

Glassdoor Feedback Analysis Project

Jessica Jazmin Castillo Rios

Table of Contents

| | | |
|----------|-----------------------------|-----------|
| 1 | Project Overview | 2 |
| 2 | Dataset and Scraping | 2 |
| 3 | Preprocessing | 5 |
| 4 | Model Training | 6 |
| 5 | Evaluation | 7 |
| 6 | Sentiment Analysis | 8 |
| 7 | MLOps with MLflow | 9 |
| 8 | Execution Guide | 10 |
| 9 | Architecture Diagram | 11 |

1 Project Overview

This project aims to extract, analyze, and classify employee reviews from Glassdoor. The reviews are classified into positive (Pro) and negative (Con) categories using natural language processing and a Logistic Regression model.

2 Dataset and Scraping

Source

Reviews are scraped directly from [glassdoor.com](https://www.glassdoor.com) using **Selenium** and **BeautifulSoup**. Up to 100 reviews per company can be extracted.

```
1
2 # Standard and third-party libraries needed for scraping, analysis, and data manipulation
3 import time
4 import json
5 import pandas as pd
6 from bs4 import BeautifulSoup
7 from langdetect import detect
8 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
9
10 # Selenium and WebDriver setup
11 from selenium import webdriver
12 from selenium.webdriver.common.by import By
13 from selenium.webdriver.chrome.options import Options
14 from selenium.webdriver.chrome.service import Service
15 from webdriver_manager.chrome import ChromeDriverManager
16
17 # Initialize the VADER sentiment analyzer
18 analyzer = SentimentIntensityAnalyzer()
19
20 # WebDriver initialization
21 def init_driver():
22     options = Options()
23
24     # NOTE: Headless is disabled to allow manual CAPTCHA solving if needed
25     # options.add_argument("--headless=new")
26
27     options.add_argument("--disable-gpu")
28     options.add_argument("--window-size=1920,1080")
29
30     # Set a custom user-agent to simulate a real browser
31     options.add_argument(
32         "user-agent=Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 "
33         "KHTML, like Gecko Chrome/123.0.0.0 Safari/537.36"
34     )
35
36     # Options to reduce Selenium detection
37     options.add_experimental_option("excludeSwitches", ["enable-automation"])
38     options.add_experimental_option("useAutomationExtension", False)
39
40     # Initialize Chrome WebDriver with the above options
41     driver = webdriver.Chrome(service=Service(ChromeDriverManager().install()), options=options)
42
43     # Hide the "webdriver" flag in JavaScript
44     driver.execute_script("Object.defineProperty(navigator, 'webdriver', {get: () => undefined})")
45
46     return driver
47
48 # Load cookies from a file to avoid manual login
49 def load_cookies(driver, cookie_file="glassdoor_cookies.json"):
50     print("Loading saved cookies...")
51     driver.get("https://www.glassdoor.com/")
52     time.sleep(3)
53
54     with open(cookie_file, "r") as f:
55         cookies = json.load(f)
56         for cookie in cookies:
57             if "sameSite" in cookie:
```

```

58         cookie["sameSite"] = 'Strict' # Selenium compatibility adjustment
59         driver.add_cookie(cookie)
60
61     # Refresh the page with applied cookies
62     driver.refresh()
63     time.sleep(4)
64
65     # Classify sentiment using VADER
66     def classify_sentiment(text):
67         score = analyzer.polarity_scores(text)["compound"]
68         if score > 0.02:
69             return "Pro"
70         elif score < -0.02:
71             return "Con"
72         else:
73             return "Neutral"
74
75     # Scrape company reviews from Glassdoor
76     def scrape_company_reviews(company_slug, max_reviews=100):
77         print(f"Scraping reviews for {company_slug}...")
78         driver = init_driver()
79         load_cookies(driver)
80
81         reviews = []
82         page = 1
83
84         while len(reviews) < max_reviews:
85             # Build the URL for the current page
86             url = f"https://www.glassdoor.com/Reviews/{company_slug.replace('.htm', f'_P{page}.htm')}"
87             driver.get(url)
88             time.sleep(5)
89
90             soup = BeautifulSoup(driver.page_source, "html.parser")
91
92             # Find PROS and CONS review sections
93             pros = soup.find_all("span", attrs={"data-test": "review-text-PROS"})
94             cons = soup.find_all("span", attrs={"data-test": "review-text-CONS"})
95
96             if not pros and not cons:
97                 print(f"Page {page} does not contain visible reviews. Stopping scraping.")
98                 break
99
100            # Pair PROS and CONS by index
101            total = min(len(pros), len(cons))
102            for i in range(total):
103                if len(reviews) >= max_reviews:
104                    break
105                try:
106                    pros_text = pros[i].get_text(strip=True)
107                    lang_pros = detect(pros_text)
108                    reviews.append({
109                        "company": company_slug.split("-Reviews")[0],
110                        "review": pros_text,
111                        "language": lang_pros,
112                        "type": "Pro"
113                    })
114                except:
115                    pass
116
117            if len(reviews) >= max_reviews:
118                break
119            try:
120                cons_text = cons[i].get_text(strip=True)
121                lang_cons = detect(cons_text)
122                reviews.append({
123                    "company": company_slug.split("-Reviews")[0],
124                    "review": cons_text,
125                    "language": lang_cons,
126                    "type": "Con"
127                })
128            except:

```

```

129         pass
130
131         print(f"Page {page}: {total * 2} reviews processed (Total: {len(reviews)})")
132         page += 1
133
134     driver.quit()
135     print(f"Total reviews extracted for {company_slug}: {len(reviews)}")
136     return reviews
137
138 # Scrape multiple companies and save results to CSV
139 def scrape_multiple(company_slugs, output_csv="data/glassdoor_reviews.csv"):
140     all_reviews = []
141     for slug in company_slugs:
142         all_reviews.extend(scrape_company_reviews(slug))
143
144     df = pd.DataFrame(all_reviews)
145     df.to_csv(output_csv, index=False)
146     print(f"Reviews saved to {output_csv}")

```

3 Preprocessing

Preprocessing includes:

- Lowercasing text, removing punctuation and digits.
- Language detection using `langdetect`.
- Removing stopwords based on the detected language (English or Spanish).

```
1
2 import pandas as pd
3 import re
4 import string
5 from langdetect import detect
6 import nltk
7 from nltk.corpus import stopwords
8
9 # Check if stopwords are already downloaded; if not, download them
10 try:
11     nltk.data.find("corpora/stopwords")
12 except LookupError:
13     nltk.download("stopwords")
14
15 # Load stopword lists for English and Spanish
16 stopwords_en = set(stopwords.words("english"))
17 stopwords_es = set(stopwords.words("spanish"))
18
19 # Function to clean the text:
20 # converts to lowercase, removes punctuation and digits
21 def clean_text(text):
22     text = text.lower() # Lowercase
23     text = re.sub(f"[{re.escape(string.punctuation)}]", "", text) # Remove punctuation
24     text = re.sub(r"\d+", "", text) # Remove digits
25     return text.strip() # Trim whitespace
26
27 # Function to detect the language of a text
28 def detect_language(text):
29     try:
30         return detect(text)
31     except:
32         return "unknown" # If detection fails, return "unknown"
33
34 # Function to remove stopwords based on the detected language
35 def remove_stopwords(text, lang):
36     tokens = text.split() # Split the text into simple tokens
37     stops = stopwords_en if lang == "en" else stopwords_es # Choose the appropriate stopword set
38     filtered = [word for word in tokens if word not in stops] # Remove stopwords
39     return " ".join(filtered) # Join the filtered words back into a string
40
41 # Function to apply all preprocessing steps to a DataFrame
42 def preprocess_dataframe(df):
43     print("[INFO] Preprocessing reviews...")
44
45     # Basic text cleaning
46     df["cleaned"] = df["review"].apply(clean_text)
47
48     # Language detection for each review
49     df["language"] = df["review"].apply(detect_language)
50
51     # Remove stopwords based on detected language
52     df["processed"] = df.apply(lambda row: remove_stopwords(row["cleaned"], row["language"]), axis=1)
53
54     return df
```

4 Model Training

The model used is a **Logistic Regression** pipeline with TF-IDF vectorization.

```
1
2 # Import necessary libraries for text vectorization, model training, and evaluation
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.pipeline import Pipeline
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import classification_report
8
9 # Function to train a text classification model using Logistic Regression
10 # df: DataFrame with preprocessed data
11 # text_column: column containing the processed text
12 # label_column: column containing the labels ("Pro", "Con")
13 def train_classifier(df, text_column="processed", label_column="type"):
14     print("Training classification model...")
15
16     # Filter the DataFrame to keep only "Pro" and "Con" classes
17     df = df[df[label_column].isin(["Pro", "Con"])]
18
19     # Split the data into training and testing sets (80%-20%)
20     # Ensures at least 2 examples in the test set
21     X_train, X_test, y_train, y_test = train_test_split(
22         df[text_column],
23         df[label_column],
24         test_size=max(2, int(len(df) * 0.2)),
25         random_state=42
26     )
27
28     # Show class distribution in the test set
29     print("Class distribution in test set:", y_test.value_counts())
30
31     # Define a pipeline: TF-IDF vectorization followed by Logistic Regression classifier
32     pipeline = Pipeline([
33         ("tfidf", TfidfVectorizer()),
34         ("clf", LogisticRegression(max_iter=500, class_weight="balanced")) # max_iter=500 to
35         ensure convergence
36     ])
37
38     # Train the pipeline with the training data
39     pipeline.fit(X_train, y_train)
40
41     # Make predictions on the test set
42     y_pred = pipeline.predict(X_test)
43
44     # Generate a classification report with per-class metrics
45     report = classification_report(y_test, y_pred, output_dict=True)
46
47     print("Classification completed.")
48
49     # Return the trained model, the report, and the evaluation data
50     return pipeline, report, X_test, y_test, y_pred
```

5 Evaluation

Evaluation metrics include:

- Precision
- Recall
- F1-Score
- Confusion Matrix

```
1
2 # Import necessary libraries for plotting and evaluating classification models
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import classification_report, confusion_matrix, ConfusionMatrixDisplay
5
6 # Function to evaluate a classification model
7 # Takes the true values (y_test), predicted values (y_pred),
8 # and the class labels to be shown in the confusion matrix
9 def evaluate_model(y_test, y_pred, labels=["Pro", "Con"]):
10     # Print a report with metrics like precision, recall, and F1-score
11     print("Model results:")
12     print(classification_report(y_test, y_pred))
13
14     # Compute the confusion matrix with the specified labels
15     cm = confusion_matrix(y_test, y_pred, labels=labels)
16
17     # Create a display object for visualizing the confusion matrix
18     disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=labels)
19
20     # Plot the confusion matrix using a blue color map
21     disp.plot(cmap=plt.cm.Blues)
22     plt.title("Confusion Matrix")
23
24     # Show the plot
25     plt.show()
```

6 Sentiment Analysis

VADER sentiment analysis is applied to reviews.

```
1
2 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
3
4 # Initialize VADER sentiment analyzer
5 analyzer = SentimentIntensityAnalyzer()
6
7 def analyze_sentiment(text, lang):
8     # Applies sentiment analysis using VADER, regardless of language.
9     # Returns the sentiment label ('POS', 'NEG', 'NEU') and the compound score.
10    score = analyzer.polarity_scores(text)
11    label = "POS" if score["compound"] > 0.05 else "NEG" if score["compound"] < -0.05 else "NEU"
12    return label, score["compound"]
```


7 MLOps with MLflow

MLflow is used to log:

- Model type and configuration
- Metrics per label (Pro/Con)
- Classification report as a JSON artifact

```
1
2 # Import necessary libraries for working with MLflow, file handling, and JSON serialization
3 import mlflow
4 import mlflow.sklearn
5 import json
6 import os
7
8 # Function to log a trained model and its evaluation metrics in MLflow
9 # model: trained sklearn model
10 # classification_report_dict: dictionary containing evaluation metrics (e.g., output of
11 #   classification_report(output_dict=True))
12 # run_name: name of the experiment run in MLflow
13 def log_with_mlflow(model, classification_report_dict, run_name="glassdoor_model_run"):
14     print("Registering experiment in MLflow...")
15
16     # Define or create an MLflow experiment to store the run data
17     mlflow.set_experiment("Glassdoor_Reviews_Analysis")
18
19     # Start a new MLflow run with the specified name
20     with mlflow.start_run(run_name=run_name):
21         # Log the trained model in MLflow
22         mlflow.sklearn.log_model(model, "model")
23
24         # Log model type as a parameter
25         mlflow.log_param("model_type", "LogisticRegression")
26
27         # Loop through evaluation metrics for each class ("Pro", "Con") and log the values
28         for label in ["Pro", "Con"]:
29             for metric in ["precision", "recall", "f1-score"]:
30                 value = classification_report_dict.get(label, {}).get(metric)
31                 if value is not None:
32                     mlflow.log_metric(f"{label}_{metric}", value)
33
34         # Create a folder to store artifacts if it doesn't exist
35         os.makedirs("mlruns_artifacts", exist_ok=True)
36
37         # Save the classification report as a JSON file
38         report_path = "mlruns_artifacts/classification_report.json"
39         with open(report_path, "w") as f:
40             json.dump(classification_report_dict, f, indent=4)
41
42         # Log the JSON report as an artifact in MLflow
43         mlflow.log_artifact(report_path)
44
45     print("MLflow tracking completed.")
```

MLflow UI

To view the results, start the MLflow UI:

```
1 mlflow ui
```

Then, open your browser and navigate to <http://localhost:5000>.

8 Execution Guide

Follow these steps to execute the project:

1. Set Up the Environment:

- Install the required dependencies:

```
1 pip install -r requirements.txt
2
```

2. Run the Script:

- Execute the main script:

```
1 python src/main.py
2
```

3. View Results:

- Check the console for metrics and visualizations.
- Use the MLflow UI to explore logged experiments and artifacts.

4. Upload Results to GitHub:

- Use the `upload_results.py` script to push results to a GitHub repository:

```
1 python upload_results.py
2
```

9 Architecture Diagram

The following diagram illustrates the architecture of the Glassdoor Feedback Analysis project. It shows the complete pipeline, starting from data extraction via web scraping, followed by preprocessing, model training and evaluation, and finally, experiment tracking using MLflow. This modular design ensures that each component is clearly separated, reusable, and easy to maintain or scale.

