into train and test sets (80% train, 20% test), stratified by 'label'
5 train, test = train_test_split(new_data, test_size=0.2, random_state=1, stratify=
6     new_data['label'])
7
8 # Further split test set into test and validation sets (50% test, 50% validation),
9     stratified by 'label'
10 test, valid = train_test_split(test, test_size=0.5, random_state=1, stratify=test['label
11     '])
12
13 # Reset index for train, validation, and test sets to ensure continuous integer indices
14 train = train.reset_index(drop=True)
15 valid = valid.reset_index(drop=True)
16 test = test.reset_index(drop=True)
17
18 # Extract comments and label_ids for train, validation, and test sets
19 x_train, y_train = train['comment'].values.tolist(), train['label_id'].values.tolist()
20 x_valid, y_valid = valid['comment'].values.tolist(), valid['label_id'].values.tolist()
21 x_test, y_test = test['comment'].values.tolist(), test['label_id'].values.tolist()
22
23 # Print shapes of train, validation, and test sets to verify sizes
24 print(train.shape)
25 print(valid.shape)
26 print(test.shape)

```

The code maps the text labels to numerical IDs and adds a new column 'label_id' to the DataFrame. The data is then split into train (80%), test (20%), and validation (50% of the test set) sets, ensuring stratification by label to maintain the distribution of classes in each subset. The indices of each set are reset, and the comments and label IDs are extracted for each set. Finally, the shapes of the train, validation, and test sets are printed to verify the sizes.

Configuration Parameters for BERT Model Training

```

1 # General configuration parameters
2 MAX_LEN = 128 # Maximum sequence length
3 TRAIN_BATCH_SIZE = 64 # Batch size for training
4 VALID_BATCH_SIZE = 64 # Batch size for validation
5 TEST_BATCH_SIZE = 64 # Batch size for testing
6
7 EPOCHS = 3 # Number of epochs for training
8 EVERY_EPOCH = 1000 # Number of steps to print progress during each epoch
9 LEARNING_RATE = 2e-5 # Learning rate for the optimizer
10 CLIP = 0.0 # Gradient clipping threshold
11
12 MODEL_NAME_OR_PATH = 'HooshvareLab/bert-fa-base-uncased' # Pre-trained model name or
   path
13 OUTPUT_PATH = '/content/bert-fa-base-uncased-sentiment-taaghceh/pytorch_model.bin' #
   Path to save the trained model
14
15 # Create the directory to save the trained model if it does not exist
16 os.makedirs(os.path.dirname(OUTPUT_PATH), exist_ok=True)

```

The general configuration parameters for training the BERT model are defined. This includes the maximum sequence length, batch sizes for training, validation, and testing, the number of epochs, learning rate, gradient clipping threshold, the pre-trained model path, and the path to save the trained model. The output directory is created if it does not already exist.

Label Mapping Dictionaries and BERT Configuration

```

1 # Create a dictionary mapping labels to numerical ids
2 label2id = {label: i for i, label in enumerate(labels)}
3
4 # Create a dictionary mapping numerical ids back to labels
5 id2label = {v: k for k, v in label2id.items()}
6
7 # Print label2id dictionary
8 print(f'label2id: {label2id}')
9
10 # Print id2label dictionary
11 print(f'id2label: {id2label}')
12
13 # Initialize a BERT tokenizer using the pre-trained model specified in
   MODEL_NAME_OR_PATH
14 tokenizer = BertTokenizer.from_pretrained(MODEL_NAME_OR_PATH, force_download=True)
15
16 # Create a BERT configuration object using the pre-trained model and additional custom
   settings
17 config = BertConfig.from_pretrained(
18     MODEL_NAME_OR_PATH, # Specify the pre-trained model name or path
19     **{ # Additional custom settings passed as keyword arguments
20         'label2id': label2id, # Mapping from labels to numerical ids
21         'id2label': id2label # Mapping from numerical ids back to labels
22     })
23
24 # Print the configuration details in JSON format
25 print(config.to_json_string())

```

Label mapping dictionaries are created to map labels to numerical IDs and vice versa. These mappings are printed for verification. A BERT tokenizer and configuration object are initialized using the specified pre-trained model path and the custom label mappings. The configuration details are printed in JSON format.

Explanations:

Preparing Datasets for BERT Model Training

```

1 import tensorflow as tf
2 import numpy as np
3 from tqdm import tqdm
4 from transformers import glue_convert_examples_to_features
5
6 class InputExample:
7     """ A single example for simple sequence classification. """
8
9     def __init__(self, guid, text_a, text_b=None, label=None):
10         """ Constructs an InputExample. """
11         self.guid = guid
12         self.text_a = text_a
13         self.text_b = text_b
14         self.label = label
15
16 def make_examples(tokenizer, x, y=None, maxlen=128, output_mode="classification",
17 is_tf_dataset=True):
18     """
19     Converts input texts and labels into examples and features suitable for BERT model
20     training.
21
22     Args:
23     tokenizer (transformers.PreTrainedTokenizer): Tokenizer object for tokenizing input
24     texts.
25     x (list): List of input texts or tuples of input texts (for sequence classification)
26     .
27     y (list, optional): List of labels corresponding to input texts. Default is None.
28     maxlen (int, optional): Maximum sequence length for input texts. Default is 128.
29     output_mode (str, optional): Output mode for the task (e.g., "classification").
30     Default is "classification".
31     is_tf_dataset (bool, optional): Whether to return a TensorFlow dataset or numpy
32     arrays. Default is True.
33
34     Returns:
35     tf.data.Dataset or tuple of numpy arrays: Depending on is_tf_dataset, returns either
36     a TensorFlow dataset
37
38     or a tuple of numpy arrays containing
39     input_ids, attention_masks,
40     token_type_ids, and labels.
41
42     list: List of features converted from InputExamples.
43     """
44     examples = []
45     y = y if isinstance(y, list) or isinstance(y, np.ndarray) else [None] * len(x)
46
47     # Create InputExamples from input texts and labels
48     for i, (_x, _y) in tqdm(enumerate(zip(x, y)), position=0, total=len(x)):
49         guid = "%s" % i
50         label = int(_y)
51
52         if isinstance(_x, str):
53             text_a = _x
54             text_b = None
55         else:
56             assert len(_x) == 2
57             text_a = _x[0]
58             text_b = _x[1]
59
60         examples.append(InputExample(guid=guid, text_a=text_a, text_b=text_b, label=
61 label))

```



```

# Convert InputExamples to features using glue_convert_examples_to_features function
features = glue_convert_examples_to_features(
    examples,
    tokenizer,
    maxlen,
    output_mode=output_mode,
    label_list=list(np.unique(y)))

all_input_ids = []
all_attention_masks = []
all_token_type_ids = []
all_labels = []

# Process features into input arrays for TensorFlow dataset or numpy arrays
for f in tqdm(features, position=0, total=len(examples)):
    if is_tf_dataset:
        all_input_ids.append(tf.constant(f.input_ids))
        all_attention_masks.append(tf.constant(f.attention_mask))
        all_token_type_ids.append(tf.constant(f.token_type_ids))
        all_labels.append(tf.constant(f.label))
    else:
        all_input_ids.append(f.input_ids)
        all_attention_masks.append(f.attention_mask)
        all_token_type_ids.append(f.token_type_ids)
        all_labels.append(f.label)

if is_tf_dataset:
    # Create TensorFlow dataset from input arrays
    dataset = tf.data.Dataset.from_tensor_slices(({
        'input_ids': all_input_ids,
        'attention_mask': all_attention_masks,
        'token_type_ids': all_token_type_ids
    }, all_labels))

    return dataset, features

# Return tuple of numpy arrays if is_tf_dataset=False
xdata = [np.array(all_input_ids), np.array(all_attention_masks), np.array(
    all_token_type_ids)]
ydata = all_labels

return [xdata, ydata], features

# Create training dataset and examples using make_examples function
train_dataset_base, train_examples = make_examples(tokenizer, x_train, y_train, maxlen
=128)

# Create validation dataset and examples using make_examples function
valid_dataset_base, valid_examples = make_examples(tokenizer, x_valid, y_valid, maxlen
=128)

# Create test dataset and examples using make_examples function
test_dataset_base, test_examples = make_examples(tokenizer, x_test, y_test, maxlen=128)

# Create test dataset and examples as numpy arrays (not TensorFlow dataset)
[xtest, ytest], test_examples = make_examples(tokenizer, x_test, y_test, maxlen=128,
is_tf_dataset=False)

# Iterate over the first example in the train_dataset_base
for value in train_dataset_base.take(1):
    # Print input_ids tensor
    print(f'input_ids: {value[0]["input_ids"]}')
    # Print attention_mask tensor
    print(f'attention_mask: {value[0]["attention_mask"]}')
    # Print token_type_ids tensor
    print(f'token_type_ids: {value[0]["token_type_ids"]}')

```

```

4 # Print target (label) tensor
5 print(f'          target: {value[1]}')

```

The code defines an ‘InputExample’ class for creating input examples for sequence classification and a ‘make_examples’ function to convert texts and labels into examples and features suitable for BERT model training. It creates TensorFlow datasets for training, validation, and testing, and prints the first example in the training dataset to verify correctness.

Building and Training the BERT Model

```

1 def get_training_dataset(dataset, batch_size):
2     """
3     Creates a training dataset pipeline.
4
5     Args:
6     dataset (tf.data.Dataset): TensorFlow dataset containing training examples.
7     batch_size (int): Batch size for training.
8
9     Returns:
10    tf.data.Dataset: Processed training dataset ready for model training.
11    """
12    # Repeat the dataset indefinitely
13    dataset = dataset.repeat()
14    # Shuffle the dataset with a buffer size of 2048
15    dataset = dataset.shuffle(2048)
16    # Batch the dataset with the specified batch size
17    dataset = dataset.batch(batch_size)
18
19    return dataset
20
21 def get_validation_dataset(dataset, batch_size):
22     """
23     Creates a validation dataset pipeline.
24
25     Args:
26     dataset (tf.data.Dataset): TensorFlow dataset containing validation examples.
27     batch_size (int): Batch size for validation.
28
29     Returns:
30    tf.data.Dataset: Processed validation dataset ready for model evaluation.
31    """
32    # Batch the dataset with the specified batch size
33    dataset = dataset.batch(batch_size)
34
35    return dataset
36
37 # Create training dataset using get_training_dataset function
38 train_dataset = get_training_dataset(train_dataset_base, TRAIN_BATCH_SIZE)
39
40 # Create validation dataset using get_validation_dataset function
41 valid_dataset = get_validation_dataset(valid_dataset_base, VALID_BATCH_SIZE)
42
43 # Calculate steps per epoch for training and validation
44 train_steps = len(train_examples) // TRAIN_BATCH_SIZE
45 valid_steps = len(valid_examples) // VALID_BATCH_SIZE
46
47 # Print the calculated steps for training and validation
48 print(train_steps, valid_steps)
49
50 def build_model(model_name, config, learning_rate=3e-5):
51     """
52     Builds a TensorFlow model for sequence classification using a pre-trained BERT model
53     .
54
55     Args:
56     model_name (str): Pre-trained model name or path.

```

```

36     config (transformers.PretrainedConfig): Configuration object for the BERT model.
37     learning_rate (float, optional): Learning rate for optimizer. Default is 3e-5.
38
39     Returns:
40     tf.keras.Model: Compiled BERT-based model for sequence classification.
41     """
42     # Load pre-trained BERT model for sequence classification
43     model = TFBertForSequenceClassification.from_pretrained(model_name, force_download=
44         True, config=config)
45
46     # Define optimizer, loss function, and metrics for the model
47     optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate)
48     loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
49     metric = tf.keras.metrics.SparseCategoricalAccuracy('accuracy')
50
51     # Compile the model with optimizer, loss function, and metrics
52     model.compile(optimizer=optimizer, loss=loss, metrics=[metric])
53
54     return model
55
56 # Build BERT-based model for sequence classification
57 model = build_model(MODEL_NAME_OR_PATH, config, learning_rate=LEARNING_RATE)
58
59 %%time
60
61 # Train the BERT-based model
62 r = model.fit(
63     train_dataset,                    # Training dataset
64     validation_data=valid_dataset,    # Validation dataset
65     steps_per_epoch=train_steps,       # Number of steps per epoch for training
66     batch_size=128,                   # Batch size for training
67     validation_steps=valid_steps,      # Number of steps per epoch for validation
68     epochs=EPOCHS,                   # Number of epochs
69     verbose=1                         # Verbosity mode (1 for progress bar)
70 )
71
72 # Extract validation accuracy from training history
73 final_accuracy = r.history['val_accuracy']
74 print('FINAL ACCURACY MEAN: ', np.mean(final_accuracy))
75
76 # Save the trained model
77 model.save_pretrained(os.path.dirname(OUTPUT_PATH))

```

The 'get_training_dataset' and 'get_validation_dataset' functions create data pipelines for training and validation. The 'build_model' function builds and compiles a BERT-based model for sequence classification. The model is then trained, and its validation accuracy is printed and the trained model is saved.

0.0.2 Results

Results are generated using the following code

Explanations:

```

1  def evaluate_model(y_true, y_pred, class_names):
2      """
3      Evaluates the performance of a classification model using various metrics and
4      visualizations.
5
6      Args:
7      - y_true (array-like): True labels of the data.
8      - y_pred (array-like): Predicted labels of the data.
9      - class_names (list): List of class names in the same order as the confusion
10         matrix.
11
12     Returns:

```

```

1 - pd.DataFrame: DataFrame containing the classification report.
2 """
3 # Generate and print the classification report
4 report = classification_report(y_true, y_pred, target_names=class_names,
5                               output_dict=True)
6 report_df = pd.DataFrame(report).transpose()
7 print("Classification Report:\n", report_df)
8
9 # Generate and display the confusion matrix as a heatmap
10 cm = confusion_matrix(y_true, y_pred)
11 plt.figure(figsize=(10, 7))
12 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=class_names,
13            yticklabels=class_names)
14 plt.xlabel('Predicted')
15 plt.ylabel('True')
16 plt.title('Confusion Matrix')
17 plt.show()
18
19 # Calculate and print overall metrics: accuracy, precision, recall, and F1 score
20 accuracy = accuracy_score(y_true, y_pred)
21 precision = precision_score(y_true, y_pred, average='weighted')
22 recall = recall_score(y_true, y_pred, average='weighted')
23 f1 = f1_score(y_true, y_pred, average='weighted')
24
25 metrics = {
26     "Accuracy": accuracy,
27     "Precision": precision,
28     "Recall": recall,
29     "F1 Score": f1
30 }
31
32 print("\nOverall Metrics:")
33 for metric, value in metrics.items():
34     print(f"{metric}: {value:.4f}")
35
36 return report_df
37
38 # Perform predictions on the test dataset
39 predictions = model.predict(xtest)
40
41 # Extract predicted labels from predictions
42 ypred = predictions[0].argmax(axis=-1).tolist()
43
44 # Evaluate the model
45 report_df = evaluate_model(ytest, ypred, ["Negative", "Neutral", "Positive"])

```

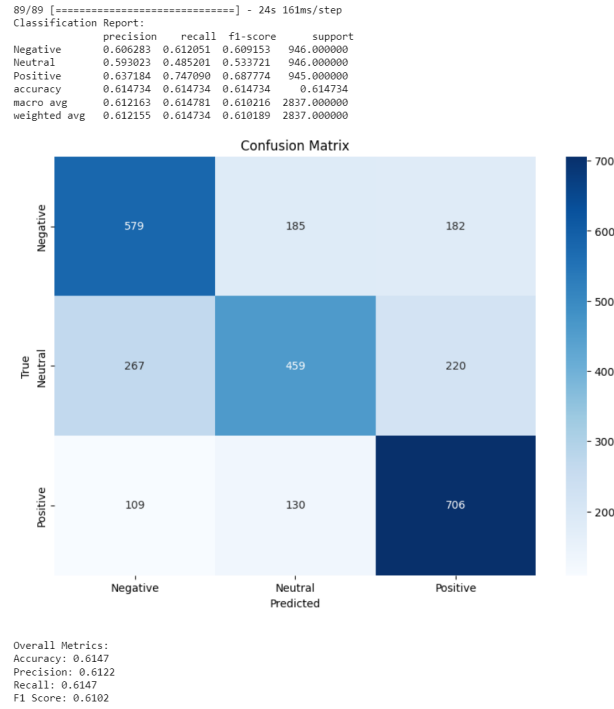
The 'evaluate_model' function assesses the performance of a classification model using various metrics and visualizations. It generates a classification report, displays a confusion matrix as a heatmap, and calculates overall metrics such as accuracy, precision, recall, and F1 score. The function prints these metrics and returns the classification report as a DataFrame.

The code performs predictions on the test dataset using the trained BERT model, extracts the predicted labels, and evaluates the model's performance by calling the 'evaluate_model' function. The function is provided with the true labels, predicted labels, and class names for evaluation.

The following steps are performed:

1. The 'evaluate_model' function generates a classification report and displays it.
2. A confusion matrix is generated and visualized as a heatmap.
3. Overall metrics (accuracy, precision, recall, and F1 score) are calculated and printed.
4. The classification report is returned as a DataFrame.

Figure 6: Results of running Transformer on Dataset



Word Classifier - Base Model

The base model is in `WordClassifier_Base.ipynb`. We investigate each segment in the following sections.

Initial Preprocessing

Explanations:

```
1 prepared_data: pd.DataFrame = pd.read_csv(url + 'datasets/taghche.csv')
2 prepared_data = prepared_data[['comment', 'bookname', 'bookID']]
```

The script reads a CSV file containing text data and selects specific columns: `comment`, `bookname`, and `bookID`.

```
1 chars_stop_words = ''
2 with open(url + 'stopwords/chars (without digits).txt', 'r', encoding='utf-8') as file:
3     :
4     chars_stop_words = ''.join(file.read().splitlines())
5
6 chars_stop_words = chars_stop_words.replace('[', '\\[')
7 chars_stop_words = chars_stop_words.replace(']', '\\]')
8 chars_pattern = re.compile(f'[{chars_stop_words}]')
9 chars_pattern
```

The script reads a file containing stop words (characters to be removed) and creates a regular expression pattern to match these characters.

```
1 emojis_pattern = re.compile("[
2     u"\U0001F600-\U0001F64F" # emoticons
3     u"\U0001F300-\U0001F5FF" # symbols & pictographs
4     u"\U0001F680-\U0001F6FF" # transport & map symbols
5     u"\U0001F1E0-\U0001F1FF" # flags (iOS)
6     ]+")
```

The script defines a regular expression pattern to match emojis using Unicode ranges.

```
1 def elementary_preprocess(text):
2     global chars_pattern, emojis_pattern
3
```

```

4 text = str(text)
5 text = chars_pattern.sub(r' ', text)
6 text = emojis_pattern.sub(r' ', text)
7 return text.translate(str.maketrans('0123456789', ''))

```

The `elementary_preprocess` function removes unwanted characters and emojis from the text and translates digits to Persian numerals.

```

1 def higher_preprocess(text, is_informal=False):
2     global normalizer
3
4     text = str(text)
5
6     if is_informal:
7         text = informal_normalizer_function(text)
8         # progress_bar.update(1)
9     else:
10        text = normalizer.normalize(text)
11
12    text = word_tokenize(text)
13    return text

```

The `higher_preprocess` function normalizes the text. If `is_informal` is `True`, it uses an informal normalizer; otherwise, it uses a standard normalizer. The function tokenizes the normalized text.

```

1 def informal_normalizer_function(text):
2     global informal_normalizer
3     text = str(text)
4
5     informal_normalizer = InformalNormalizer()
6     text = Normalizer.normalize(informal_normalizer, text)
7     sents = [
8         informal_normalizer.word_tokenizer.tokenize(sentence)
9         for sentence in informal_normalizer.sent_tokenizer.tokenize(text)
10    ]
11
12    normalized = [[informal_normalizer.normalized_word(word)[0] for word in sent] for
13                  sent in sents]
14    normalized = np.array(normalized, dtype=object)
15    return np.hstack(normalized)
16
17    normalizer = Normalizer()
18    informal_normalizer = InformalNormalizer()

```

The `informal_normalizer_function` customizes the normalization process for informal text. It tokenizes sentences and words, then normalizes each word.

```

1 for column in prepared_data.columns:
2     if column == 'bookID':
3         continue
4
5     print(f'Column: {column}')
6     prepared_data[column] = prepared_data[column].progress_apply(elementary_preprocess)

```

The script applies the `elementary_preprocess` function to each column of the data except `bookID`.

```

1 for column in prepared_data.columns:
2     if column == 'bookID':
3         continue
4
5     print(f'Column: {column}')
6     prepared_data[column] = prepared_data[column].progress_apply(higher_preprocess)

```

Next, the script applies the `higher_preprocess` function to each column of the data except `bookID`.

```

1 before_dropping = len(prepared_data)
2 prepared_data = prepared_data[prepared_data['comment'].apply(lambda x: len(x) != 0)]

```

```

3     print(f'Dropped {before_dropping - len(prepared_data)} rows with empty comment.')
4
5     before_dropping = len(prepared_data)
6     prepared_data = prepared_data.dropna(subset=['bookID'])
7     print(f'Dropped {before_dropping - len(prepared_data)} rows with NaN bookID.')

```

Finally, the script removes rows with empty comments and rows with missing bookID values, printing the number of dropped rows.

Using Crawled Data

Explanations:

The following Python script processes and normalizes a dataset containing book data. It combines multiple CSV files into a single DataFrame, sorts and processes author names, removes duplicates and rows with missing IDs, and normalizes text data.

```

1  ALL_PARTS_LEN = 19
2  crawled_data: pd.DataFrame = pd.read_csv(url + 'datasets/books data/books_data_part_1.csv'
3  )
4  for i in range(2, ALL_PARTS_LEN + 1):
5      crawled_data = pd.concat([crawled_data, pd.read_csv(url + f'datasets/books data/
        books_data_part_{i}.csv')],
        ignore_index=True)

```

The script first reads the initial CSV file into a DataFrame called `crawled_data`. It then iteratively reads and concatenates additional CSV files into this DataFrame. The variable `ALL_PARTS_LEN` determines the number of parts to be combined.

```

1  new_author_function = lambda x: ' $ '.join(sorted(str(x).split(' $ ')))
2  crawled_data['author'] = crawled_data['author'].apply(new_author_function)

```

This section of the code defines a lambda function to sort the author names within each entry. The authors are separated by the delimiter ' \$ '. The lambda function sorts the names alphabetically and joins them back together with the same delimiter.

```

1  before_dropping = len(crawled_data)
2  crawled_data = crawled_data.drop_duplicates()
3  print(f'Dropped {before_dropping - len(crawled_data)} duplicates.')

```

The script removes duplicate rows from `crawled_data` and prints the number of duplicates dropped.

```

1  before_dropping = len(crawled_data)
2  crawled_data = crawled_data.dropna(subset=['id'])
3  print(f'Dropped {before_dropping - len(crawled_data)} rows with NaN id.')

```

Rows with missing id values are dropped, and the number of such rows is printed.

```

1  new_author_function = lambda x: set(x.split(' $ '))
2  crawled_data['author'] = crawled_data['author'].apply(new_author_function)

```

The script redefines the lambda function to convert the list of authors into a set, effectively eliminating duplicate authors within each entry.

```

1  crawled_data = crawled_data.explode('author')
2  crawled_data = crawled_data.reset_index(drop=True)

```

The `explode` method is used to transform each author into a separate row, enabling independent processing of each author in the comments. The index is reset to maintain a clean DataFrame.

```

1  for column in crawled_data.columns:
2      if column == 'id':
3          continue
4
5      print(f'Column: {column}')
6      crawled_data[column] = crawled_data[column].progress_apply(elementary_preprocess)

```

The script applies the `elementary_preprocess` function to each column of the `crawled_data` DataFrame, excluding the `id` column. This function removes unwanted characters and emojis and translates digits to Persian numerals.

```
1 for column in crawled_data.columns:
2     if column == 'id':
3         continue
4
5     print(f'Column: {column}')
6     crawled_data[column] = crawled_data[column].progress_apply(higher_preprocess)
```

The script then applies the `higher_preprocess` function to each column of the `crawled_data` DataFrame, again excluding the `id` column. This function normalizes and tokenizes the text.

The following Python script merges two datasets, identifies unavailable books, and saves a list of unavailable book IDs to a file. The datasets are preprocessed to remove duplicates and missing values before being merged.

```
1 data: pd.DataFrame = pd.merge(prepared_data, crawled_data, left_on='bookID', right_on='id')
2 data = data.drop(columns=['bookID'])
3
4 print(f'Prepared data: {len(prepared_data)}\nCrawled data: {len(crawled_data)}\nMerged data: {
5     len(data)}')
6 data.head()
```

The script merges the `prepared_data` and `crawled_data` DataFrames using a common key: `bookID` from `prepared_data` and `id` from `crawled_data`. After merging, the `bookID` column is dropped from the merged DataFrame. The script then prints the lengths of the prepared, crawled, and merged datasets and displays the first few rows of the merged dataset.

```
1 crawled_books = set(crawled_data['id'].values)
2 prepared_books = prepared_data[['bookID']].copy()
3
4 unavailable_books = prepared_books[~prepared_books['bookID'].apply(lambda x: x in
5     crawled_books)]
6 unavailable_books = unavailable_books.drop_duplicates()
7 print(f'Unavailable books (The page has 404 error): {len(unavailable_books)}')
8 unavailable_books
```

The script identifies books in the `prepared_data` DataFrame that are not present in the `crawled_data` DataFrame. It creates a set of book IDs from `crawled_data` and compares it to the `bookID` column in `prepared_data`. Books not found in `crawled_data` are considered unavailable. The script removes duplicate entries and prints the number of unavailable books.

```
1 with open('unavailable_books_list.txt', 'w', encoding='utf-8') as file:
2     unavailable_books_list = unavailable_books['bookID'].values.flatten()
3     unavailable_books_list = unavailable_books_list[~np.isnan(unavailable_books_list)]
4     unavailable_books_list = unavailable_books_list.astype(int)
5     unavailable_books_list = unavailable_books_list.tolist()
6     unavailable_books_list = sorted(list(set(unavailable_books_list)))
7     file.write(str(unavailable_books_list))
```

The script writes the list of unavailable book IDs to a text file. It flattens the array of book IDs, removes any NaN values, converts the IDs to integers, and sorts them. Finally, it saves the sorted list of unique unavailable book IDs to a file named `unavailable_books_list.txt`.

Explanations:

The following Python script processes a dataset to label parts of text related to books, authors, translators, and publishers. It then visualizes the distribution of these labels and removes rows with insufficient labels.

```
1 labeled_data = data.copy()
2 labeled_data['label'] = [[0]] * len(labeled_data)
```



```

3 labeled_data.head()
4
5 tags = ['name', 'author', 'translator', 'publisher']

```

The script begins by copying the `data` DataFrame to `labeled_data` and initializes a new column, `label`, with zero values. It also defines a list of tags representing different entities in the dataset.

```

1 def get_label(tag):
2     if tag == 'name':
3         return 'Book'
4     elif tag == 'author':
5         return 'Author'
6     elif tag == 'translator':
7         return 'Translator'
8     elif tag == 'publisher':
9         return 'Publisher'
0     else:
1         return None

```

The `get_label` function maps each tag to its corresponding entity label. This function helps in converting tag names to more descriptive labels.

```

1 def convert_index_to_label(index, tag):
2     labels = []
3
4     for i in range(len(index)):
5         if index[i] == -1:
6             labels.append('O')
7         else:
8             if i == 0 or index[i - 1] == -1 or index[i] - index[i - 1] != 1:
9                 labels.append(f'B-{get_label(tag)}')
0             else:
1                 labels.append(f'I-{get_label(tag)}')
2
3     return labels

```

The `convert_index_to_label` function converts a list of indexes to labels. It uses the BIO (Beginning, Inside, Outside) tagging scheme to mark the beginning (B-) and continuation (I-) of entities. If the index is -1, it assigns the label O (Outside).

```

1 def combine_labels(labels):
2     global tags
3
4     result = []
5
6     for i in range(len(labels[tags[0]])):
7         name = labels[tags[0]][i] if len(labels[tags[0]]) > 0 else 'O'
8         author = labels[tags[1]][i] if len(labels[tags[1]]) > 0 else 'O'
9         translator = labels[tags[2]][i] if len(labels[tags[2]]) > 0 else 'O'
0         publisher = labels[tags[3]][i] if len(labels[tags[3]]) > 0 else 'O'
1
2         if name != 'O':
3             result.append(name)
4         elif author != 'O':
5             result.append(author)
6         elif translator != 'O':
7             result.append(translator)
8         elif publisher != 'O':
9             result.append(publisher)
0         else:
1             result.append('O')
2
3     return result

```

The `combine_labels` function combines labels for different tags into a single list. It prioritizes the `name` tag, followed by `author`, `translator`, and `publisher`, and assigns the label O if no other label is present.

```

1 def get_labels(row):
2     indexes = {
3         tag: []
4         for tag in tags
5     }
6     labels = {
7         tag: []
8         for tag in tags
9     }
10
11     for tag in tags:
12         cell = row[tag]
13         if cell == {np.nan}:
14             continue
15
16         filled_indexes = set()
17         for word in row['comment']:
18             try:
19                 current_index = cell.index(word)
20                 if current_index in filled_indexes:
21                     raise ValueError
22                 indexes[tag].append(current_index)
23                 filled_indexes.add(current_index)
24             except ValueError:
25                 indexes[tag].append(-1)
26
27         labels[tag] = convert_index_to_label(indexes[tag], tag)
28
29     return combine_labels(labels)
30
31 labeled_data['label'] = labeled_data.progress_apply(get_labels, axis=1)

```

The `get_labels` function generates labels for each row in the dataset. It finds the indexes of words in the `comment` column for each tag and converts them to labels using the `convert_index_to_label` function. The labels are then combined using the `combine_labels` function. The `progress_apply` method is used to apply this function to each row in `labeled_data`.

```

1 percentage_of_o = [len([tag for tag in label if tag == 'O']) / len(label) for label in
2     labeled_data['label']]
3 fig = px.histogram(percentage_of_o, title='Percentage of O in Each Label List Histogram')
4 fig.update_layout(showlegend=False)
5 fig.show()

```

The script calculates the percentage of 0 labels in each label list and creates a histogram to visualize this distribution using Plotly.

```

1 before_dropping = len(labeled_data)
2 labeled_data = labeled_data[[percentage < 0.99 for percentage in percentage_of_o]]
3 print(f'Dropped {before_dropping - len(labeled_data)} rows with more than 99% O.')
4 print(f'New length: {len(labeled_data)}')

```

Finally, the script removes rows from `labeled_data` where the percentage of 0 labels exceeds 99

Training

Explanations:

```

1 import tensorflow as tf
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import LabelEncoder
4 from sklearn.metrics import classification_report
5 from tensorflow.keras.preprocessing.sequence import pad_sequences
6 from tensorflow.keras.models import Sequential

```

```

7 from tensorflow.keras.layers import Embedding, LSTM, Bidirectional, SpatialDropout1D,
   InputLayer
8 from tensorflow.keras.callbacks import ModelCheckpoint, EarlyStopping, TensorBoard
9 from livelossplot.tf_keras import PlotLossesCallback
10
11 words = list(set([word for comment in labeled_data['comment'] for word in comment]))
12 words.append('STARTPAD')
13 words_size = len(words)
14 words_size
15
16 tags = sorted(list(set([tag for label in labeled_data['label'] for tag in label])))
17 tags_size = len(tags)
18 tags_size, tags
19
20 sentences = labeled_data.progress_apply(
21     lambda row: [(word, label) for word, label in zip(row['comment'], row['label'])],
22     axis=1)
23 sentences = sentences.progress_apply(lambda sentence: sentence + [('STARTPAD', 'O')])
24
25 word2idx = {word: idx for idx, word in enumerate(words)}
26 tag2idx = {tag: idx for idx, tag in enumerate(tags)}
27 idx2tag = {idx: tag for tag, idx in tag2idx.items()}
28
29 fig = px.histogram([len(sentence) for sentence in sentences], title='Sentence Length
   Histogram')
30 fig.update_layout(showlegend=False)
31 fig.show()

```

The script begins by importing necessary libraries and defining words and tags used in the labeled data. It calculates the unique words and tags and their respective sizes. It then creates a list of sentences from the labeled data and appends a special token 'STARTPAD' to each sentence. It also maps words and tags to unique indices.

```

1 # Based on the last hist.
2 max_length = 210
3
4 comments_less_than_max_length = len([sentence for sentence in sentences if len(
   sentence) <= max_length])
5 print(f'Number of comments less than {max_length}: {comments_less_than_max_length}')
6 print(f'Remainder: {len(sentences) - comments_less_than_max_length}')

```

The script determines the maximum sentence length based on the histogram of sentence lengths and counts the number of sentences that are shorter than this maximum length.

```

1 X = [[word2idx[word[0]] for word in sentence] for sentence in sentences]
2 X = pad_sequences(X, maxlen=max_length, padding='post', value=words_size-1)
3
4 y = [[tag2idx[word[1]] for word in sentence] for sentence in sentences]
5 y = pad_sequences(y, maxlen=max_length, padding='post', value=tag2idx['O'])
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state
   =42)
8 X_val, X_test, y_val, y_test = train_test_split(X_test, y_test, test_size=0.5,
   random_state=42)

```

Next, the script converts the sentences into sequences of word indices and pads them to ensure uniform length. It also converts the tags into sequences of tag indices and pads them similarly. The data is then split into training, validation, and test sets.

```

1 ## Create Bidirectional LSTM Model
2
3 model = Sequential()
4 model.add(InputLayer((max_length)))
5 model.add(Embedding(input_dim=words_size, output_dim=max_length, input_length=
   max_length))
6 model.add(SpatialDropout1D(0.1))

```

```

7 model.add(Bidirectional(LSTM(units=300, return_sequences=True)))
8
9 model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['
    accuracy'])
10 model.summary()
11
12 tf.keras.utils.plot_model(
13     model, to_file='resources/model.png', show_shapes=True, show_dtype=False,
14     show_layer_names=True, rankdir='LR', expand_nested=True, dpi=300,
15 )

```

The script defines a Bidirectional LSTM model using TensorFlow's Keras API. The model consists of an embedding layer, a spatial dropout layer, and a bidirectional LSTM layer. The model is compiled with the Adam optimizer and sparse categorical cross-entropy loss. A summary of the model architecture is displayed, and the model is visualized and saved as an image.

```

1 logdir = 'logs/'
2 tensorboard_callback = TensorBoard(log_dir=logdir)
3
4 model_checkpoint = ModelCheckpoint('resources/model.keras', monitor='val_loss',
5     verbose=1, save_best_only=True,
6     save_weights_only=True)
7 early_stopping = EarlyStopping(monitor='val_accuracy', min_delta=0, patience=1,
8     verbose=1, mode='max', baseline=None)
9 plot_losses = PlotLossesCallback()
10
11 history = model.fit(X_train[:10], y_train[:10], validation_data=(X_val[:10], y_val
12     [:10]), batch_size=32, epochs=3, verbose=1,
13     callbacks=[tensorboard_callback, model_checkpoint, early_stopping,
14     plot_losses])

```

The script sets up callbacks for TensorBoard logging, model checkpointing, early stopping, and live loss plotting. Finally, the model is trained on a subset of the data with the specified callbacks, batch size, and number of epochs. The training process includes validation on a separate validation set.

Evaluation

Explanations:

The following Python script is used to make predictions with a trained Bidirectional LSTM model and evaluate its performance. It includes functions for predicting and printing predictions, as well as evaluating the model on a test dataset.

```

1 def pred_y(x):
2     p = model.predict(np.array([x]))
3     p = np.argmax(p, axis=-1)
4     return p
5
6 def print_prediction(x, y):
7     print("{:15}{:5}\t {}".format("Word", "True", "Pred"))
8     print("-" * 30)
9     for w, true, pred in zip(x, y, pred_y(x)):
10         try:
11             print("{:15}{:5}\t {}".format(words[w-1], tags[true], tags[pred]))
12         except IndexError:
13             print("{:15}{:5}\t {}".format(words[w-1], tags[true], 'O'))

```

The `pred_y` function takes an input sequence `x`, makes a prediction using the trained model, and returns the predicted tags. The `print_prediction` function prints the words along with their true and predicted tags for a given input sequence.

```

1 sample_X = 'some persian text'
2 sample_X = higher_preprocess(elementary_preprocess(sample_X))
3 sample_X = [[word2idx[word] if word in word2idx else -1 for word in sample_X]]

```

```

4 sample_X = pad_sequences(sample_X, maxlen=max_length, padding='post', value=words_size-1)
5
6 p = pred_y(sample_X[0])
7 print("{:15}\t {}".format("Word", "Pred"))
8 print("-" * 30)
9 for w, pred in zip(sample_X[0], p[0]):
10     try:
11         print("{:15}\t {}".format(words[w-1], tags[pred]))
12     except IndexError:
13         print("{:15}\t {}".format(words[w-1], 'O'))

```

A sample Persian sentence is preprocessed using the `higher_preprocess` and `elementary_preprocess` functions. The sentence is converted to a sequence of word indices and padded to the maximum sequence length. Predictions are made for the sample sentence, and the words along with their predicted tags are printed.

```

1 i = np.random.randint(0, X_test.shape[0])
2 print("This is sentence:", i)
3 p = model.predict(np.array([X_test[i]]))
4 p = np.argmax(p, axis=-1)
5
6 print("{:15}{:5}\t {}".format("Word", "True", "Pred"))
7 print("-" * 30)
8 for w, true, pred in zip(X_test[i], y_test[i], p[0]):
9     try:
10         print("{:15}{:5}\t {}".format(words[w-1], tags[true], tags[pred]))
11     except IndexError:
12         print("{:15}{:5}\t {}".format(words[w-1], tags[true], 'O'))

```

A random sentence from the test set is selected, and predictions are made for it. The words along with their true and predicted tags are printed.

```

1 y_pred = model.predict(X_test)
2
3 y_pred_index = np.argmax(y_pred, axis=-1)
4
5 X_test_list = X_test.tolist()
6 y_test_list = y_test.tolist()
7 y_pred_list = y_pred_index.tolist()
8
9 start_pad_index = word2idx['STARTPAD']
10 # progress_bar = tqdm(range(len(X_test)))
11 for i in range(len(X_test_list)):
12     x_start_pad = X_test_list[i].index(start_pad_index)
13     X_test_list[i] = X_test_list[i][x_start_pad:]
14     y_test_list[i] = y_test_list[i][x_start_pad:]
15
16     y_pred_list[i] = y_pred_list[i][x_start_pad:]
17     y_pred_list[i] = [tag if tag <= 8 else 8 for tag in y_pred_list[i]]

```

The script then makes predictions for the entire test set. It converts the predictions to tag indices and removes the padding from the sequences. It ensures that any tag indices greater than the maximum tag index are set to the maximum tag index.

```

1 y_test_flatten = []
2 for y in y_test_list:
3     y_test_flatten.extend(y)
4
5 y_pred_flatten = []
6 for y in y_pred_list:
7     y_pred_flatten.extend(y)
8
9 print(classification_report(y_test_flatten, y_pred_flatten))

```

Finally, the true and predicted tags are flattened, and a classification report is generated to evaluate the model's performance.

Figure 7: Percentage of O in Each Label List Histogram

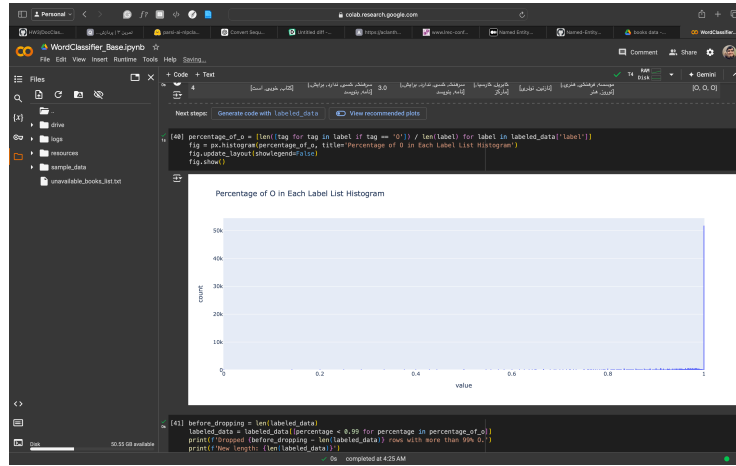


Figure 8: Tag Distribution with O



Figure 9: Tag distribution without O

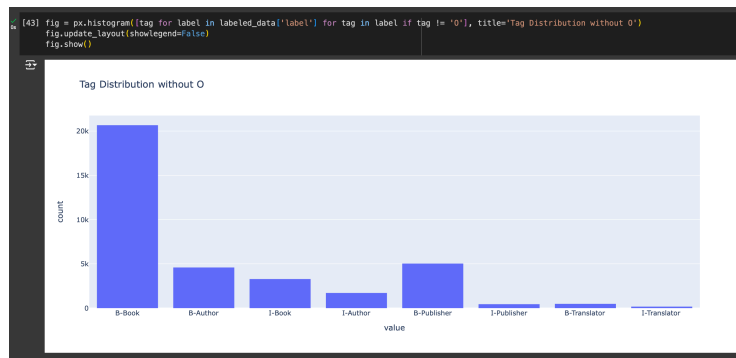


Figure 10: Sentence (Comment) Length Histogram

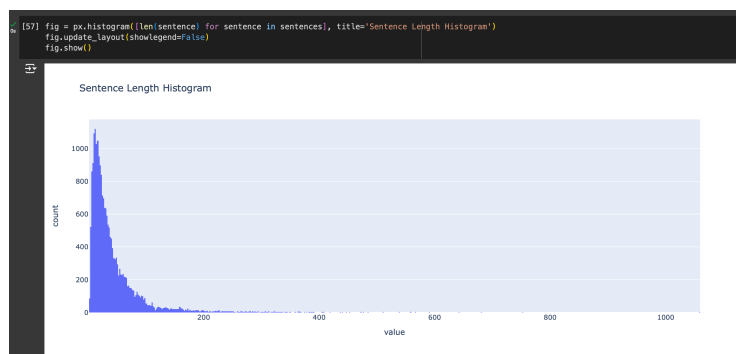


Figure 11: LSTM Model

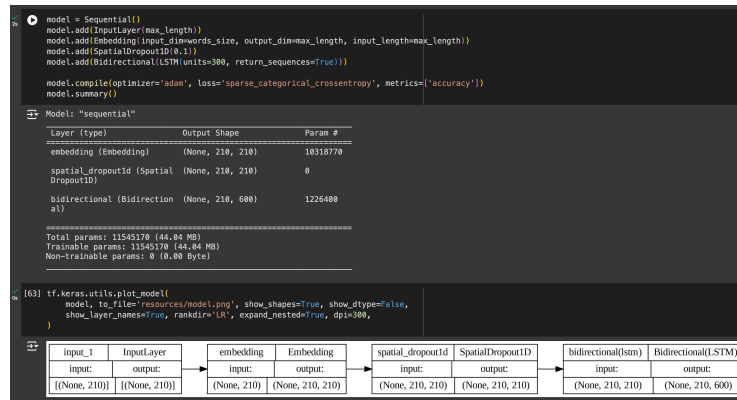


Figure 12: Accuracy and Loss During Training

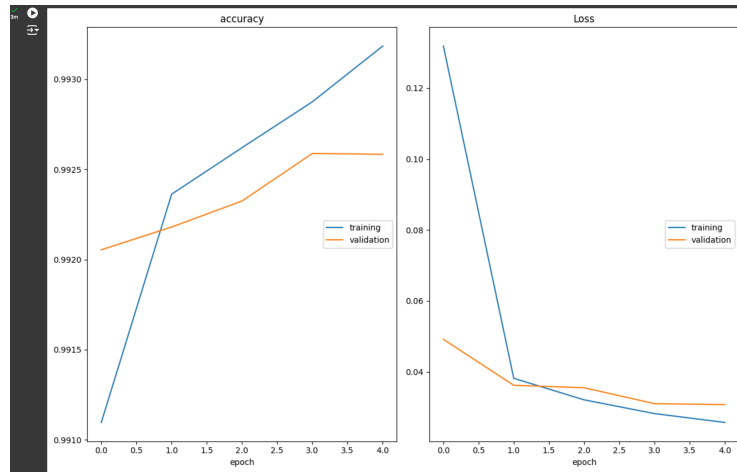
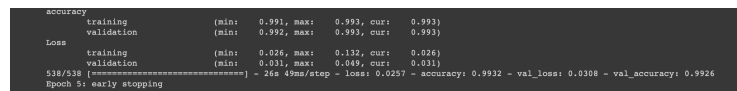


Figure 13: Text of Running the Model



Word Classifier - Transformer Model

The first part is exactly similar to the previous section. The only difference is that BERT is used as a normalizer. So, we explain only the transformer model in this section.

Explanations:

```

1 from transformers import BertConfig
2 from transformers import TFBertForTokenClassification, InputFeatures
3 import tensorflow as tf
4 from livelossplot.tf_keras import PlotLossesCallback
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import classification_report, accuracy_score, confusion_matrix,
  precision_score, recall_score, f1_score

```

First, the necessary libraries are imported. We use the Hugging Face Transformers library for BERT, TensorFlow for training, and various scikit-learn utilities for data handling and evaluation.

```

1 tags = sorted(list(set([tag for label in labeled_data['label'] for tag in label])))
2 tags_size = len(tags)
3 tags_size, tags
4
5 tag2idx = {tag: idx for idx, tag in enumerate(tags)}
6 idx2tag = {idx: tag for tag, idx in tag2idx.items()}

```

Figure 14: The Sentence in the Doc

[illegible]

Figure 15: Classification Report

```

precision    recall    f1-score   support

 0         0.65         0.64         0.65         442
 1         0.59         0.11         0.18        2033
 2         0.53         0.18         0.27         483
 3         0.80         0.00         0.00          48
 4         0.50         0.20         0.26         155
 5         0.39         0.15         0.22         319
 6         0.80         0.00         0.00          42
 7         0.80         0.00         0.00          12
 8         0.97         0.99         0.98       76751

 accuracy          0.96      80277
 macro avg         0.39      0.26      0.29      80277
 weighted avg      0.95      0.96      0.95      80277

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning:
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning:
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning:
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

```

The tags used for labeling are extracted and encoded into indices. The `tag2idx` and `idx2tag` dictionaries map tags to indices and vice versa.

```
1 X_train, X_test, y_train, y_test = train_test_split(labeled_data['comment'], labeled_data[
    'label'], test_size=0.2)
2 X_val, X_test, y_val, y_test = train_test_split(X_test, y_test, test_size=0.5)
3
4 train = pd.DataFrame({'comment': X_train, 'label': y_train})
5 val = pd.DataFrame({'comment': X_val, 'label': y_val})
```


Figure 16: Confusion matrix

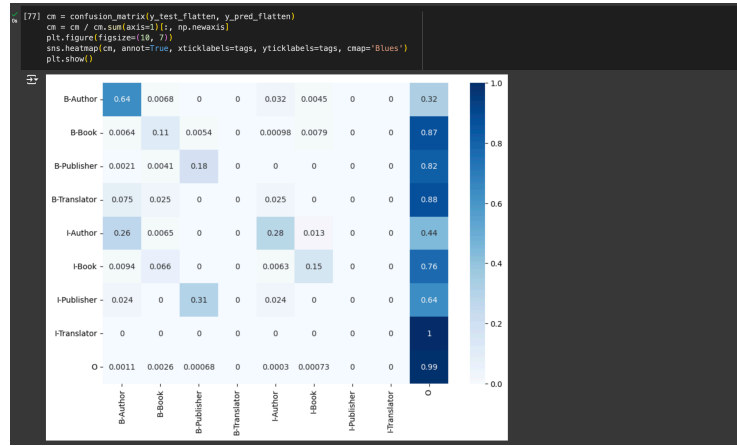


Figure 17: Metrics

```
Accuracy: 0.9594279886171132
Micro Precision: 0.9594279886171132
Macro Precision: 0.3934988376208633
Weighted Precision: 0.9455743243739555

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning:
Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning:
Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

Micro Recall: 0.9594279886171132
Macro Recall: 0.26183399222539545
Weighted Recall: 0.9594279886171132

Micro F1-score: 0.9594279886171132
Macro F1-score: 0.2949953149539064
Weighted F1-score: 0.9484172831756426
```

```
6 test = pd.DataFrame({'comment': X_test, 'label': y_test})
7
8 train = train[:64]
9 val = val[:64]
10 test = test[:64]
11
12 train = train.reset_index(drop=True)
13 val = val.reset_index(drop=True)
14 test = test.reset_index(drop=True)
```

The labeled data is split into training, validation, and test sets using an 80-20 split, with the test set further split in half for validation. For demonstration purposes, only the first 64 rows of each dataset are used.

```
1 def add_padding(data: pd.DataFrame, max_length: int):
2     progress_bar = tqdm(range(len(data)))
3
4     for i in progress_bar:
5         data['comment'][i] = data['comment'][i][:max_length]
6         data['label'][i] = data['label'][i][:max_length]
7
8         if len(data['comment'][i]) < max_length:
9             data['comment'][i] += ['PAD'] * (max_length - len(data['comment'][i]))
10            data['label'][i] += ['O'] * (max_length - len(data['label'][i]))
11
12    return data
13
14    max_length = 128
15
16    train = add_padding(train, max_length)
17    val = add_padding(val, max_length)
18    test = add_padding(test, max_length)
```

The `add_padding` function ensures all comments and labels in the data are padded to a fixed length (`max_length`). Padding is necessary for BERT to handle inputs of consistent length.

```
1 ## Convert to GLUE Format
```

```

3 class InputExample:
4     def __init__(self, guid, text_a, text_b=None, label=None):
5         self.guid = guid
6         self.text_a = text_a
7         self.text_b = text_b
8         self.label = label

```

The `InputExample` class is defined to structure each data example with a unique ID, the text, and the label.

```

1 def convert_data_to_examples(data: pd.DataFrame):
2     examples = []
3
4     for i in range(len(data)):
5         guid = i
6         text_a = data['comment'][i]
7         label = data['label'][i]
8         examples.append(InputExample(guid=guid, text_a=text_a, label=label))
9
10    return examples
11
12
13 def convert_examples_to_features(examples, tokenizer, max_length, task=None):
14     features = []
15
16     for example in examples:
17         input_dict = tokenizer.encode_plus(
18             example.text_a,
19             add_special_tokens=True,
20             max_length=max_length,
21             return_token_type_ids=True,
22             return_attention_mask=True,
23             padding='max_length',
24             truncation=True
25         )
26
27         input_ids = input_dict['input_ids']
28         attention_mask = input_dict['attention_mask']
29         token_type_ids = input_dict['token_type_ids']
30
31         label = example.label
32
33         if task is not None:
34             label = [tag2idx[tag] for tag in label]
35
36         features.append(
37             InputFeatures(
38                 input_ids=input_ids,
39                 attention_mask=attention_mask,
40                 token_type_ids=token_type_ids,
41                 label=label
42             )
43         )
44
45     return features
46
47
48 def convert_data_to_features(data: pd.DataFrame, tokenizer, max_length, task=None):
49     examples = convert_data_to_examples(data)
50     return convert_examples_to_features(examples, tokenizer, max_length, task)

```

Functions are provided to convert data into examples and then into features suitable for BERT. These functions tokenize the text, add special tokens, and ensure padding and truncation to the maximum length.

```

1 def convert_data_to_tf_dataset(data: pd.DataFrame, tokenizer, max_length, task=None):
2     features = convert_data_to_features(data, tokenizer, max_length, task)
3
4     all_input_ids = []

```

```

5 all_attention_masks = []
6 all_token_type_ids = []
7 all_labels = []
8
9 for feature in features:
10     all_input_ids.append(feature.input_ids)
11     all_attention_masks.append(feature.attention_mask)
12     all_token_type_ids.append(feature.token_type_ids)
13     all_labels.append(feature.label)
14
15 return (
16     tf.data.Dataset.from_tensor_slices((
17         {
18             'input_ids': all_input_ids,
19             'attention_mask': all_attention_masks,
20             'token_type_ids': all_token_type_ids
21         },
22         all_labels
23     ))
24 )

```

The `convert_data_to_tf_dataset` function converts the features into a TensorFlow dataset suitable for training.

```

1 max_length = 128
2 task = 'ner'
3
4 train_dataset = convert_data_to_tf_dataset(train, tokenizer, max_length, task=task)
5 val_dataset = convert_data_to_tf_dataset(val, tokenizer, max_length, task=task)
6 test_dataset = convert_data_to_tf_dataset(test, tokenizer, max_length, task=task)
7
8 train_dataset = train_dataset.shuffle(100).batch(32).repeat(2)
9 val_dataset = val_dataset.batch(64)
10 test_dataset = test_dataset.batch(64)
11
12 train_dataset.element_spec, val_dataset.element_spec, test_dataset.element_spec

```

The datasets for training, validation, and testing are prepared, shuffled, and batched appropriately.

```

1 config = BertConfig.from_pretrained(MODEL_NAME_OR_PATH, num_labels=tags_size, id2label=
    idx2tag, label2id=tag2idx)
2
3 model = TFBertForTokenClassification.from_pretrained(MODEL_NAME_OR_PATH, config=config)
4 model.summary()
5
6 multi_label_loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
7 model.compile(optimizer='adam', loss=multi_label_loss, metrics=['accuracy'])
8
9 plot_losses = PlotLossesCallback()
10
11 model.fit(train_dataset, epochs=1, validation_data=val_dataset, callbacks=[plot_losses])

```

The BERT model is configured and initialized for token classification. The model is compiled with a loss function and accuracy metric, and training is initiated with the `fit` method, including a callback for plotting the training losses.

Results and Imporatnt Notes

We must note that we tried to use Hazm's Informal Normalizer but it was very slow. So, e customized the Hazm Informal Normalizer. It returns a 3D array in which the third dimension is a candidate for the formal word. We picked the first candidate directly by customizing the Hazm library. Now, it returns a 2D array, and its performance has improved.

For labeling, we iterate over comments and check if the word is in, as an example, in the book name or not. If it was in the book names list, we appended the index of the word in the book names list to a list. So now we have a list of indexes. Then, we tried to find all consecutive ascending sequences in the list. Then we tag them by BIO, or if they are -1, we tag them by O.

Another important note is related to the distribution of tags. The "O" tag completely dominates other tags by count, and therefore, it impacts our results. You can see the distribution in the following figures.

Figure 18: Percentage of O in Each Label List Histogram

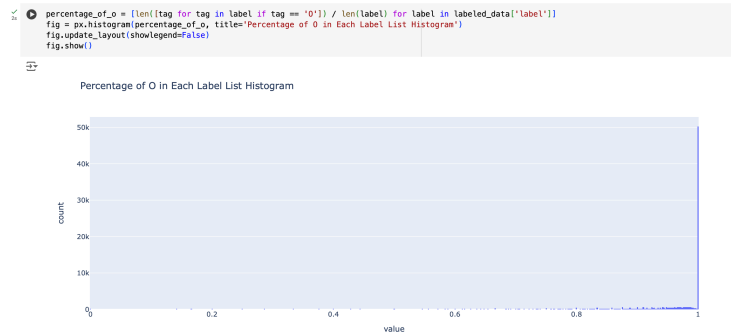


Figure 19: Tag Distribution with O

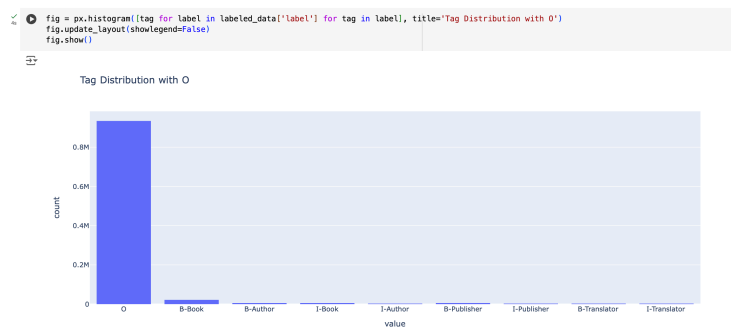


Figure 20: Tag distribution without O

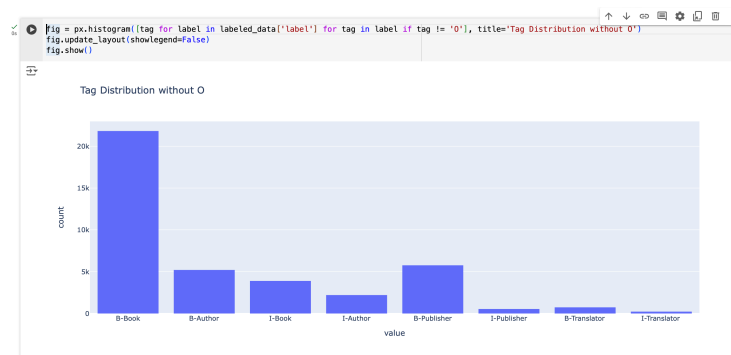


Figure 21: Transfomer Model

```
[71] config = BertConfig.from_pretrained(MODEL_NAME_OR_PATH, num_labels=tags_size, id2label=id2tag2, label2id=label2tag2)
/usr/local/lib/python3.10/dist-packages/huggingface_hub/file_download.py:1132: FutureWarning: 'resume_download' is deprecated and will be removed in warnings.warn()

[72] model = TFBertForTokenClassification.from_pretrained(MODEL_NAME_OR_PATH, config=config)
model.summary()
```

f_model.hf: 100% ██████████ 965M/965M [01:08<00:00, 15.0MB/s]

All model checkpoint layers were used when initializing TFBertForTokenClassification.

Some layers of TFBertForTokenClassification were not initialized from the model checkpoint at Hooosvarab/tbert-fa-base-uncased and are newly initialized. You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Model: "tf_bert_for_token_classification"

Layer (type)	Output Shape	Param #

bert (TFBertMainLayer)	multiple	162250752
dropout_37 (Dropout)	multiple	0 (unused)
classifier (Dense)	multiple	6921

Total params: 162250752 (610.96 MB)
Trainable params: 162250752 (610.96 MB)
Non-trainable params: 0 (0.00 Byte)

Figure 22: Accuracy and Loss During Training

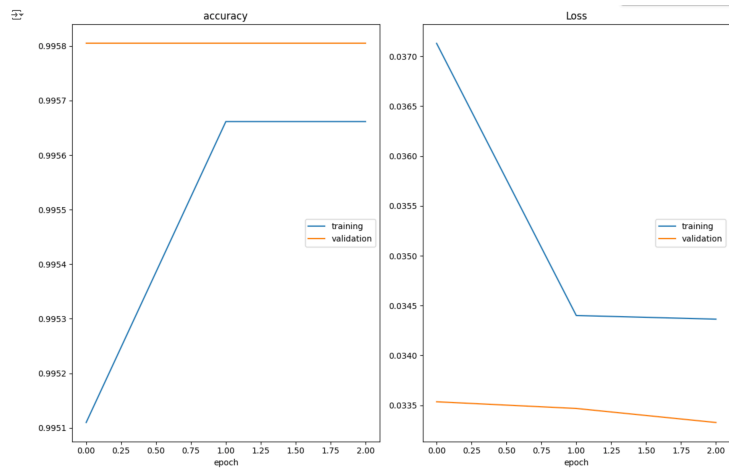


Figure 23: Some Evaluations

```

accuracy
  training (min: 0.995, max: 0.996, cur: 0.996)
  validation (min: 0.996, max: 0.996, cur: 0.996)
Loss
  training (min: 0.034, max: 0.037, cur: 0.034)
  validation (min: 0.033, max: 0.034, cur: 0.033)
3720/3720 [====] - 3200 steps - loss: 0.0344 - accuracy: 0.9957 - val_loss: 0.0333 - val_accuracy: 0.9958
tf.keras.src.callbacks.History at 0x7f17613b2d20

```

Figure 24: Classification Report

```
[77] model.evaluate(test_dataset)

117/117 [=====] - 60s 514ms/step - loss: 0.0341 - accuracy: 0.9957
0.8340860114337444305, 0.9956721067428589)

[78] y_pred = model.predict(test_dataset)
y_pred = np.argmax(y_pred.logits, axis=-1)
y_pred = [idx*tag[idx] for idx in row] for row in y_pred]
y_true = test['label'].values

117/117 [=====] - 70s 514ms/step

y_pred = [tag for row in y_pred for tag in row]
y_true = [tag for row in y_true for tag in row]

print(classification_report(y_true, y_pred))

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-defined and be
warn_pr(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-defined and be
warn_pr(average, modifier, msg_start, len(result))

          precision    recall  f1-score   support

B-Author      0.00      0.00      0.00        586
B-Book        0.00      0.00      0.00       2282
B-Publisher   0.00      0.00      0.00        638
B-Translator  0.00      0.00      0.00         74
I-Author      0.00      0.00      0.00        286
I-Book        0.00      0.00      0.00        368
I-Publisher   0.00      0.00      0.00         36
I-Translator  0.00      0.00      0.00         19
0          1.00      1.00      1.00      948071

 accuracy      0.99          1.00          1.00      952192
 macro avg     0.11      0.11      0.11      952192
 weighted avg  0.99      1.00      0.99      952192

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-defined and be
```

Figure 25: Confusion matrix

