

$\begin{array}{c} {\rm Introduction\ to\ Artificial\ Intelligence} \\ {\rm SWE} 3011_41 \\ {\rm Fall\ } 2023 \end{array}$

Project Anonymous ACL submission

Sungkyunkwan University

Abstract

Within the context of this assignment, a substantial hands-on encounter with applications of artificial intelligence will be attained, with a particular focus on two pivotal tasks: traditional supervised learning methodologies and prompt engineering. Upon the culmination of this assignment, an adept understanding of the practical implementation of artificial intelligence techniques for the purpose of classification shall be cultivated. Furthermore, through collaborative participation in the curation of prompts and the execution of experiments, the acquisition of insights and skills that are notably advantageous for a group project is anticipated.

1 Introduction

In this lab journal, the primary objective is to offer a comprehensive practical experience in the field of artificial intelligence, with a specific focus on supervised learning. Upon the conclusion of this project, a solid and pragmatic understanding of traditional machine learning techniques for text classification is expected to be acquired. Furthermore, the utilization of Langchain for prompt engineering tasks, especially in the context of Large Language Models (LLMs), will have been explored. This multifaceted approach will equip us with a versatile skill set and a deeper insight into the practical aspects of artificial intelligence.

This assignment categorizes tasks into three main problems:

1. Supervised Text Classification using Traditional Machine Learning Methods: The task involves implementing a text classification process using traditional machine learning methods. A dataset and skeleton code utilizing the scikit-learn library, with specific sections left blank for completion, will be provided. The responsibility is to finalize the code, conduct experiments, and evaluate to derive results.

- 2. Prompt Engineering for LLMs via Langchain Framework: Using Langchain, a tool designed for crafting prompts for language models, a series of prompts will be created to elicit specific and varied responses. Baseline code for conducting prompt engineering experiments with ChatGPT and Large Language Models on Hugging Face Hub is provided. The task involves completing the code and optimizing the prompts for sentiment classification.
- 3. LATEX: The third focus of this assignment is report writing using LATEX, a document preparation system that aids in creating structured documents. This part helps students develop skills in professional communication and documentation within the AI field.

This journal serves as a documentation of the academic journey in addressing these AI challenges and acquiring practical skills in the field of artificial intelligence.

2 Dataset

The GLUE benchmark stands as a compendium of diverse Natural Language Processing (NLP) tasks, employed to assess the efficacy of language models. A significant component of the GLUE benchmark, known as Stanford Sentiment Treebank V2 (SST-2), is dedicated to evaluating the performance of models in sentiment classification on the SST-2 dataset and other related tasks, serving as a litmus test for their broader language comprehension and reasoning capabilities (Wang et al., 2018).

The SST-2 dataset is a renowned benchmark dataset in the realm of NLP, frequently employed for sentiment analysis undertakings. An extension of the Stanford Sentiment Treebank (SST), which originally featured fine-grained sentiment labels encompassing categories such as very positive, positive, neutral, negative, and very negative. SST-2 consists of sentences sourced from movie

reviews, each meticulously annotated with binary labels signifying either positive or negative sentiment (Socher et al., 2013).

3 Methology

For both Task 1 and Task 2, distinct methods have been chosen to address their specific objectives.

For Task 1, which centers on supervised text classification, the chosen trio of methods comprises logistic regression, random forest, and naive Bayes classifiers. These algorithms were meticulously selected for their suitability and proven effectiveness in tackling text classification challenges.

For Task 2, which focuses on prompt engineering for Large Language Models, LangChain was chosen as the tool to design both zero-shot and few-shot prompts. These prompts were to be tested using ChatGPT and the Large Language Models available on the Hugging Face Hub in order to solve the sentiment classification problem on the SST-2 dataset.

It's important to note that, for this task, the decision was made to abstain from using OpenAI models. Instead, the focus was solely on experimenting with freely available models on the Hugging Face Hub. This approach allowed exploration of the capabilities of large language models within the constraints of free resources.

