


# SENTIMENT ANALYSIS PLATFORM FOR MEDIA PRODUCTIONS

*Sentiment Analysis and Rating Prediction of Movie Reviews*


*Presented by: Lawrence(Zhiwei) Qiu and Gurnoor Singh*

 NISARG CHODVADIYA · UPDATED 3 YEARS AGO

21New NotebookDownload (54 MB)

## IMDB Movie Reviews with ratings

50K labeled review with ratings and sentiment, 50K reviews unlabeled



[Data Card](#) [Code \(6\)](#) [Discussion \(0\)](#)


### About Dataset

Preprocessed Stanford IMDB movie review dataset in CSV format with sentiment and rating column also to predict ratings. Preprocessed from folder structure to CSV by file handling. To know how to extract ratings read the read me file of the original dataset.

**Usability** ⓘ  
8.82

**License**  
Other (specified in description)


<https://www.kaggle.com/datasets/nisargchodavadiya/imdb-movie-reviews-with-ratings-50k>

 LAKSHMIPATHI N · UPDATED 5 YEARS AGO

1036New NotebookDownload (27 MB)

## IMDB Dataset of 50K Movie Reviews

Large Movie Review Dataset



[Data Card](#) [Code \(888\)](#) [Discussion \(9\)](#)

## Development Environment & Tools

- *Python*
- *Logistic Regression & Naive Bayes*
- *FastText*
- *Linear Regression*
- *spaCy*
- *TfidfVectorizer*
- *Jupyter Notebook*

These two 3D scatter plots show the frequency distribution of positive and negative words respectively under different movie ratings.

# Part - 1

## Data Processing and Feature Extraction

```
# Clean the review text
def clean_text(text):
    text = re.sub(r'<.*?>', ' ', text)
    text = re.sub(r'[-_~@-9\s]', ' ', text)
    return text
imdb_data['clean_review'] = imdb_data['review'].apply(clean_text)
print("\nSample review before cleaning:", imdb_data['review'][0])
print("Sample review after cleaning:", imdb_data['clean_review'][0])
```

Sample review before cleaning: One of the other reviewers has mentioned that after war right, as this is exactly what happened with me. <br /><br />The first thing that struck me was the violence, which set in right from the word GO. Trust me, this is not a show for ladies or gentlemen. It is hardcore, in the classic use of the nickname given to the Oswald Maximum Security State Penitentiary. It focuses mainly on prison where all the cells have glass fronts and face inwards, so privacy is not high on the list. Gangsters, Latinos, Christians, Italians, Irish and more... so scuffles, death struggles, and so on. <br /><br />I would say the main appeal of the show is due to the fact that it shows gritty pictures painted for mainstream audiences, forget charm, forget romance... It does not make me feel as so nasty it was surreal. I couldn't say I was ready for it, but as I watched it, I was drawn to the high levels of graphic violence. Not just violence, but injustice (crooked cops, corrupt judges, etc.). All kill on order and get away with it, well mannered, middle class inmates being taught skills or prison experience. Watching it, you may become comfortable with what is going on with your darker side.

Sample review after cleaning: One of the other reviewers has mentioned that after war

```
# Load the dataset
imdb_data = pd.read_csv("review_sent.csv")
print("Dataset Loaded")
print("Dataset shape:", imdb_data.shape)
print("First few rows of the dataset:\n", imdb_data.head())
```

Dataset Loaded

Dataset shape: (50000, 2)

First few rows of the dataset:

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production.   The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive

```
# TF-IDF vectorization
tfidf_vectorizer = TfidfVectorizer(max_features=5000, stop_words='english', ngram_range=(1, 2))
features = tfidf_vectorizer.fit_transform(imdb_data['clean_review'])
print("\nTF-IDF Vectorization complete")
print("Shape of the TF-IDF matrix:", features.shape)
```

