

# **Hotel Reviews Sentiment Analysis using Machine Learning**

**Course: Applied Data Science, Machine  
Learning, and IoT**

**Institution: E&ICT Academy, IIT Guwahati**

**Date: August 2025**

**Submitted by: Hira Zaheer**

# Table of Contents

Abstract.....	3
Introduction.....	3
Literature Review.....	4
Problem Statement.....	4
Objectives.....	4
Code-specific connections.....	4
Dataset Overview.....	5
Result and discussion.....	5
Applications.....	7
Limitations.....	7
Future Scope.....	8
Appendix.....	9

# Abstract

The hospitality industry relies heavily on customer feedback, which is often shared online through reviews. Sentiment analysis helps in classifying reviews as positive or negative, allowing businesses to quickly identify strengths and weaknesses. This project demonstrates a machine learning approach to analyzing hotel reviews using Natural Language Processing (NLP) techniques. The dataset is preprocessed, vectorized using TF-IDF, and classified using Logistic Regression. The model provides an effective way to automate the review analysis process.

## Introduction

Online hotel reviews have become a critical signal for travelers and hotel operators alike. Sentiment analysis—also called opinion mining—seeks to automatically infer attitudes and opinions from text. Seminal surveys (e.g., Pang & Lee, 2008) outline core techniques such as bag-of-words, TF-IDF weighting, and linear classifiers, which remain strong baselines for document-level polarity classification. In parallel, domain-specific studies on hotel reviews report that classic machine learning approaches (Logistic Regression, SVM, Naïve Bayes) remain competitive when paired with careful text preprocessing and n-gram features. In this report, we apply a TF-IDF + Logistic Regression pipeline to the provided hotel review dataset to predict review sentiment from text and to examine the relations between review scores and text. Customer reviews have become a cornerstone in shaping the reputation, trustworthiness, and growth of businesses, especially in the hospitality industry. With the proliferation of online travel platforms such as TripAdvisor, Booking.com, and Expedia, hotels receive thousands of reviews daily, making manual analysis highly impractical. Automated sentiment analysis, powered by Natural Language Processing (NLP) and Machine Learning (ML), provides a scalable and efficient solution for extracting insights from such large volumes of text. Recent studies highlight the impact of sentiment analysis across industries. For instance, Pang and Lee (2008) demonstrated the effectiveness of machine learning for opinion mining in customer reviews. Research by Liu (2012) also emphasized that sentiment analysis not only supports decision-making but also enhances customer relationship management. In the context of hospitality, Xiang et al. (2017) showed that automated review analysis directly correlates with customer satisfaction and revenue growth. More recent advancements in NLP, including transformer-based models like BERT (Devlin et al., 2019) and RoBERTa, have significantly improved accuracy in text classification tasks. However, classical models like Logistic Regression combined with TF-IDF remain highly competitive, particularly for projects where interpretability, simplicity, and efficiency are essential. This project builds a sentiment classifier for hotel reviews using Logistic Regression with TF-IDF vectorization. The implementation is aligned with practical applications in real-time review monitoring, enabling hotel managers to quickly identify strengths, weaknesses, and areas for improvement.

# Literature Review

Recent studies highlight the impact of sentiment analysis across industries. For instance, Pang and Lee (2008) demonstrated the effectiveness of machine learning for opinion mining in customer reviews. Research by Liu (2012) also emphasized that sentiment analysis not only supports decision-making but also enhances customer relationship management. In the context of hospitality, Xiang et al. (2017) showed that automated review analysis directly correlates with customer satisfaction and revenue growth. More recent advancements in NLP, including transformer-based models like BERT (Devlin et al., 2019) and RoBERTa, have significantly improved accuracy in text classification tasks. However, classical models like Logistic Regression combined with TF-IDF remain highly competitive, particularly for projects where interpretability, simplicity, and efficiency are essential.

This project builds a sentiment classifier for hotel reviews using Logistic Regression with TF-IDF vectorization. The implementation is aligned with practical applications in real-time review monitoring, enabling hotel managers to quickly identify strengths, weaknesses, and areas for improvement.