4 Experiments

4.0.1 Task 1

In the experimental phase of this project, a comprehensive analysis was undertaken to evaluate the effectiveness of three distinct machine learning models: Logistic Regression, Random Forest, and the Naive Bayes Classifier. The approach adopted was two-fold, entailing not only model selection but also the intricate process of hyperparameter tuning and cross-validation. In this context, the Logistic Regression model demonstrated remarkable performance, optimizing its hyper-

parameters to yield the most accurate re-Specifically, the observation revealed that a regularization strength parameter, denoted as 'C,' was fine-tuned to the value of 100, and the regularization term 'penalty' In essence, 'C' was conwas set to 'l2'. trolled to determine the degree of regularization, with larger values signifying less regularization and a potentially higher sensitivity to the training data, while 'penalty' indicated the type of regularization applied, with 'l2' corresponding to the L2 penalty. On the other hand, when the Random Forest model was employed, the hyperparameter optimization revealed that the ideal configuration consisted of a maximum depth set to 'None' and an ensemble size, or 'n_estimators,' fixed at 200. Here, 'max_depth' signified the depth of each tree in the Random Forest, with 'None' indicating no restrictions on the depth. As for 'n_estimators,' it represented the number of trees within the ensemble. Lastly, in the case of the Naive Bayes Classifier, it was determined that the hyperparameter 'alpha' was most effective at a value of 0.1, indicating the Laplace smoothing factor for the model. Through conducting this detailed experimentation and elucidating the significance of these parameter choices, the objective was not only to optimize model performance but also to gain a profound understanding of the underlying mechanisms that govern these machine learning algorithms.

4.0.2 Task 2

For the experimental part of Task 2, both zero-shot and few-shot prompts were prepared in order to utilize the Hugging Face Hub's Large Language Models' ability to perform sentiment analysis of the SST-2 dataset. "Zero-shot prompt" refers to prompt engineering where the model is not trained with examples on how to perform the given task, relying only on its existing knowledge, while "few-shot prompt" refers to prompts where the model is trained with specific examples on how to complete the task. For both the zero-shot and few-shot prompts, the model was

given the following prefix, as the instruction for its task of sentiment classification:

"You are a sentiment classifier. Classify the given text input as "positive" or "negative" based on its tone. Please read each sentence carefully and determine its sentiment based on common definitions.

- <positive>: A positive sentiment indicates a favorable or optimistic tone. It expresses approval, satisfaction, or positive emotions.
- <negative>: A negative sentiment indicates an unfavorable or pessimistic tone. It expresses disapproval, dissatisfaction, or negative emotions.

Please read each sentence and classify it as either "positive" or "negative" based on the provided definitions."

Apart from the prefix that serves as the task description, for the few-shot prompt, the model was trained with a provided dataset that shares similar characteristics to the SST-2 dataset. Due to the constraints of using only free resources, it was not possible to provide the Large Language Model with more than 27 training examples per execution. Since the training dataset was larger than that, a different set of examples was used for each execution. Each of the examples in the set was selected randomly out of all the available ones in the dataset; however, in order to maximize the accuracy of the Large Language Model predictions, a constraint was applied to maintain consistent proportions of positive and negative instances in each set. This approach was taken to prevent any imbalance between positive and negative examples in the generated sets, which could impact the model's training. After executing the prompts multiple times, it became apparent that the accuracy of few-shot prompts varied with each iteration due to changing examples influencing the model. Consequently, three iterations were conducted to capture the accuracy range, constrained by limited resources. In the case of zero-shot prompts, a single iteration was performed for assessment. The resulting accuracy of the zero-shot prompt and the accuracy ranges of the few-shot prompts were compared in chapters "Results" and "Discussions".

5 Results

5.0.1 Task 1

For the logistic regression model, the achieved accuracy was 0.77. The classification report demonstrated that for class 0, the precision was 0.83, recall was 0.72, and the F1-score was 0.77, with a support of 54 instances. For class 1, the precision was 0.72, recall was 0.83, and the F1-score was 0.77, with a support of 46 instances. The overall accuracy for this model was 0.77.

