# Advanced Certification in Applied Data Science, Machine Learning & IoT

By E&ICT Academy, IIT Guwahati

## Project: Hotel Reviews Sentiment Analysis

**Name:** Rajesh Bisht

**Email:** rbisht.india@gmail.com

**Batch:** 8 (Eight)

**Git Link:** RajeshBisht28/IIT_ML_HotelReviewSentiment

**Description:** ML project for Hotel Review Sentiment Analysis.

In [ ]:

In [ ]:

# 🏨 Hotel Reviews sentiment analysis.

## 📌 Dataset - Column Explanation:

The dataset contains hotel reviews with the following columns:

| Column Name | Description |
| --- | --- |
| Hotel_Address | Address of the hotel. |
| Additional_Number_of_Scoring | Additional number of reviews for the hotel. |
| Review_Date | Date when the review was posted. |
| Average_Score | Average score of the hotel based on all reviews. |
| Hotel_Name | Name of the hotel. |
| Reviewer_Nationality | Nationality of the reviewer. |
| Negative_Review | Text of the negative review (if any). |
| Review_Total_Negative_Word_Counts | Word count of the negative review. |
| Total_Number_of_Reviews | Total number of reviews for the hotel. |
| Positive_Review | Text of the positive review (if any). |
| Review_Total_Positive_Word_Counts | Word count of the positive review. |
| Total_Number_of_Reviews_Reviewer_Has_Given | Total reviews the reviewer has given across hotels. |
| Reviewer_Score | Score given by the reviewer. |
| Tags | Tags related to the review (e.g., trip type, room type). |

| Column Name | Description |
| --- | --- |
| **days_since_review** | Days since the review was posted. |
| **lat** | Latitude of the hotel. |
| **lng** | Longitude of the hotel. |

## ❌ Columns **NOT** Important for Sentiment Analysis:

These columns do not directly impact sentiment classification:

- 🏠 **Hotel_Address** – Location is not relevant.
- 📊 **Additional_Number_of_Scoring** – Unrelated to individual reviews.
- 📅 **Review_Date** – Date does not influence sentiment.
- ⭐ **Average_Score** – A general rating, not review-specific.
- 🏨 **Hotel_Name** – Hotel name does not impact sentiment.
- 🌍 **Reviewer_Nationality** – Nationality is not crucial for sentiment.
- 🔢 **Total_Number_of_Reviews** – Does not affect individual sentiment.
- 🔄 **Total_Number_of_Reviews_Reviewer_Has_Given** – Unrelated to review sentiment.
- 🏷️ **Tags** – Useful for filtering, not sentiment.
- ⏳ **days_since_review** – Time does not affect sentiment directly.
- 📍 **lat, lng** – Location coordinates are irrelevant.

## ✅ Important Columns for Sentiment Analysis:

These columns are key for extracting sentiment:

- 📝 **Negative_Review** – Essential for detecting negative sentiment.
- 😊 **Positive_Review** – Essential for detecting positive sentiment.
- 🔢 **Review_Total_Negative_Word_Counts** – Indicates negativity intensity.
- 🔢 **Review_Total_Positive_Word_Counts** – Indicates positivity intensity.

- ⭐ **Reviewer_Score** – Can be used to validate sentiment classification.

## Overview

# 🧠 Tool: Sentiment Analysis with Python

## 📝 Sentiment Analysis?

Sentiment analysis is a crucial part of **Natural Language Processing (NLP)**. It involves **extracting emotions** from raw text to determine whether a statement conveys a **positive, negative, or neutral** sentiment.

This technique is widely used for:

- 🔻 **Social media monitoring** – Understanding public opinion.
- ⭐ **Customer reviews analysis** – Measuring customer satisfaction.
- 📊 **Brand reputation tracking** – Identifying trends in feedback.

The goal of this study is to demonstrate how sentiment analysis can be performed using **Python**.

---

## 🛠️ Libraries We Will Use:

For this task, we will leverage the following key libraries:

| 📦 Library | 🔍 Purpose |
| --- | --- |
| **NLTK** | The most popular Python library for NLP techniques. |
| | POS, tokenizer, lemmatizer, wordnet. |
| **Gensim** | A toolkit for topic modeling and vector space modeling. |
| **Scikit-learn** | The most widely used machine learning library in Python. |

---

# 🏨 Dataset: Hotel Reviews

We will analyze a dataset containing **hotel reviews** from various customers. Each observation consists of:

- **A textual review** – The customer's experience in their own words.
- **An overall rating** – A numerical score given by the customer.

