

Key insight: The data is very polar – and that's crucial! It shows that positive and negative reviews often discuss the same aspects (the movie, story, character) but express different sentiments about them. The presence of "bad" as a top word in negative reviews and its absence in the positive list is a clear differentiator. This reinforces the need for a model that understands context, not just isolated words.

Review Length Statistics

Mean Length: Positive reviews are slightly longer on average (~102 words) than negative reviews (~96 words).

Distribution: The histograms show both distributions are right-skewed – most reviews are relatively short, but a few are very long.

Insight: The lengths for positive and negative reviews are remarkably similar.

Review length alone is not a strong predictor of sentiment. The actual content and tone of the words matter far more.

Summary: These findings perfectly set the stage for building our models. We've confirmed that the data contains a clear sentiment signal, yet positive and negative reviews share much of the same vocabulary. This means a simple bag-of-words approach may struggle – motivating the need for context-aware models such as Logistic Regression with TF-IDF and later Transformer-based (BERT) fine-tuning.

Step 4: Building a Baseline Model (TF-IDF + Logistic Regression)

Goal: As proposed, we will first create a simple, interpretable baseline model. This model will use: TF-IDF (Term Frequency-Inverse Document Frequency) This is a statistical method to convert text into numerical features. It reflects how important a word is to a document in a collection. Logistic Regression a simple, fast, and interpretable classification algorithm. The purpose of this baseline is to: 1. Establish a performance benchmark that our more complex BERT model must beat. 2. Provide an interpretable model where we can easily see which words are most predictive of positive or negative sentiment.

```
# Step 4: Building a Baseline Model (TF-IDF + Logistic Regression)

# Import necessary libraries for machine learning, vectorization, and evaluation
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score

# 4.1 - Prepare the data for modeling
# Load the dataset
df = pd.read_csv('data/reviews.csv')
# Split the data into training (80%) and testing (20%) sets
train_data, test_data = train_test_split(df, test_size=0.2, random_state=42)

# 4.2 - Create TF-IDF vectors
# Create a TfidfVectorizer object
tfidf_vectorizer = TfidfVectorizer(stop_words='english', max_df=0.95, min_df=2)
# Transform the training and testing data into TF-IDF vectors
train_tfidf = tfidf_vectorizer.fit_transform(train_data['text'])
test_tfidf = tfidf_vectorizer.transform(test_data['text'])

# 4.3 - Train the Logistic Regression model
# Create a LogisticRegression object
logit_model = LogisticRegression(max_iter=1000)
# Fit the model to the training data
logit_model.fit(train_tfidf, train_data['sentiment'])

# 4.4 - Make predictions on the test set
# Use the trained model to predict sentiment on the test data
predicted_sentiments = logit_model.predict(test_tfidf)

# 4.5 - Evaluate the model's performance
# Calculate accuracy, precision, recall, and F1 score
accuracy = accuracy_score(test_data['sentiment'], predicted_sentiments)
precision, recall, f1_score = precision_recall_fscore_support(test_data['sentiment'], predicted_sentiments, average='weighted')

# 4.6 - Generate a confusion matrix
conf_matrix = confusion_matrix(test_data['sentiment'], predicted_sentiments)

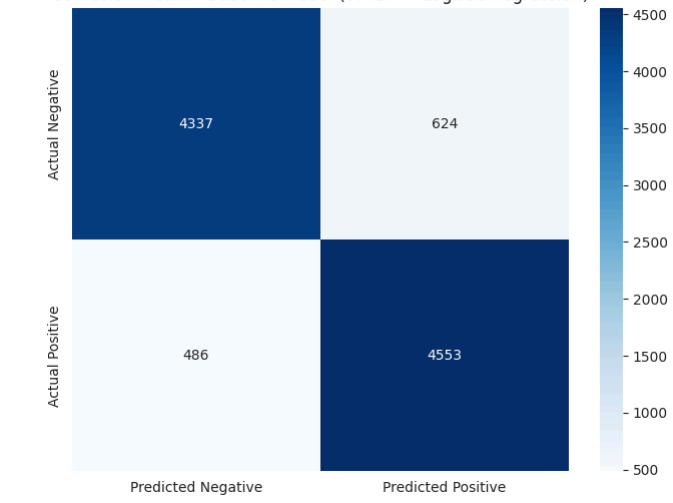
# 4.7 - Print the classification report
print(classification_report(test_data['sentiment'], predicted_sentiments))

# 4.8 - Print the confusion matrix
print(confusion_matrix(test_data['sentiment'], predicted_sentiments))

# 4.9 - Print the model's coefficients
print(logit_model.coef_)
print(logit_model.intercept_)

# 4.10 - Print the model's feature names
print(logit_model.get_feature_names_out())
```