Accuracy: 0.7				
Classification	on Report:			
	precision	recall	f1-score	support
0	0.83	0.72	0.77	54
1	0.72	0.83	0.77	46
accuracy			0.77	100
macro avg	0.77	0.77	0.77	100
weighted avg	0.78	0.77	0.77	100

Figure 1: Result of Logistic Regression model

In the case of the random forest model, the accuracy obtained was 0.71. The classification report provided the following results: for class 0, the precision was 0.75, recall was 0.70, and the F1-score was 0.72, with a support of 54 instances. For class 1, the precision was 0.67, recall was 0.72, and the F1-score was 0.69, with a support of 46 instances. The overall accuracy for this model was 0.71.

Accuracy: 0	71				
Classificat		•			
	pr	ecision	recall	f1-score	support
	0	0.75	0.70	0.72	54
	1	0.67	0.72	0.69	46
accurac	У			0.71	100
macro av	/g	0.71	0.71	0.71	100
weighted av	/g	0.71	0.71	0.71	100

Figure 2: Result of Random Forest model

Regarding the naive Bayes model, the achieved accuracy was 0.70. The classification report revealed that for class 0, the precision was 0.74, recall was 0.69, and the F1-score was 0.71, with a support of 54 instances. For class 1, the precision was 0.66,

recall was 0.72, and the F1-score was 0.69, with a support of 46 instances. The overall accuracy for this model was 0.70.

Accuracy: 0.7				
Classificatio			C4	
	precision	recall	f1-score	support
0	0.74	0.69	0.71	54
1	0.66	0.72	0.69	46
accuracy			0.70	100
macro avg	0.70	0.70	0.70	100
weighted avg	0.70	0.70	0.70	100

Figure 3: Result of Naive Bayes model

5.0.2 Task 2

There was a significant difference between the accuracy of sentiment classification obtained from the Large Language Models from the zero-shot and few-shot prompts. The zero-shot prompt had a prediction accuracy of 51%, while the few-shot prompt yielded a prediction accuracy in the range of 82-89%. The most important thing to note is the greater accuracy of the few-shot prompt for the task of sentiment classification. Additionally, the figures are included in the appendix for more detailed examination.



Figure 4: Result of prediction using zero-shot prompt

0, 1, 1,	0, 0, 0, 1	, 0, 0, 1,		-1, 0, 0, 0, 0, 0, 1, 0, -1, 1, 0, 0, 1, 1, 0, 0,

Figure 5: Result of prediction with the first execution of few-shot prompt

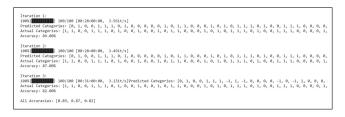


Figure 6: Result of prediction with subsequent executions of few-shot prompt

6 Discussions

6.0.1 Task 1

Accuracy was chosen as the primary evaluation metric for Task 1, aligning with the approach taken in the GLUE benchmark. To evaluate the models, the provided test_dataset was employed.

It is evident that the logistic regression model outperformed the other two models, namely random forest and naive Bayes, with the highest accuracy of 0.77. The logistic regression model exhibited a balanced performance in terms of precision, recall, and F1-scores for both classes 0 and 1. This indicates its effectiveness in handling text classification tasks. On the other hand, the random forest model achieved an accuracy of 0.71. While it provided competitive results, its performance was slightly lower compared to logistic regression. The classification report demonstrated reasonable precision, recall, and F1-scores for both classes, making it a viable choice but not the top-performing one.

Lastly, the naive Bayes model exhibited an accuracy of 0.70, which was consistent with random forest but slightly lower than logistic regression. The classification report also showed balanced precision, recall, and F1-scores for both classes.

6.0.2 Task 2

Accuracy was also chosen as the primary evaluation metric for Task 2, in order to stay consistent with the findings of Task 1, as well as provide insight into the capability of the Large Language Model to solve the sentiment classification problem on the SST-2 dataset.