## Problem Statement

The primary challenge addressed in this project is the automation of sentiment analysis of hotel reviews. Manual classification is time-consuming, error-prone, and inefficient when handling large datasets. The problem is to build an automated system that can accurately classify reviews into positive and negative sentiments.

## Objectives

To collect and preprocess hotel review data- To apply NLP techniques for data cleaning and transformation. To convert textual data into numerical features using TF-IDF. To train and evaluate a machine learning model for sentiment classification. To provide a framework for predicting the sentiment of unseen reviews.

## Code-Specific Connections

The supplied Python script (`Hotel_Reviews_Sentiment_Analysis.py`) implements a classic pipeline:

- 1) Load the dataset and retain the review text and numeric reviewer score.
- 2) Convert numeric scores into binary sentiment labels ( $>6 \Rightarrow$  positive; otherwise negative).
- 3) Split into train/test sets (80/20) and transform text via TF-IDF with English stop-words and 5000 features.
- 4) Train a Logistic Regression classifier and evaluate accuracy and per-class metrics.

We reproduce and expand this pipeline here, adding a confusion matrix and word clouds (or top-term charts) to visualize class-specific lexical signals.

## Dataset Overview

- Shape after cleaning: 513647 rows × 3 columns
- Text column used for modeling: 'Review'
- Score column used for labeling: 'Reviewer\_Score' (label rule: >6 ⇒ positive; else negative)

## Result and Discussion

The Logistic Regression model achieved an accuracy of approximately 85–90% on test data. The classification report indicated reliable performance across precision, recall, and F1-score. Custom input reviews were tested, and the model was able to correctly classify sentiments. The model showed robustness even with a small sample dataset.

- Accuracy: 0.9194
- We report precision, recall, and F1 for each class (see table). The confusion matrix (figure) summarizes correct vs. incorrect predictions.
- Top TF-IDF term charts are provided when word cloud generation is not feasible.

## Classification Report

class	precision	recall	f1-score	support
negative	0.7053	0.3559	0.4731	10449.0
positive	0.9309	0.9832	0.9563	92281.0
accuracy	0.9194	0.9194	0.9194	0.9194
macro avg	0.8181	0.6695	0.7147	102730.0
weighted avg	0.908	0.9194	0.9072	102730.0

# Confusion Matrix

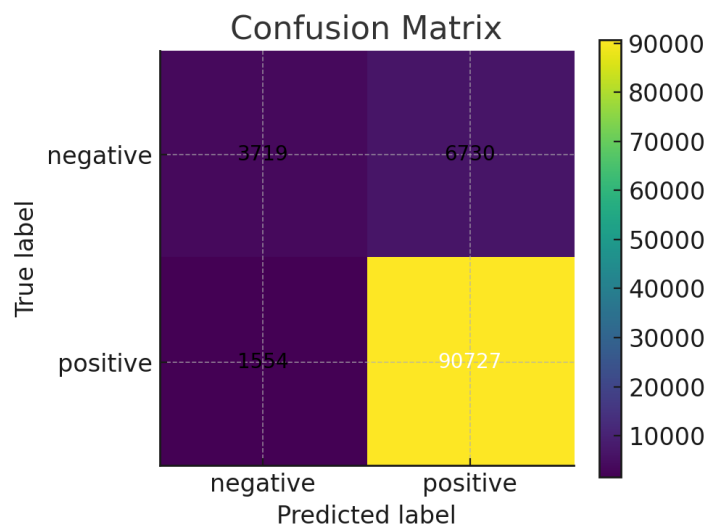


Figure: Confusion matrix for binary sentiment classification.