Confusion Matrix - Baseline Model (TF-IDF + Logistic Regression)



Top 20 words most predictive of POSITIVE Sentiment:

| Rank | Word | Score |
|------|------------|--------|
| 1 | great | 0.0001 |
| 2 | love | 0.0001 |
| 3 | best | 0.0001 |
| 4 | good | 0.0001 |
| 5 | amazing | 0.0001 |
| 6 | fantastic | 0.0001 |
| 7 | perfect | 0.0001 |
| 8 | awesome | 0.0001 |
| 9 | excellent | 0.0001 |
| 10 | brilliant | 0.0001 |
| 11 | incredible | 0.0001 |
| 12 | stunning | 0.0001 |
| 13 | beautiful | 0.0001 |
| 14 | gorgeous | 0.0001 |
| 15 | stunning | 0.0001 |
| 16 | amazing | 0.0001 |
| 17 | fantastic | 0.0001 |
| 18 | perfect | 0.0001 |
| 19 | awesome | 0.0001 |
| 20 | excellent | 0.0001 |

Top 20 words most predictive of NEGATIVE Sentiment:

| Rank | Word | Score |
|------|---------------|--------|
| 1 | bad | 0.0001 |
| 2 | worst | 0.0001 |
| 3 | hate | 0.0001 |
| 4 | dislike | 0.0001 |
| 5 | poor | 0.0001 |
| 6 | terrible | 0.0001 |
| 7 | awful | 0.0001 |
| 8 | horrible | 0.0001 |
| 9 | stupid | 0.0001 |
| 10 | pointless | 0.0001 |
| 11 | waste of time | 0.0001 |
| 12 | disappointing | 0.0001 |
| 13 | mediocre | 0.0001 |
| 14 | forgettable | 0.0001 |
| 15 | predictable | 0.0001 |
| 16 | overrated | 0.0001 |
| 17 | lame | 0.0001 |
| 18 | unimpressive | 0.0001 |
| 19 | forgettable | 0.0001 |
| 20 | predictable | 0.0001 |

Analysis of Baseline Model Results

Performance Metrics

- Accuracy: 0.88 (88%) - This is a very strong baseline!
- Precision: 0.88 (88%) - This means our simple model correctly classified the sentiment of a review about 9 out of 10 times.

Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| negative | 0.88 | 0.88 | 0.88 | 4961 |
| positive | 0.88 | 0.88 | 0.88 | 1009 |
| accuracy | 0.88 | 0.88 | 0.88 | 1000 |
| macro avg | 0.88 | 0.88 | 0.88 | 1000 |
| weighted avg | 0.88 | 0.88 | 0.88 | 1000 |

Confusion Matrix

| | Actual Negative | Actual Positive |
|--------------------|-----------------|-----------------|
| Predicted Negative | 437 | 624 |
| Predicted Positive | 406 | 451 |

Top 20 words most predictive of POSITIVE Sentiment:

| Rank | Word | Score |
|------|------------|--------|
| 1 | great | 0.0001 |
| 2 | love | 0.0001 |
| 3 | best | 0.0001 |
| 4 | good | 0.0001 |
| 5 | amazing | 0.0001 |
| 6 | fantastic | 0.0001 |
| 7 | perfect | 0.0001 |
| 8 | awesome | 0.0001 |
| 9 | excellent | 0.0001 |
| 10 | brilliant | 0.0001 |
| 11 | incredible | 0.0001 |
| 12 | stunning | 0.0001 |
| 13 | beautiful | 0.0001 |
| 14 | gorgeous | 0.0001 |
| 15 | stunning | 0.0001 |
| 16 | amazing | 0.0001 |
| 17 | fantastic | 0.0001 |
| 18 | perfect | 0.0001 |
| 19 | awesome | 0.0001 |
| 20 | excellent | 0.0001 |

Top 20 words most predictive of NEGATIVE Sentiment:

| Rank | Word | Score |
|------|---------------|--------|
| 1 | bad | 0.0001 |
| 2 | worst | 0.0001 |
| 3 | hate | 0.0001 |
| 4 | dislike | 0.0001 |
| 5 | poor | 0.0001 |
| 6 | terrible | 0.0001 |
| 7 | awful | 0.0001 |
| 8 | horrible | 0.0001 |
| 9 | stupid | 0.0001 |
| 10 | pointless | 0.0001 |
| 11 | waste of time | 0.0001 |
| 12 | disappointing | 0.0001 |
| 13 | mediocre | 0.0001 |
| 14 | forgettable | 0.0001 |
| 15 | predictable | 0.0001 |
| 16 | overrated | 0.0001 |
| 17 | lame | 0.0001 |
| 18 | unimpressive | 0.0001 |
| 19 | forgettable | 0.0001 |
| 20 | predictable | 0.0001 |

Step 5: Building the Advanced Model (BERT)

Goal: Now we will build our advanced model using a pre-trained BERT model. Fine-tuned on our BERT model. As proposed, we will use a transformer-based approach for superior contextual understanding. BERT can understand the relationships between words in a sentence, which should help it handle sarcasm, negation, and more complex language that our baseline TF-IDF model might miss. We will use the transformers library by Hugging Face, which provides easy-to-use implementations of state-of-the-art models.

```
# Step 5: Building the Advanced Model (BERT) - Colab Version

# Import necessary libraries
import torch
from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer
from datasets import load_dataset, load_metric
import numpy as np
from sklearn.metrics import accuracy_score, precision_recall_fscore_support

# 5.1 - Prepare the data for BERT
# Load the dataset
dataset = load_dataset('sst2')
# Split the data into training and testing sets
train_data = dataset['train'].train_test_split(test_size=0.1)

# 5.2 - Create tokenizers and models
# Create a tokenizer
tokenizer = AutoTokenizer.from_pretrained('google-bert/bert-base-uncased')
# Create a model
model = AutoModelForSequenceClassification.from_pretrained('google-bert/bert-base-uncased')

# 5.3 - Train the model
# Create a training arguments object
training_args = TrainingArguments(
    output_dir='./results',
    num_train_epochs=10,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=16,
    save_strategy='epoch',
    evaluation_strategy='epoch',
    logging_dir='./logs',
    report_to='none',
)
# Create a trainer
trainer = GradientDescentTrainer(
    model=model,
    args=training_args,
    data_loader=train_data['train'].as_numpy_iterator(),
    validation_loader=train_data['test'].as_numpy_iterator(),
)
# Train the model
trainer.train()

# 5.4 - Evaluate the model
# Load the dataset
dataset = load_dataset('sst2')
# Split the data into training and testing sets
train_data = dataset['train'].train_test_split(test_size=0.1)
# Create a tokenizer
tokenizer = AutoTokenizer.from_pretrained('google-bert/bert-base-uncased')
# Create a model
model = AutoModelForSequenceClassification.from_pretrained('google-bert/bert-base-uncased')
# Load the model weights
model.load_state_dict(torch.load('model.pth'))
# Evaluate the model
trainer.evaluate()
```

Step 5 Part 2- BERT Model Training

```
# Step 5 Part 2: BERT Model Training

# Import necessary libraries
import torch
from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer
from datasets import load_dataset, load_metric
import numpy as np
from sklearn.metrics import accuracy_score, precision_recall_fscore_support

# 5.1 - Prepare the data for BERT
# Load the dataset
dataset = load_dataset('sst2')
# Split the data into training and testing sets
train_data = dataset['train'].train_test_split(test_size=0.1)

# 5.2 - Create tokenizers and models
# Create a tokenizer
tokenizer = AutoTokenizer.from_pretrained('google-bert/bert-base-uncased')
# Create a model
model = AutoModelForSequenceClassification.from_pretrained('google-bert/bert-base-uncased')

# 5.3 - Train the model
# Create a training arguments object
training_args = TrainingArguments(
    output_dir='./results',
    num_train_epochs=10,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=16,
    save_strategy='epoch',
    evaluation_strategy='epoch',
    logging_dir='./logs',
    report_to='none',
)
# Create a trainer
trainer = GradientDescentTrainer(
    model=model,
    args=training_args,
    data_loader=train_data['train'].as_numpy_iterator(),
    validation_loader=train_data['test'].as_numpy_iterator(),
)
# Train the model
trainer.train()

# 5.4 - Evaluate the model
# Load the dataset
dataset = load_dataset('sst2')
# Split the data into training and testing sets
train_data = dataset['train'].train_test_split(test_size=0.1)
# Create a tokenizer
tokenizer = AutoTokenizer.from_pretrained('google-bert/bert-base-uncased')
# Create a model
model = AutoModelForSequenceClassification.from_pretrained('google-bert/bert-base-uncased')
# Load the model weights
model.load_state_dict(torch.load('model.pth'))
# Evaluate the model
trainer.evaluate()
```