The most important result was the difference in accuracy when employing few-shot prompts instead of zero-shot prompts. This was consistent with our expectation that sentiment classification yields better results when the model is trained with specific examples rather than when the task is simply performed based on the model's pre-existing knowledge. This poses the idea that few-shot prompts are superior for solving sentiment classification problems. Because of this, choosing a proper

dataset and improving the quality, diversity and proportionality of positive and negative emotion examples is an important part of solving problems like this. Lastly, while there is a positive correlation between the amount of examples provided and the accuracy of the predictions, it is important to not overtune the model with an excessive amount of examples.

7 Conclusion

7.0.1 Task 1

In conclusion, this project delved into the application of various machine learning models, including logistic regression, random forest, and naive Bayes, for the task of text classification. Through rigorous experimentation and evaluation, it was observed that the logistic regression model emerged as the top-performing model, achieving the highest accuracy and demonstrating balanced precision, recall, and F1-scores for both classes. This emphasizes the significance of selecting an appropriate machine learning model tailored to the task at hand.

Despite the success of the logistic regression model, there are potential areas for improvement. First and foremost, further exploration of feature engineering techniques and text preprocessing methods could enhance the performance of all models. Additionally, the incorporation of more advanced machine learning models or deep learning architectures may lead to improved results in text classification tasks. Furthermore, the dataset size and diversity could be increased to bolster the models' ability to generalize across different textual data.

Moreover, an investigation into hyperparameter tuning techniques, such as grid search or Bayesian optimization, could be beneficial to fine-tune the models further. Finally, ensembling methods, which combine the strengths of multiple models, may provide an opportunity to boost overall performance.

In summary, while logistic regression proved to be the most effective model in this specific context, there remains room for enhancements through various avenues, including feature engineering, model selection, and hyperparameter tuning. This project serves as a valuable foundation for future research and applications in the field of text classification and natural language processing.

7.0.2 Task 2

In conclusion, Task 2 of this project allowed us to explore prompt engineering and how it is a fundamental part of AI projects. In this case, a sentiment classification problem was solved through zero-shot and few-shot prompts. It was undeniable that the few-shot prompt performed significantly better than the zero-shot prompt for this task. Not only that, but is also relevant to point out that the few-shot prompt achieved a range of accuracy at sentiment classification superior to that achieved by the machine learning methods utilized in Task 1. This could be attributed to many factors, most notably the nature of the SST-2 dataset. This dataset is text based, which makes it well suited to the use of prompt engineering. The machine learning methods explored in Task 1 are better suited for dataframes with attributes, especially Logistic Regression, which had the biggest accuracy rate in Task 1 and works better with numerical datasets.

References

Momojit Biswas. 2020. Prompt engineering: Unleashing the power of few-shot learning for accurate spam detection. Medium.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment tree-bank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.