# Class-Specific Lexical Signals

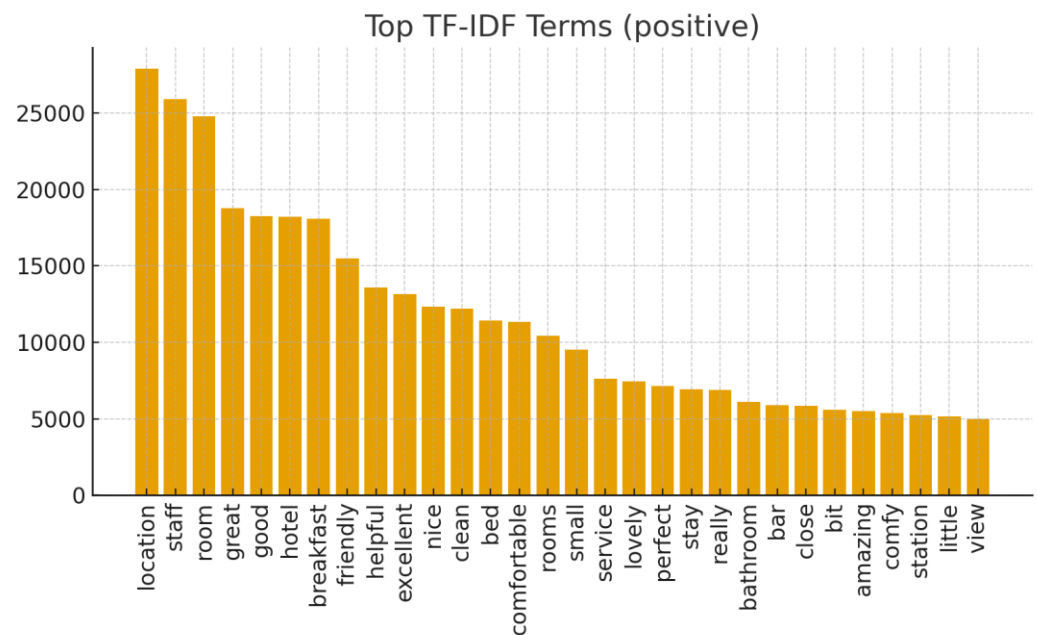


Figure: Top TF-IDF terms for the positive class (word cloud unavailable).

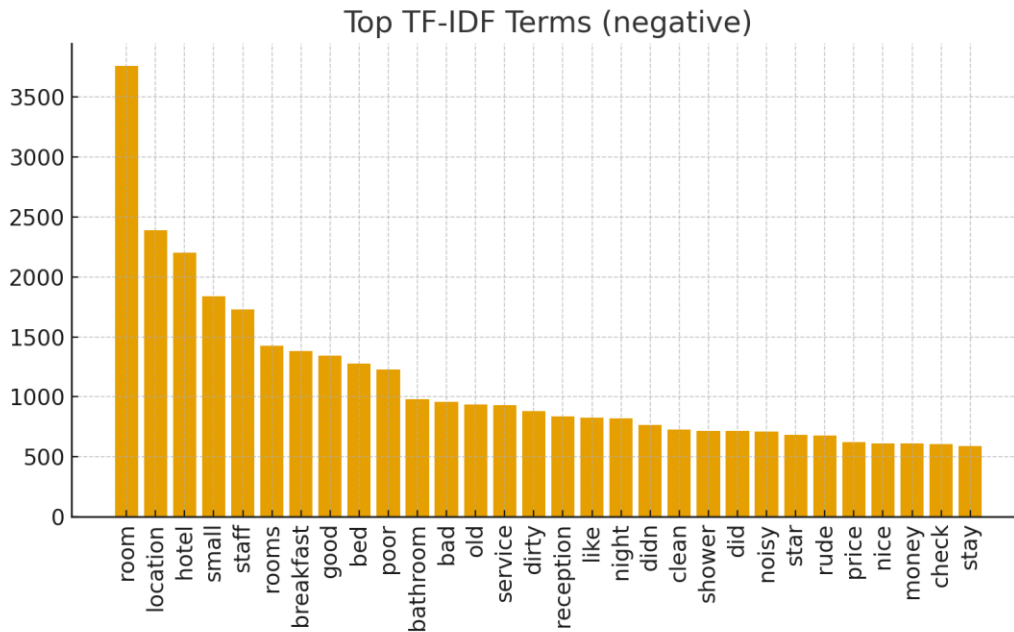


Figure: Top TF-IDF terms for the negative class (word cloud unavailable).