TF-IDF Vectorization complete  
Shape of the TF-IDF matrix: (50000, 5000)

```
# Splitting the dataset
X_train, X_test, y_train, y_test = train_test_split(features, imdb_data['sentiment'], test_size=0.2, random_state=42)
print("\nTrain-Test Split complete")
print("X_train shape:", X_train.shape)
print("X_test shape:", X_test.shape)
print("y_train shape:", y_train.shape)
print("y_test shape:", y_test.shape)
```

Train-Test Split complete  
X\_train shape: (40000, 5000)  
X\_test shape: (10000, 5000)  
y\_train shape: (40000,)  
y\_test shape: (10000,)

## Model Training & Evaluation

```
# Logistic Regression Model
logistic_model = LogisticRegression(max_iter=1000)
logistic_model.fit(X_train, y_train)
```

```
LogisticRegression
LogisticRegression(max_iter=1000)
```

```
# Naive Bayes Model
naive_bayes_model = MultinomialNB()
naive_bayes_model.fit(X_train, y_train)
```

```
MultinomialNB
MultinomialNB()
```

```
# Evaluation
evaluation = {
    'Model': ['Logistic Regression', 'Naive Bayes'],
    'Accuracy': [accuracy_score(y_test, logistic_predictions), accuracy_score(y_test, naive_bayes_predictions)],
    'Precision': [precision_score(y_test, logistic_predictions, pos_label='positive'), precision_score(y_test, naive_bayes_predictions, pos_label='positive')],
    'Recall': [recall_score(y_test, logistic_predictions, pos_label='positive'), recall_score(y_test, naive_bayes_predictions, pos_label='positive')],
    'F1-Score': [f1_score(y_test, logistic_predictions, pos_label='positive'), f1_score(y_test, naive_bayes_predictions, pos_label='positive')]
}

evaluation_df = pd.DataFrame(evaluation)
print(evaluation_df)
```

	Model	Accuracy	Precision	Recall	F1-Score
0	Logistic Regression	0.8900	0.879261	0.906132	0.892494
1	Naive Bayes	0.8541	0.845426	0.869419	0.857255

## Results of Sentiment Analysis Model

```
# Extracting Coefficients
feature_names = tfidf_vectorizer.get_feature_names_out() # Updated method
coefficients = model.coef_.flatten()

# Combining coefficients with feature names
feature_importance = pd.DataFrame({'Feature': feature_names, 'Coefficient': coefficients})

# Sorting by absolute coefficients to see the most influential words
top_influential_words = feature_importance.reindex(feature_importance.Coefficient.abs().sort_values(ascending=False))

# Display top 20 words
print(top_influential_words.head(20))
```

	Feature	Coefficient
4942	worst	-9.417322
345	awful	-7.310270
515	boring	-6.866844
4824	waste	-6.865693
1509	excellent	6.759788
352	bad	-6.745389
1978	great	6.744836
3387	poor	-5.731439
4451	terrible	-5.642180
1351	dull	-5.274667
3388	poorly	-5.129762
4941	worse	-5.077302
432	best	5.075981
4915	wonderful	5.017101
207	amazing	4.962230
1230	disappointment	-4.935761
3288	perfect	4.903917
2181	horrible	-4.746227
553	brilliant	4.642281
1229	disappointing	-4.598668



# Part - 2

## Sentiment Analysis with FastText and Rating Predictions

## Data Preprocessing and Feature Extraction for Rating Prediction

Before:

After:

Before:

After:

Before:

After:

Before:

After:

[illegible][illegible]

**Top 10 Words in Positive Reviews**

Word	Frequency (approx.)
film	40,000
movie	37,000
one	26,000
like	17,000
good	15,000
great	13,000
story	13,000
time	12,000
see	12,000
well	11,000

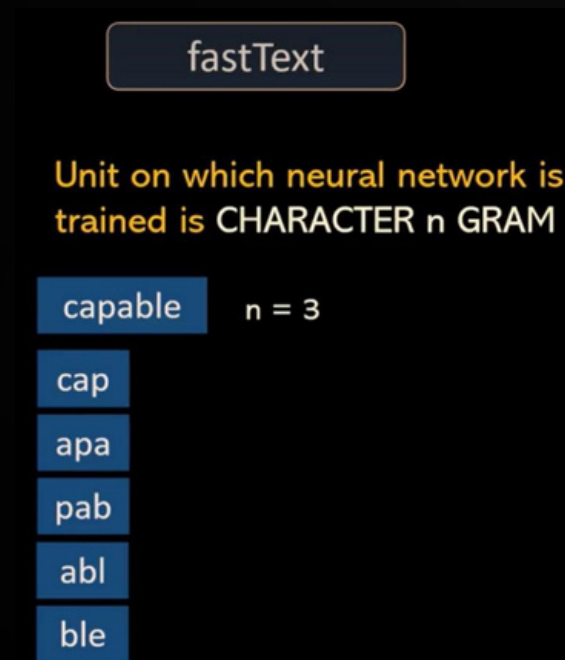
**Top 10 Words in Negative Reviews**

Word	Frequency (approx.)
movie	48,000
film	36,000
one	25,000
like	22,000
even	15,000
bad	14,000
good	14,000
would	13,000
really	12,000
time	11,000

The figure consists of two 3D scatter plots side-by-side. The left plot is titled 'Positive Adjectives' and the right plot is titled 'Negative Adjectives'. Both plots share the same axes: the x-axis is 'Rating' (ranging from 2 to 10), the y-axis is 'Frequency' (ranging from 0 to 5000), and the z-axis is 'Random Z Value' (ranging from 0.0 to 0.8). In the 'Positive Adjectives' plot, words like 'good', 'great', 'excellent', 'amazing', 'best', and 'awesome' are scattered across the plot area. In the 'Negative Adjectives' plot, words like 'terrible', 'horrible', 'bad', 'awful', 'poor', and 'disastrous' are scattered across the plot area. The distribution of points in both plots appears relatively uniform across the 3D space.

# Model Training and Evaluation

## Classifying Movie Reviews with FastText



Average deviation of our predictions from the actual ratings

## Problem

Before

```
# Make predictions on the test set and calculate evaluation metrics
test_df['predict_label'] = [classifier.predict(' ' + text)[0][0].replace('__label__', '') for text in test_df['Clean_review']]
true_labels = test_df['Sentiment']
predicted_labels = test_df['predict_label']
# Calculate and output performance metrics
accuracy = accuracy_score(true_labels, predicted_labels)
precision = precision_score(true_labels, predicted_labels, pos_label='positive')
recall = recall_score(true_labels, predicted_labels, pos_label='positive')
f1 = f1_score(true_labels, predicted_labels, pos_label='positive')
print(f'Accuracy: {accuracy}\nPrecision: {precision}\nRecall: {recall}\nF1 Score: {f1}')
```

Accuracy: 0.8815  
Precision: 0.0  
Recall: 0.0  
F1 Score: 0.0

After

```
# Generate a detailed classification report
report = classification_report(true_labels, predicted_labels)
print("Classification Report:\n", report)
```

Model Accuracy: 0.9103  
Model Precision: 0.9098901098901099  
Model Recall: 0.9108  
Model F1 Score: 0.9103448275862069  
Classification Report:

	precision	recall	f1-score	support
negative	0.91	0.91	0.91	5000
positive	0.91	0.91	0.91	5000
accuracy			0.91	10000
macro avg	0.91	0.91	0.91	10000
weighted avg	0.91	0.91	0.91	10000

```
# Evaluate the model
# Calculate Mean Squared Error (MSE) - lower values are better
mse = mean_squared_error(y_test, y_pred)
print(f'Mean Squared Error: {mse}')
```

```
# Calculate R-squared score - values closer to 1 indicate better fit
r2 = r2_score(y_test, y_pred)
print(f'R-squared: {r2}')
```

Mean Squared Error: 4.805162198134217  
R-squared: 0.6007543201980878



# **CONCLUSION & FUTURE WORK**