Code

$Task_{-}1$

```
# -*- coding: utf-8 -*-
"""SWE3011_41_Task1.ipynb
2
3
    Automatically generated by Colaboratory.
4
5
    Original file is located at
6
        https://colab.research.google.com/drive/1RKBXBhSWJVU9MqRcyZBCIKWBRx6SMf_d
9
    # SWE3011_41 Task1
10
    **Supervised Text Classification using traditional machine learning methods**
11
12
13
    1. Complete all the functions given.
    2. Conduct various experiments including hyper-parameter tuning, cross validation, etc.
14
15
    3. Write a report on the analysis of experiment results.
16
    **0. Installation**
17
18
19
    **1. Load Dataset**
20
21
22
    pip install datasets
23
24
25
    Evaluation should be done using **provided test dataset **""
26
^{27}
    from datasets import load_dataset
28
    train_ds = load_dataset("glue", "sst2", split="train")
29
30
31
     # Evaluation should be done using test_ds
32
    test_ds = load_dataset("csv", data_files="./test_dataset.csv")['train']
33
    """**2. Preparing Dataset**""
34
35
    from sklearn.feature_extraction.text import TfidfVectorizer
36
37
38
    def transform_data(X_train, X_test):
39
40
        Input:
41
        - X_train, X_test: Series containing the text data for training and testing respectively.
42
43
44
         - X_train_tfidf, X_test_tfidf: Transformed text data in TF-IDF format for training and testing
        \hookrightarrow respectively.
45
        - vectorizer: Fitted TfidfVectorizer object.
46
47
        # TODO: Convert the text data to TF-IDF format and return the transformed data and the vectorizer
48
         # Create a TfidfVectorizer
49
50
        vectorizer = TfidfVectorizer()
51
52
        # Fit and transform X_train
        X_train_tfidf = vectorizer.fit_transform(X_train)
53
54
55
          Transform X_test
56
        X_test_tfidf = vectorizer.transform(X_test)
57
        58
        return X_train_tfidf, X_test_tfidf, vectorizer
59
    X train, v train = train ds['sentence'], train ds['label']
60
    X_test, y_test = test_ds['sentence'], test_ds['label']
61
    X_train_tfidf, X_test_tfidf, vectorizer = transform_data(X_train, X_test)
62
63
    """**3. Train**"""
64
65
    from sklearn.model_selection import GridSearchCV
66
    from sklearn.linear_model import LogisticRegression
67
    from sklearn.ensemble import RandomForestClassifier
68
69
    from sklearn.naive_bayes import MultinomialNB
70
71
    def logistic_regression(X_train_tfidf, y_train):
72
73
        Input:
74
         - X_train_tfidf: Transformed text data in TF-IDF format for training.
75
        - y_train: Series containing the labels for training.
76
77
        Output:
        - clf: Trained Logistic Regression model.
78
79
        80
        # Define a logistic regression classifier with max_iter
```

```
82
         clf = LogisticRegression(max_iter=1000)
83
          # Define hyperparameters to tune
84
         param_grid = {
   'C': [0.001, 0.01, 0.1, 1, 10, 100],
   'penalty': ['11', '12']
85
 86
 87
 88
 89
90
          {\it \# Perform grid search with cross-validation to find the best hyperparameters}
         grid_search = GridSearchCV(clf, param_grid, cv=5, scoring='accuracy')
91
92
         grid_search.fit(X_train_tfidf, y_train)
 93
94
          # Print the best parameters
         print("Best Parameters: ", grid_search.best_params_)
95
96
          # Get the best model with the optimal hyperparameters
97
98
         clf = grid_search.best_estimator_
99
100
          # Train the final model with the entire training data
101
          clf.fit(X_train_tfidf, y_train)
102
          103
          return clf
104
105
     def random_forest(X_train_tfidf, y_train):
106
107
108
          - X_{\rm train} tfidf: Transformed text data in TF-IDF format for training. - y_{\rm train}: Series containing the labels for training.
109
110
111
112
          Output:
          - clf: Trained Random Forest classifier.
113
114
          115
          # Define a Random Forest classifier
116
          clf = RandomForestClassifier()
117
118
119
          # Define hyperparameters to tune
120
          param_grid = {
              'n_estimators': [100, 200, 300], # Number of trees in the forest 'max_depth': [None, 10, 20], # Maximum depth of the tree
121
122
123
125
          # Perform grid search with cross-validation to find the best hyperparameters
126
          grid_search = GridSearchCV(clf, param_grid, cv=5, scoring='accuracy')
          grid_search.fit(X_train_tfidf, y_train)
127
128
129
          # Print the best parameters
130
         print("Best Parameters: ", grid_search.best_params_)
131
132
          # Get the best model with the optimal hyperparameters
133
          clf = grid_search.best_estimator_
134
          # Train the final model with the entire training data
135
          clf.fit(X_train_tfidf, y_train)
136
137
138
          130
          return clf
140
     def naive bayes classifier(X train tfidf, y train):
141
142
143
144
          - X_train_tfidf: Transformed text data in TF-IDF format for training.
145
          - y_train: Series containing the labels for training.
146
          - clf: Trained Multinomial Naive Bayes classifier.
147
          Output:
148
149
150
          151
          # Define a Multinomial Naive Bayes classifier
          clf = MultinomialNB()
152
153
          # Define hyperparameters to tune
154
155
         param_grid
156
              'alpha': [0.1, 0.5, 1.0, 2.0] # Smoothing parameter (Laplace/Lidstone smoothing)
157
158
          # Perform grid search with cross-validation to find the best hyperparameters
159
          grid_search = GridSearchCV(clf, param_grid, cv=5, scoring='accuracy')
160
         grid_search.fit(X_train_tfidf, y_train)
161
162
         # Print the best parameters
print("Best Parameters: ", grid_search.best_params_)
163
164
165
          # Get the best model with the optimal hyperparameters
166
167
          clf = grid_search.best_estimator_
168
```

```
169
          # Train the final model with the entire training data
170
          clf.fit(X_train_tfidf, y_train)
171
          172
          return clf
173
174
175
      clf = logistic_regression(X_train_tfidf, y_train)
176
      clf_rf = random_forest(X_train_tfidf, y_train)
177
178
179
      clf_nb = naive_bayes_classifier(X_train_tfidf, y_train)
180
      """**4. Evaluation**""
181
182
      from sklearn.metrics import accuracy_score, classification_report
183
184
185
      def evaluate_model(clf, X_test_tfidf, y_test):
187
188
          - clf: Trained Logistic Regression model.
          - X_test_tfidf: Transformed text data in TF-IDF format for testing.
- y_test: Series containing the labels for testing.
189
190
191
192
          - None (This function will print the evaluation results.)
193
194
195
          # TODO: Evaluate the model and print the results
# Predict the labels using the trained classifier
196
197
          y_pred = clf.predict(X_test_tfidf)
198
199
200
          # Calculate the accuracy
          accuracy = accuracy_score(y_test, y_pred)
201
202
203
          print(f"Accuracy: {accuracy:.2f}")
print("Classification Report:")
204
205
206
          print(classification_report(y_test, y_pred))
207
      evaluate_model(clf, X_test_tfidf, y_test)
208
209
210
      evaluate_model(clf_rf, X_test_tfidf, y_test)
211
212
      evaluate_model(clf_nb, X_test_tfidf, y_test)
```

Task 2

(Biswas, 2020)

```
# -*- coding: utf-8 -*-
"""SWE3011_41_Task2.ipynb
2
 3
     Automatically generated by Colaboratory.
 4
 5
     Original file is located at
         https://colab.research.google.com/drive/12C7W-1nFm4VgbgysgY1E3tIqFN99MLtf
9
     # SWE3011_41 Task2
10
     **Prompt Engineering with Langchain for LLMS**
11
12
     1. Based on the given code, you need to extend (modify) to a classification task using Langchain.
13
     2. Since large language model baseline is required for the task, you may choose any model from OpenAI API and
14
     \hookrightarrow HuggingfaceHub.
15
      **OpenAI**: Only available for paid users.
16
17
18
      **HuggingfaceHUB**: Free to use with usage limits (reset hourly).
19
20
     3. Conduct experiments and document the results in the report. Here, you should consider what kind of prompt
     → design you use so please find some tutorials/resources in our homework description to obtain more

    information.

21
22
     **0. Installation**
23
^{24}
25
     pip install datasets
26
27
     pip install openai
28
29
     from getpass import getpass
30
     import os
31
     {\tt\#\ get\ a\ token:\ https://huggingface.co/docs/api-inference/quicktour\#get-your-api-token\ HUGGINGFACEHUB\_API\_TOKEN\ =\ getpass()}
32
33
34
35
     os.environ["HUGGINGFACEHUB_API_TOKEN"] = HUGGINGFACEHUB_API_TOKEN
36
37
     pip install huggingface_hub
38
     pip install langchain
39
40
     """**1. Load Dataset**
41
42
43
     Evaluation should be done using **provided test dataset**
44
45
46
     from datasets import load dataset
47
     import pandas as pd
48
     # You can use train_ds for few-shot examples
train_ds = load_dataset("glue", "sst2", split="train")
49
50
51
     # Convert the dataset to a Pandas DataFrame
52
53
     df_train = train_ds.to_pandas()
54
55
     # Display the DataFrame
56
     print(df_train.head())
57
     # Evaluation should be done using test_ds
58
     test_ds = load_dataset("csv", data_files="./test_dataset.csv")['train']
59
60
61
     # Load the dataset from CSV
     dataset_path = "/content/test_dataset.csv" # Replace with the actual path to your CSV file
62
63
     df_test = pd.read_csv(dataset_path)
64
     # Display the DataFrame
65
     print(df_test.head())
66
67
68
     import random
69
     import numpy as np
70
71
     def gen_data(df, flag=True, samples=27):
72
         examples = []
73
74
         n = len(df) if flag else samples
75
         # Count occurrences of each category
category_counts = df["label"].value_counts()
76
77
78
         for i in range(0, n):
```

```
# Choose category proportionally based on occurrences
81
              category = np.random.choice(category_counts.index.tolist(), p=(category_counts /
              \hookrightarrow category_counts.sum()).tolist())
82
 83
               # Get a random row with the chosen category
              category_df = df[df["label"] == category]
random_row = category_df.sample()
 85
 86
              obj = {
87
 88
                   "messages": random_row["sentence"].iloc[0],
                   "categories": random_row["label"].iloc[0]
89
 90
91
              examples.append(obj)
92
93
          random.shuffle(examples)
94
95
          return examples
96
97
      """**2. Preparing Prompt**""
98
     from langchain import PromptTemplate
99
100
      prompt_template = '''
101
102
      Messages: {messages}
103
      Categories: {categories}
104
105
     example_prompt = PromptTemplate(
    input_variables=["messages", "categories"],
106
107
          template=prompt_template)
108
109
110
     import re
111
112
      def trim(s):
          if "positive" in s.lower():
113
              return 1
114
          elif "negative" in s.lower():
116
             return 0
117
              return -1 # or any default value for the case where no match is found
118
119
      from sklearn.metrics import accuracy_score
120
      from tqdm import tqdm
122
      import time
123
124
      def eval(extraction_chain, example):
          preds = []
actual = []
125
126
127
          for x in tqdm(example):
129
              p = extraction_chain.predict(messages=x["messages"])
              p = trim(p)
130
              a = x["categories"]
131
              preds.append(p)
132
133
              actual.append(a)
134
135
          accuracy = accuracy_score(actual, preds)
136
          return preds, actual, accuracy
137
      from langchain.llms import HuggingFaceHub
138
      from langchain.chains import LLMChain
139
140
141
      def initialize_llm(model_name, api_key=None):
142
143
          Initialize the model using the langchain library.
144
145
          11m = HuggingFaceHub(repo_id=model_name, model_kwargs={"temperature": 0.5, "max_length": 64})
146
147
148
      """**Few-Shot**"""
149
150
     from langchain import FewShotPromptTemplate
151
152
153
      examples = gen_data(df_train, False)
154
     print (examples)
155
      # Create the FewShotPromptTemplate
156
      final_prompt = FewShotPromptTemplate(
157
          examples=examples,
158
159
          example_prompt=example_prompt,
160
          suffix="Messages: {messages} \nCategories: ",
161
          input_variables=["messages"],
          prefix="""You are a sentiment classifier. Classify the given text input as "positive" or "negative" based
162
          \hookrightarrow on its tone.
163
                   Please read each sentence carefully and determine its sentiment based on common definitions.
164
```