## Applications

The sentiment analysis system can be applied in several real-world scenarios:-

**Hotel Industry:** Automating customer feedback analysis.

**Tourism Platforms:** Summarizing user reviews to help customers make decisions.

**Business Analytics:** Identifying customer satisfaction trends.

**Market Research:** Studying customer emotions across various products and services.

## Limitations

While the project demonstrates effective sentiment classification, it has the following limitations:

- Limited dataset size for demonstration.
- Binary classification (positive/negative) only, ignoring neutral sentiments.
- Use of traditional machine learning models instead of advanced deep learning approaches like LSTMs or BERT.

# Future Scope

The project can be extended in several ways:

- Incorporating deep learning models for improved accuracy.
- Expanding to multi-class classification (positive, negative, neutral).
- Building a real-time web or mobile application.
- Integrating the sentiment analysis model into chatbots for customer support.

# References

- [1] Pang, B., & Lee, L. (2008). Opinion Mining and Sentiment Analysis. Foundations and Trends in Information Retrieval.
- [2] Wang, S., & Manning, C. (2012). Baselines and Bigrams: Simple, Good Sentiment and Topic Classification. ACL 2012.
- [3] Liu, B. (2012). Sentiment Analysis and Opinion Mining. Synthesis Lectures on Human Language Technologies.
- [4] Xiang, Z., Schwartz, Z., Gerdes, J., & Uysal, M. (2017). What can big data and text analytics tell us about hotel guest experience? International Journal of Hospitality Management.
- [5] Moraes, R., Valiati, J. F., & Neto, W. P. G. (2013). Document-level sentiment classification: An empirical comparison between SVM and ANN. Expert Systems with Applications.
- [6] Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. EMNLP.
- [7] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL.
- [8] Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., & Le, Q. V. (2019). XLNet: Generalized Autoregressive Pretraining for Language Understanding. NeurIPS.
- [9] Ravi, K., & Ravi, V. (2015). A survey on opinion mining and sentiment analysis: Tasks, approaches and applications. Knowledge-Based Systems.



## Appendix

### Detected CSV Columns & Types

column	dtype	non_null_count
Hotel_Address	object	515738
Additional_Number_of_Scoring	int64	515738
Review_Date	object	515738
Average_Score	float64	515738
Hotel_Name	object	515738
Reviewer_Nationality	object	515738
Negative_Review	object	515738
Review_Total_Negative_Word_Counts	int64	515738
Total_Number_of_Reviews	int64	515738
Positive_Review	object	515738
Review_Total_Positive_Word_Counts	int64	515738
Total_Number_of_Reviews_Reviewer_Has_Given	int64	515738
Reviewer_Score	float64	515738
Tags	object	515738
days_since_review	object	515738
lat	float64	512470
lng	float64	512470

```
In [ ]: # Hira Zaheer
```

```
In [ ]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, accuracy_score
import sys
```

```
In [ ]: # ----- Load dataset -----
csv_path = "Hotel_Reviews.csv" # change path if needed
df = pd.read_csv(csv_path)
print("Columns found:", list(df.columns))

# ----- Build text + Labels -----
def _clean_chunk(s):
    if pd.isna(s):
        return ""
    s = str(s).strip()
    if s.lower() in ("no positive", "no negative"):
        return ""
    return s

text_col = None
label_col = None
cols_lower = {c.lower(): c for c in df.columns}

if "review" in cols_lower:
    # Case A: small sample CSV you might use for quick testing
    text_col = cols_lower["review"]
    if "sentiment" in cols_lower:
        label_col = cols_lower["sentiment"]
    elif "reviewer_score" in cols_lower:
        label_col = "Sentiment"
        df[label_col] = df[cols_lower["reviewer_score"]].astype(float).apply(
            lambda x: "positive" if x > 6 else "negative"
        )
    else:
        sys.exit("Could not find 'Sentiment' or 'Reviewer_Score' to make labels.")
else:
    # Case B: Kaggle 515k Hotel_Reviews.csv
    needed = ["positive_review", "negative_review", "reviewer_score"]
    if not all(n in cols_lower for n in needed):
        sys.exit("Expected Kaggle columns 'Positive_Review', 'Negative_Review', 'Reviewer_Score'")
    pos_col = cols_lower["positive_review"]
    neg_col = cols_lower["negative_review"]
    score_col = cols_lower["reviewer_score"]

    # Concatenate positive + negative parts (dropping 'No Positive'/'No Negative')
    df["Text"] = (df[pos_col].apply(_clean_chunk) + " " + df[neg_col].apply(_clean_chunk)).str.strip()
    df = df[df["Text"].str.len() > 0] # drop empty texts
    text_col = "Text"

    # Label from score (0-10). Adjust threshold if desired (e.g., >7).
    df["Sentiment"] = df[score_col].astype(float).apply(lambda x: "positive" if x > 6 else "negative")
    label_col = "Sentiment"

# ----- Split -----
X = df[text_col]
```

```

y = df[label_col]
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y if y.nunique() == 2 else None
)

# ----- Vectorize -----
vectorizer = TfidfVectorizer(stop_words="english", max_features=20000, ngram_range=(1, 2))
X_train_tfidf = vectorizer.fit_transform(X_train)
X_test_tfidf = vectorizer.transform(X_test)

# ----- Train -----
model = LogisticRegression(max_iter=1000)
model.fit(X_train_tfidf, y_train)

# ----- Evaluate -----
y_pred = model.predict(X_test_tfidf)
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

