

# Sentiment Analysis Using Machine Learning

Davide Giuseppe Griffon

## Abstract

This document serves as the report for the third task in the "Natural Language Processing" course completed by student Davide Giuseppe Griffon at Vilnius University as part of the Master's program in Data Science.

## Contents

<b>1</b>	<b>Dataset</b>	<b>2</b>
<b>2</b>	<b>Preprocessing</b>	<b>2</b>
<b>3</b>	<b>Word Cloud</b>	<b>3</b>
<b>4</b>	<b>Vectorization</b>	<b>5</b>
<b>5</b>	<b>Machine Learning Models</b>	<b>6</b>
<b>6</b>	<b>Results and Analysis</b>	<b>6</b>
<b>7</b>	<b>Conclusions</b>	<b>7</b>
<b>8</b>	<b>Appendix - Code</b>	<b>8</b>

# 1 Dataset

For this sentiment analysis project I used the IMDB Movie Reviews dataset. The dataset was loaded using the Python's `datasets` library, which provides a convenient interface for accessing various NLP datasets. The dataset contains binary sentiment labels: positive (1) and negative (0), representing the overall sentiment of each movie review.

The dataset consists of movie reviews from the Internet Movie Database (IMDB) and is structured as follows:

- **Total Size:** 50,000 movie reviews
- **Split Distribution:**
  - Training set: 25,000 reviews
  - Testing set: 25,000 reviews
- **Class Distribution:** The dataset is perfectly balanced, with:
  - 50% negative reviews (label 0)
  - 50% positive reviews (label 1)
- **Review Length:** The average review length is 1,325 characters

The data was loaded and processed using a custom `load_imdb_dataset` function, which returns two separate pandas DataFrames: one for training and one for testing.

The balanced nature of the dataset is advantageous for training machine learning models because it eliminates the need for class weight adjustments or other imbalance-handling techniques that might otherwise be necessary. This is one of the main reasons I chose to use this particular dataset for this sentiment analysis task.

## 2 Preprocessing

### Data Loading and Sampling

The initial step involves loading the IMDB dataset using the previously mentioned `load_imdb_dataset` function. For this project, I used a subset of 10,000 reviews (5,000 for training and 5,000 for testing) to optimize computational efficiency while maintaining result quality. This sampling approach allowed for faster preprocessing and training phases without compromising the effectiveness of the sentiment analysis.

### Text Cleaning

The `clean_dataset` function implements several text cleaning steps:

- Removal of missing values and empty cells to ensure data integrity.
- Elimination of HTML tags (such as `<br>`) found in the raw text through regex patterns.
- Removal of punctuation marks, retaining only alphanumeric characters through regex patterns.

- Conversion of all text to lowercase.
- Removal of stopwords using the *spaCy* library.

## Intermediate Data Storage

To optimize the workflow and enable intermediate analysis, the preprocessing pipeline includes data persistence steps:

1. The initially cleaned data (pre-lemmatization) is saved using `save_raw_datasets_to_local` function, that stores the cleaned data in two separate CSV files: one for training and one for testing.
2. This raw preprocessed data can be retrieved using `load_local_raw_datasets` function.

This intermediate storage allows for analysis of the data before and after lemmatization.

## Lemmatization

The final preprocessing step involves lemmatization and saves the fully preprocessed data to two new CSV files using the `save_lemma_datasets_to_local` function. To load the lemmatized data, the `load_local_lemma_datasets` function can be used. The steps involved in lemmatization are as follows:

- Loads the raw preprocessed data using the `load_local_raw_datasets` function.
- Applies *spaCy*'s lemmatization to reduce words to their base form.
- Saves the fully preprocessed data to two separate CSV files: one for training and one for testing.

The separation of raw cleaning and lemmatization steps enables comparative analysis of their impact on the Word Cloud, as we'll see in the next section.

# 3 Word Cloud

In order to visualize the most frequent words in the movie reviews and understand their distribution, I generated word clouds using the *WordCloud* library. Word clouds provide an intuitive visualization where the size of each word is proportional to its frequency in the text corpus, making it easier to identify dominant themes and common expressions in the reviews.

## Implementation

The word cloud generation was implemented using a custom `generate_word_cloud` function that processes the text data and creates visualizations. This function was applied to both the raw preprocessed dataset and the lemmatized dataset to observe the effects of lemmatization on word frequencies:

- For raw data: Applied after basic cleaning but before lemmatization
- For lemmatized data: Applied after all preprocessing steps, including lemmatization

## Results Analysis



Figure 1: Word Cloud of Raw Preprocessed Reviews



Figure 2: Word Cloud of Lemmatized Reviews

Analysis of word frequencies reveals several interesting patterns:

### General Observations

- **Domain-Specific Terms:** "movie" and "film" are consistently the most frequent words across both sentiment classes and preprocessing methods, indicating their role as domain identifiers rather than sentiment indicators.
- **Common Base Words:** Terms like "like", "time", "people", and "characters" appear frequently in both positive and negative reviews, suggesting they are neutral in sentiment.

### Sentiment-Specific Patterns (after lemmatization)

- **Negative Reviews:**
  - "bad" appears significantly more frequently (4,308 vs 2939 in positive reviews)
  - Higher frequency of "not" (8,050 vs 5,419 in positive reviews)
  - Focus on technical aspects: "plot", "acting", "scene"
- **Positive Reviews:**
  - Distinctive positive terms: "great" (2,775), "love" (2,259)
  - "good" appears more frequently (4,347 vs 3,500 in negative reviews)
  - More emphasis on emotional terms: "love", "great", "best"

## Impact of Lemmatization

The lemmatization process had several notable effects:

- **Word Consolidation - examples:**
  - "movie" and "movies" were consolidated, increasing the frequency from 9,493 to 11,188 in negative reviews
  - "character" and "characters" were combined
  - "watch", "watching", and "watched" were merged into "watch"
- **Negation Handling:** The contraction "nt" was properly transformed to "not", making it more prominent in the frequency counts
- **Verb Forms:** Various forms of verbs were consolidated (e.g., "see"/"seen", "act"/"acting"), providing clearer frequency patterns

These findings suggest that while some words are strong indicators of sentiment ("bad", "great", "love"), many frequent terms are neutral and context-dependent. The lemmatization process helped clarify word usage patterns by consolidating different forms of the same word, potentially improving the accuracy of subsequent sentiment analysis steps.

## 4 Vectorization

To convert the text data into a numerical format suitable for machine learning model training, I tried three distinct vectorization techniques: Word2Vec, Bag-of-Words (BoW), and Term Frequency-Inverse Document Frequency (TF-IDF).

### Implementation Architecture

I developed a dedicated `vectorization` module to handle all text vectorization processes. The module's core function `get_vector_datasets` provides a unified interface for accessing different vectorization methods:

```
1 def get_vector_datasets(vectorization_type="tfidf"):  
2     vectorization_function = {  
3         "count": lambda: get_vector_datasets_bow("count"),  
4         "tfidf": lambda: get_vector_datasets_bow("tfidf"),  
5         "word2vec": get_vector_datasets_word2vec,  
6     }  
7     return vectorization_function[vectorization_type]()
```

### Vectorization Methods

The implementation leverages different libraries for each method:

- **Word2Vec:** Implemented using the *gensim* library
- **Bag-of-Words (BoW):** Implemented using *scikit-learn*'s `CountVectorizer`
- **TF-IDF:** Implemented using *scikit-learn*'s `TfidfVectorizer`

## Data Flow

Each vectorization method follows a consistent workflow:

1. Receives preprocessed training and testing datasets as separate pandas DataFrames
2. Fits the vectorizer on the training data only to prevent data leakage
3. Transforms both training and testing sets using the fitted vectorizer
4. Returns vectorized versions of both datasets maintaining the original split

In this way I could easily experiment with different vectorization techniques while maintaining consistent interfaces for the subsequent machine learning models. The default TF-IDF method was chosen based on preliminary experiments showing better performance with the movie review corpus.

## 5 Machine Learning Models

Since in task 2 I already used a deep learning model, for this task I chose to try classical machine learning models. I experimented with four different models: Logistic Regression, Random Forest, Gradient Boosting and Support Vector Machine. I run all of them and compared their performance using the same vectorized data.

### Implementation Details

The code resides in the `models/classical_models.py` file. All models were implemented using the *scikit-learn* library. The function `classical_models_comparison` was used to run and compare models with the following configurations:

```
1 models = {  
2     "Logistic Regression": LogisticRegression(  
3         random_state=RANDOM_STATE, max_iter=1000  
4     ),  
5     "Random Forest": RandomForestClassifier(  
6         n_estimators=100, random_state=RANDOM_STATE  
7     ),  
8     "Gradient Boosting": GradientBoostingClassifier(  
9         random_state=RANDOM_STATE  
10    ),  
11    "SVM": SVC(random_state=RANDOM_STATE),  
12 }
```

Each model was trained and evaluated using the same vectorized dataset to ensure a fair comparison of their performance.

## 6 Results and Analysis

I evaluated all models with each vectorization method, conducting a total of  $4 \times 3 = 12$  experiments. The TF-IDF vectorization method demonstrated the best overall performance, and its results are presented below for each model.

In table 1, I present the performance metrics for each model using the TF-IDF vectorization method. The metrics include training accuracy, test accuracy, precision, and F1-score.

Table 1: Model Performance Comparison

Model	Training Accuracy	Test Accuracy	Precision	F1-Score
Logistic Regression	0.8902	0.8308	0.83	0.83
Random Forest	1.0000	0.8187	0.82	0.82
Gradient Boosting	0.8353	0.7989	0.80	0.80
SVM	0.9751	0.8400	0.84	0.84

## Results Summary

Among all tested models, SVM achieved the highest test accuracy (0.84), followed closely by Logistic Regression (0.83). Random Forest showed signs of overfitting with perfect training accuracy (1.0) but lower test accuracy (0.82). Gradient Boosting demonstrated the lowest performance with a test accuracy of 0.80. All models maintained balanced performance across classes, as evidenced by similar precision, recall, and F1-scores for both positive and negative reviews.

## 7 Conclusions

The sentiment analysis task on IMDB movie reviews yielded several key findings:

- TF-IDF vectorization consistently outperformed other vectorization methods
- SVM achieved the best performance with 84% test accuracy
- All models maintained balanced performance between positive and negative classes

The results demonstrate that classical machine learning models can effectively perform sentiment analysis on movie reviews, despite modern deep learning models would probably achieve better results.

## 8 Appendix - Code

All the code is available in the GitHub repository <https://github.com/Griffosx/nlp> under the src/task\_3 folder. For completeness, I include here the code for the main files used in this project.

File preprocessing.py

```
1 from collections import Counter
2 import pandas as pd
3 from wordcloud import WordCloud
4 import matplotlib.pyplot as plt
5 from datasets import load_dataset
6 from task_3.constants import (
7     nlp,
8     POSITIVE_LABEL,
9     TRAIN_RAW_DATA_PATH,
10    TEST_RAW_DATA_PATH,
11    TRAIN_LEMMA_DATA_PATH,
12    TEST_LEMMA_DATA_PATH,
13 )
14
15
16 def load_imdb_dataset(
17     num_samples=None, print_stats=False
18 ) -> tuple[pd.DataFrame, pd.DataFrame]:
19     """
20     Load the IMDB movie reviews dataset for sentiment analysis.
21     """
22     # Load the dataset
23     dataset = load_dataset("imdb")
24
25     # Convert to pandas DataFrames
26     train_df = pd.DataFrame(dataset["train"])
27     test_df = pd.DataFrame(dataset["test"])
28
29     # Sample if specified
30     if num_samples:
31         train_df = train_df.sample(min(num_samples, len(train_df)),
32                                     random_state=42)
33         test_df = test_df.sample(min(num_samples, len(test_df)),
34                                  random_state=42)
35
36     if print_stats:
37         # Add some basic statistics
38         print(f"Dataset Statistics:")
39         print(f"Training samples: {len(train_df)}")
40         print(f"Testing samples: {len(test_df)}")
41         print(f"\nClass distribution in training:")
42         print(train_df["label"].value_counts(normalize=True))
43         # Calculate average review length
44         train_df["review_length"] = train_df["text"].str.len()
45         print(
46             f"\nAverage review length: {train_df['review_length'].mean():.0f} characters"
47         )
48
49     return train_df, test_df
```



```

48
49
50 def clean_dataset(dataset: pd.DataFrame) -> pd.DataFrame:
51     # Create a copy of the dataset
52     cleaned_dataset = dataset.copy()
53
54     # Remove missing values and empty cells
55     cleaned_dataset = cleaned_dataset.dropna(subset=["text"])
56     cleaned_dataset = cleaned_dataset[cleaned_dataset["text"].str.strip
57         ().astype(bool)]
58
59     # Remove HTML tags
60     cleaned_dataset["text"] = cleaned_dataset["text"].str.replace(
61         r"<[^>*>", "", regex=True
62     )
63
64     # Remove punctuation, leave only alphanumeric characters using regex
65     cleaned_dataset["text"] = cleaned_dataset["text"].str.replace(
66         r"[^\w\s]", "", regex=True
67     )
68
69     # Convert to lowercase
70     cleaned_dataset["text"] = cleaned_dataset["text"].str.lower()
71
72     # Remove stopwords using spaCy
73     def clean_text(text):
74         doc = nlp(text)
75         # Keep only non-stopword tokens and strip spaces
76         return " ".join(token.text.strip() for token in doc if not
77             token.is_stop)
78
79     cleaned_dataset["text"] = cleaned_dataset["text"].apply(clean_text)
80
81     return cleaned_dataset
82
83 def load_and_clean_imdb_dataset(
84     num_samples=None, print_stats=False
85 ) -> tuple[pd.DataFrame, pd.DataFrame]:
86     train_data, test_data = load_imdb_dataset(num_samples, print_stats)
87     return clean_dataset(train_data), clean_dataset(test_data)
88
89 def save_raw_datasets_to_local(num_samples=None, print_stats=False):
90     train_data, test_data = load_and_clean_imdb_dataset(num_samples,
91         print_stats)
92     train_data.to_csv(TRAIN_RAW_DATA_PATH, index=False)
93     test_data.to_csv(TEST_RAW_DATA_PATH, index=False)
94
95 def load_local_raw_datasets() -> tuple[pd.DataFrame, pd.DataFrame]:
96     train_data = pd.read_csv(TRAIN_RAW_DATA_PATH)
97     test_data = pd.read_csv(TEST_RAW_DATA_PATH)
98     return train_data, test_data
99
100
101 def save_lemma_datasets_to_local():
102     train_data, test_data = load_local_raw_datasets()