```
165
                   <positive>: A positive sentiment indicates a favorable or optimistic tone. It expresses approval,

→ satisfaction, or positive emotions.

166
                   <negative>: A negative sentiment indicates an unfavorable or pessimistic tone. It expresses
167

→ disapproval, dissatisfaction, or negative emotions.

168
169
                   Please read each sentence and classify it as either "positive" or "negative" based on the provided

→ definitions.

170
      )
171
172
173
      # Initialize the language model
      model_name = "google/flan-t5-xxl"
174
      llm = initialize_llm(model_name)
175
      val = gen_data(df_test)
176
177
      # Create the LLMChain
178
      extraction_chain = LLMChain(llm=llm, prompt=final_prompt, output_key="categories")
180
181
      # Evaluate and print results without a new line after a comma
      preds, actual, acc = eval(extraction_chain, val)
print("\nPredicted Categories:", preds)
182
183
184
185
      print("Actual Categories:", actual)
186
187
      print(f"Accuracy: {acc * 100:.2f}%")
188
      from langchain import FewShotPromptTemplate
189
190
      # Initialize the language model
191
      model_name = "google/flan-t5-xxl"
192
193
      llm = initialize_llm(model_name)
194
195
      # List to store accuracies
      accuracies = []
196
197
      # Perform 10 iterations
198
199
      for iteration in range(1, 4):
200
          print(f"\nIteration {iteration}:")
201
           # Generate new examples for each iteration
202
203
          examples = gen_data(df_train, flag = False, samples=27)
204
205
           # Define your FewShotPromptTemplate
206
          final_prompt = FewShotPromptTemplate(
207
               examples=examples,
208
               example prompt=example prompt,
               suffix="Messages: {messages} \nCategories: ",
input_variables=["messages"],
209
210
211
               prefix="""You are a sentiment classifier. Classify the given text input as "positive" or "negative"

    ⇒ based on its tone.

212
                        Please read each sentence carefully and determine its sentiment based on common definitions.
213
                        <positive>: A positive sentiment indicates a favorable or optimistic tone. It expresses
214
                        → approval, satisfaction, or positive emotions.
215
216
                        <negative>: A negative sentiment indicates an unfavorable or pessimistic tone. It expresses
                        \,\hookrightarrow\, disapproval, dissatisfaction, or negative emotions.
217
                       Please read each sentence and classify it as either "positive" or "negative" based on the
218
                        \hookrightarrow provided definitions.
               ....
219
220
221
          # Create the LLMChain with the updated prompt
extraction_chain = LLMChain(llm=llm, prompt=final_prompt, output_key="categories")
222
223
224
225
           # Evaluate and print results without a new line after a comma
226
          preds, actual, acc = eval(extraction_chain, val)
          accuracies.append(acc) # Store accuracy in the list
print("Predicted Categories:", preds)
227
228
229
          print("Actual Categories:", actual)
230
231
232
          print(f"Accuracy: {acc * 100:.2f}%")
233
      # Print all accuracies at the end
print("\nAll Accuracies:", accuracies)
234
235
236
      """**Zero-Shot**"""
237
238
239
      from langchain import PromptTemplate
240
241
      # Create the zero-shot prompt template
      zero_shot_template = PromptTemplate(
242
          input_variables=["messages"],
243
```

```
244
            template="""You are a sentiment classifier. Classify the given text input as "positive" or "negative"
            \hookrightarrow based on its tone. Please read each sentence carefully and determine its sentiment based on common definitions.
245
246
                            <positive>: A positive sentiment indicates a favorable or optimistic tone. It expresses
247
                            \stackrel{\smile}{\hookrightarrow} approval, satisfaction, or positive emotions.
248
                           <negative>: A negative sentiment indicates an unfavorable or pessimistic tone. It expresses \leftrightarrow disapproval, dissatisfaction, or negative emotions.
249
250
                            Please read each sentence and classify it as either "positive" or "negative" based on the
251
                            \hookrightarrow provided definitions.
            ....
252
253
254
      # Initialize the language model
model_name = "google/flan-t5-xx1"
llm = initialize_llm(model_name)
255
256
257
258
       val = gen_data(df_test)
^{259}
       # Create the LLMChain
260
       extraction_chain = LLMChain(llm=llm, prompt=zero_shot_template, output_key="categories")
261
262
       # Evaluate and print results
preds, actual, acc = eval(extraction_chain, val)
print("\nPredicted Categories:", preds)
263
264
265
266
       print("Actual Categories:", actual)
267
268
       print(f"Accuracy: {acc * 100:.2f}%", end=" ")
269
```

A Appendix

Figure 7: Result of prediction using zero-shot prompt

```
100%| 100/100 [01:07<00:00, 1.49it/s]
Predicted Categories: [1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, -1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0
```

Figure 8: Result of prediction with the first execution of few-shot prompt

Figure 9: Result of prediction with subsequent executions of few-shot prompt