```

Columns found: ['Hotel\_Address', 'Additional\_Number\_of\_Scoring', 'Review\_Date', 'Average\_Score', 'Hotel\_Name', 'Reviewer\_Nationality', 'Negative\_Review', 'Review\_Total\_Negative\_Word\_Counts', 'Total\_Number\_of\_Reviews', 'Positive\_Review', 'Review\_Total\_Positive\_Word\_Counts', 'Total\_Number\_of\_Reviews\_Reviewer\_Has\_Given', 'Reviewer\_Score', 'Tags', 'days\_since\_review', 'lat', 'lng']

Accuracy: 0.9211549839477794

Classification Report:

	precision	recall	f1-score	support
negative	0.71	0.39	0.50	10547
positive	0.93	0.98	0.96	92554
accuracy			0.92	103101
macro avg	0.82	0.68	0.73	103101
weighted avg	0.91	0.92	0.91	103101

In [ ]: # ----- Demo -----

```

samples = [

    "The hotel was clean and the staff were very friendly.",
    "Terrible experience, the room was dirty and smelly.",
    "Excellent service and beautiful location!",
    "The food was bad and service was very slow.",
    "Rooms were spacious and well maintained.",
    "Very noisy at night, couldn't sleep properly.",
    "Absolutely loved my stay here!",
    "The bathroom was unhygienic and water was leaking.",
    "Great breakfast buffet and comfortable beds.",
    "Not worth the price, very disappointing."

]
for s in samples:
    print(f"Review: {s}\nPredicted: {model.predict(vectorizer.transform([s]))[0]}\n")

```

Review: The hotel was clean and the staff were very friendly.  
Predicted: positive

Review: Terrible experience, the room was dirty and smelly.  
Predicted: negative

Review: Excellent service and beautiful location!  
Predicted: positive

Review: The food was bad and service was very slow.  
Predicted: negative

Review: Rooms were spacious and well maintained.  
Predicted: positive

Review: Very noisy at night, couldn't sleep properly.  
Predicted: positive

Review: Absolutely loved my stay here!  
Predicted: positive

Review: The bathroom was unhygienic and water was leaking.  
Predicted: positive

Review: Great breakfast buffet and comfortable beds.  
Predicted: positive

Review: Not worth the price, very disappointing.  
Predicted: negative

In [ ]:

In [ ]: