

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

```

## Explanations:

The main part of extracorp.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

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

### 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,
      positive_threshold, neutral_threshold))
20
21    # Combine the 'comment' and 'sentiment' columns into a single DataFrame
22    df = pd.concat([labeled_data['comment'], labeled_data['sentiment']], axis=1)
23
24    # Separate the DataFrame into three classes based on sentiment
25    positive = df[df['sentiment'] == 'positive']
26    neutral = df[df['sentiment'] == 'neutral']
27    negative = df[df['sentiment'] == 'negative']
28
29    # 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, 2), 'positive_threshold': 4,
'neutral_threshold': 2}
Test accuracy of Logistic Regression model: 0.5416666666666666

```

	precision	recall	f1-score	support
negative	0.47	0.56	0.52	117
neutral	0.54	0.59	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: 0.4713541666666667

```

	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: 0.5083333333333333

```

	precision	recall	f1-score	support
negative	0.53	0.59	0.56	1020
neutral	0.52	0.52	0.52	907
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

Note: Some parts of explanations are written with the help of ChatGPT to help us write clear and unambiguous descriptions.