```

```

103 train_data["text"] = train_data["text"].apply(
104     lambda x: " ".join([token.lemma_ for token in nlp(x)])
105 )
106 test_data["text"] = test_data["text"].apply(
107     lambda x: " ".join([token.lemma_ for token in nlp(x)])
108 )
109 train_data.to_csv(TRAIN_LEMMA_DATA_PATH, index=False)
110 test_data.to_csv(TEST_LEMMA_DATA_PATH, index=False)
111
112
113 def load_local_lemma_datasets() -> tuple[pd.DataFrame, pd.DataFrame]:
114     train_data = pd.read_csv(TRAIN_LEMMA_DATA_PATH)
115     test_data = pd.read_csv(TEST_LEMMA_DATA_PATH)
116     return train_data, test_data
117
118
119 def generate_wordcloud(lemmatisation=True):
120     """
121     For each label in the dataset, generate and plot a wordcloud and
122     print the top 20 most frequent words.
123     """
124     if lemmatisation:
125         dataset, _ = load_local_lemma_datasets()
126     else:
127         dataset, _ = load_local_raw_datasets()
128
129     # Get unique labels from the dataset
130     unique_labels = dataset["label"].unique()
131
132     # Create a figure with subplots for each label
133     fig, axes = plt.subplots(1, len(unique_labels), figsize=(15, 5))
134
135     for idx, label in enumerate(unique_labels):
136         # Filter text for current label
137         texts = dataset[dataset["label"] == label]["text"]
138
139         text = " ".join(texts)
140
141         # Generate word frequency distribution
142         words = text.split()
143         word_freq = Counter(words)
144
145         # Get top 20 words and their frequencies
146         top_20 = word_freq.most_common(20)
147
148         # Print top 20 words for current label
149         print(
150             f"\nTop 20 most frequent words for {'positive' if label == POSITIVE_LABEL else 'negative'} reviews:"
151         )
152         print("{:<15} {:<10}".format("Word", "Frequency"))
153         print("-" * 25)
154         for word, freq in top_20:
155             print("{:<15} {:<10}".format(word, freq))
156
157         # Generate wordcloud
158         wordcloud = WordCloud(
159             width=800, height=400, background_color="white", stopwords=

```

```

159         set()
160     ).generate(text)
161
162     # Plot wordcloud
163     if len(unique_labels) > 1:
164         axes[idx].imshow(wordcloud)
165         axes[idx].axis("off")
166         axes[idx].set_title(
167             f"Label: {'positive' if label == POSITIVE_LABEL else 'negative'}"
168         )
169     else:
170         axes.imshow(wordcloud)
171         axes.axis("off")
172         axes.set_title(
173             f"Label: {'positive' if label == POSITIVE_LABEL else 'negative'}"
174         )
175
176 plt.tight_layout()
plt.show()

```

#### File vectorization/bow.py

```

1 from sklearn.feature_extraction.text import CountVectorizer,
   TfIdfVectorizer
2 import pandas as pd
3 from task_3.preprocessing import load_local_lemma_datasets
4
5
6 def create_bow_datasets(
7     train_dataset,
8     test_dataset,
9     vectorizer_type="count",
10    max_features=1000,
11    ngram_range=(1, 1),
12 ):
13     """
14     Create BOW vectors using either CountVectorizer or TfIdfVectorizer
15     """
16     # Choose vectorizer
17     if vectorizer_type == "tfidf":
18         vectorizer = TfIdfVectorizer(
19             max_features=max_features, ngram_range=ngram_range, min_df=
20             =2
21         ) # Ignore terms that appear in less than 2 documents
22     else:
23         vectorizer = CountVectorizer(
24             max_features=max_features, ngram_range=ngram_range, min_df=
25             =2
26         )
27
28     # Fit and transform training data
29     X_train = vectorizer.fit_transform(train_dataset["text"])
30
31     # Transform test data
32     X_test = vectorizer.transform(test_dataset["text"])
33
34     # Convert to DataFrames

```

```

33 feature_names = vectorizer.get_feature_names_out()
34
35 train_df = pd.DataFrame(X_train.toarray(), columns=feature_names)
36 train_df["sentiment"] = train_dataset["label"]
37
38 test_df = pd.DataFrame(X_test.toarray(), columns=feature_names)
39 test_df["sentiment"] = test_dataset["label"]
40
41 # Print some information about the vectorization
42 print(f"Vocabulary size: {len(vectorizer.vocabulary_)}")
43 print(f"Feature matrix shape: {X_train.shape}")
44 print("\nMost common terms:")
45 if vectorizer_type == "count":
46     term_frequencies = X_train.sum(axis=0).A1
47     top_terms = sorted(
48         zip(vectorizer.get_feature_names_out(), term_frequencies),
49         key=lambda x: x[1],
50         reverse=True,
51     )[:10]
52     for term, freq in top_terms:
53         print(f"{term}: {freq}")
54
55 return train_df, test_df
56
57
58 def get_vector_datasets(vectorization_type="count"):
59     train_dataset, test_dataset = load_local_lemma_datasets()
60
61     train_dataset, test_dataset = create_bow_datasets(
62         train_dataset, test_dataset, vectorization_type
63     )
64
65     return train_dataset, test_dataset

```

File models/classical\_models.py

```

1 from sklearn.ensemble import RandomForestClassifier,
  GradientBoostingClassifier
2 from sklearn.svm import SVC
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score, classification_report
5 from sklearn.preprocessing import StandardScaler
6 from task_3.vectorization import get_vector_datasets
7
8
9 RANDOM_STATE = 42
10
11
12 def classical_models_comparison():
13     # Get data
14     train_df, test_df = get_vector_datasets("word2vect")
15
16     # Prepare features and labels
17     X_train = train_df.drop("sentiment", axis=1).values
18     y_train = train_df["sentiment"].values
19     X_test = test_df.drop("sentiment", axis=1).values
20     y_test = test_df["sentiment"].values
21
22     # Scale features

```

```

23 scaler = StandardScaler()
24 X_train_scaled = scaler.fit_transform(X_train)
25 X_test_scaled = scaler.transform(X_test)
26
27 # Define models to try
28 models = {
29     "Logistic Regression": LogisticRegression(
30         random_state=RANDOM_STATE, max_iter=1000
31     ),
32     "Random Forest": RandomForestClassifier(
33         n_estimators=100, random_state=RANDOM_STATE
34     ),
35     "Gradient Boosting": GradientBoostingClassifier(random_state=
36         RANDOM_STATE),
37     "SVM": SVC(random_state=RANDOM_STATE),
38 }
39
40 # Try each model
41 for name, model in models.items():
42     print(f"\nTraining {name}...")
43     model.fit(X_train_scaled, y_train)
44
45     # Make predictions
46     train_pred = model.predict(X_train_scaled)
47     test_pred = model.predict(X_test_scaled)
48
49     # Print results
50     print(f"{name} Results:")
51     print(f"Training Accuracy: {accuracy_score(y_train, train_pred)
52         :.4f}")
53     print(f"Test Accuracy: {accuracy_score(y_test, test_pred):.4f}"
54         )
55     print("\nDetailed Test Report:")
56     print(classification_report(y_test, test_pred))
57
58     # For Random Forest and Gradient Boosting, print feature
59     importance
60     if hasattr(model, "feature_importances_"):
61         importances = model.feature_importances_
62         top_features = sorted(
63             zip(range(len(importances)), importances),
64             key=lambda x: x[1],
65             reverse=True,
66         )[:10]
67         print("\nTop 10 Most Important Features:")
68         for idx, importance in top_features:
69             print(f"Feature {idx}: {importance:.4f}")

```