📁 **Dataset Source:**
🔗 515K Hotel Reviews Data (Europe)

---

# 🎯 Problem Statement

For each customer review, our goal is to predict whether it is **positive (good experience)** or **negative (bad experience)** based only on the raw textual feedback.

## 📊 Sentiment Classification:

We will categorize reviews into **two sentiment classes** based on their overall ratings:

- 👎 **Negative Reviews** → Ratings **< 5**
- 👍 **Positive Reviews** → Ratings **≥ 5**

# 📥 Loading the Raw Data

Let's sentiment analysis, we first **load the raw dataset** containing hotel reviews.

## 🔍 Data Structure:

Each **textual review** is divided into two parts:

- **Positive Review** ✨ – Highlights what the customer liked.

- **Negative Review** ❌ – Mentions aspects that the customer disliked.

## 🛠️ Data Preprocessing:

To simplify our analysis, we will:

1. **Merge** both the positive and negative review sections.
2. **Create a single text column** containing the full customer feedback.
3. **Remove unnecessary metadata** to focus only on the raw text data.

By doing this, we ensure that our analysis is based purely on textual content without any additional numerical scores or labels.

In [1]:
```python
import pandas as pd
# Load data from CSV
main_df = pd.read_csv("Hotel_Reviews.csv")
# Concateneate neative and positive reviews.
main_df["review"] = main_df["Negative_Review"] + main_df["Positive_Review"]
# create the label
main_df["is_bad_review"] = main_df["Reviewer_Score"].apply(lambda x: 1 if x < 5 else 0)
# select only relevant columns
main_df = main_df[["review", "is_bad_review"]]
main_df.head()
```

Out[1]:

| | review | is_bad_review |
|---|---|---|
| 0 | I am so angry that i made this post available… | 1 |
| 1 | No Negative No real complaints the hotel was g… | 0 |
| 2 | Rooms are nice but for elderly a bit difficul… | 0 |
| 3 | My room was dirty and I was afraid to walk ba… | 1 |
| 4 | You When I booked with your company on line y… | 0 |

## Data set Rows, Columns Counts.

```
In [2]:    no_rows, no_cols = main_df.shape
           print(f"Rows Counts : {no_rows} , Columns Counts: {no_cols}")
```

Rows Counts : 515738 , Columns Counts: 2

## 🏨 Sample Data ~ 10% Of Dataset for reliable computation

```
In [4]:    sample_df = main_df.sample(frac = 0.1, replace = False, random_state=42)
           no_rows, no_cols = sample_df.shape
           print(f"Rows Counts : {no_rows} , Columns Counts: {no_cols}")
```

Rows Counts : 51574 , Columns Counts: 2

```
In [5]:    sample_df.isnull().sum()
```

```
Out[5]:    review            0
           is_bad_review     0
           dtype: int64
```

No Null values exist in sample data

## 🧹 Text Preprocessing: Cleaning the Data

### 🛑 Removing Default Placeholder Comments

In our dataset, customers who did not leave specific feedback have the following placeholders:

- **"No Negative"** → Appears when no negative feedback is provided.
- **"No Positive"** → Appears when no positive feedback is given.

These placeholders do not carry any real sentiment and **must be removed** from our text data.

---

### 🔧 Cleaning Text Data

Once we remove placeholder comments, the next step involves **text preprocessing** to ensure better accuracy in sentiment analysis.

We will perform the following operations:

- ✅ **Lowercasing** – Convert all text to lowercase for uniformity.
- ✅ **Removing Punctuation** – Eliminate special characters and symbols.
- ✅ **Tokenization** – Split sentences into individual words (tokens).
- ✅ **Stopword Removal** – Remove common words like "*the*", "*is*", "*and*" that do not add sentiment value.
- ✅ **Lemmatization/Stemming** – Reduce words to their base form (e.g., "*running*" → "*run*").

These steps are crucial for improving our **Natural Language Processing (NLP) model's performance** and ensuring more accurate sentiment predictions.

🚀 **Let's dive into the implementation!**

## Preprocess Text process :: for sample_df , all basic processing of NLP

```python
# return the wordnet values according to the POS tag
from nltk.corpus import wordnet

def fetch_wordnet_pos(pos_tag):
    if pos_tag.startswith('J'):
        return wordnet.ADJ
    elif pos_tag.startswith('V'):
        return wordnet.VERB
    elif pos_tag.startswith('N'):
        return wordnet.NOUN
    elif pos_tag.startswith('R'):
        return wordnet.ADV
    else:
        return wordnet.NOUN


import string
from nltk import pos_tag
from nltk.corpus import stopwords
from nltk.tokenize import WhitespaceTokenizer
from nltk.stem import WordNetLemmatizer
```

```python
def preprocess_text(text):
    # lower text
    text = text.lower()
    # tokenize text and remove puncutation
    text = [word.strip(string.punctuation) for word in text.split(" ")]
    # remove words that contain numbers
    text = [word for word in text if not any(c.isdigit() for c in word)]
    # remove stop words
    stop = stopwords.words('english')
    text = [x for x in text if x not in stop]
    # remove empty tokens
    text = [t for t in text if len(t) > 0]
    # find pos tag text
    pos_tags = pos_tag(text)
    # Lemmatization on a List of part-of-speech (POS) tagged words
    text = [WordNetLemmatizer().lemmatize(t[0], fetch_wordnet_pos(t[1])) for t in pos_tags]
    # remove words with only one letter
    text = [t for t in text if len(t) > 1]
    # join together
    text = " ".join(text)
    return(text)
```

In [8]:
```python
# remove 'No Negative' or 'No Positive' from text
sample_df["review"] = sample_df["review"].apply(lambda x: x.replace("No Negative", "").replace("No Positive", ""))
```

In [9]:
```python
### Preprocess sample data
sample_df["review_clean"] = sample_df["review"].apply(lambda x: preprocess_text(x))
```

# ✨ Feature engineering

## Adding Sentiment Analysis to Reviews

### 📌 Overview

To analyze customer feedback, we use **Sentiment Analysis**, which helps classify reviews as **positive, negative, or neutral** based on their text content.

In this step, we utilize **NLTK's** `SentimentIntensityAnalyzer` **(VADER)** to compute sentiment scores for each review.

In [10]:
```python
# Import the VADER Sentiment Analyzer from NLTK
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# Initialize Sentiment Intensity Analyzer
sid = SentimentIntensityAnalyzer()

# Apply Sentiment Analysis to each review text
sample_df["sentiments"] = sample_df["review"].apply(lambda x: sid.polarity_scores(x))

# Expand sentiment dictionary into separate columns and merge with the main DataFrame
sample_df = pd.concat([sample_df.drop(['sentiments'], axis=1), sample_df['sentiments'].apply(pd.Series)], axis=1)
```

In [11]:
```python
reviews_df = sample_df.copy()
```

## Creating Doc2Vec Vector Columns for Text Data

In [14]:
```python
## Creating Doc2Vec Vector Columns for Text Data
from gensim.test.utils import common_texts
## This code transforms text data into numerical vectors using Doc2Vec from the gensim library.
from gensim.models.doc2vec import Doc2Vec, TaggedDocument
## Convert Text Data into TaggedDocument Format
documents = [TaggedDocument(doc, [i]) for i, doc in enumerate(reviews_df["review_clean"].apply(lambda x: x.split(" ")
# train a Doc2Vec model
model = Doc2Vec(documents, vector_size=5, window=2, min_count=1, workers=4)
### Convert Each Document into a Vector
doc2vec_df = reviews_df["review_clean"].apply(lambda x: model.infer_vector(x.split(" "))).apply(pd.Series)
## Merge the Vector Data Back into the Original DataFrame
doc2vec_df.columns = ["doc2vec_vector_" + str(x) for x in doc2vec_df.columns]
```

## Doc2Vec Generated Vectors Table

| doc2vec_vector_0 | doc2vec_vector_1 | doc2vec_vector_2 | doc2vec_vector_3 | doc2vec_vector_4 |
|---|---|---|---|---|
| 0.12 | -0.45 | 0.78 | 0.34 | -0.23 |

| -0.21 | 0.67 | -0.88 | 0.13 | 0.45 |

## Tagged Documents Representation

```
[
    TaggedDocument(words=['great', 'product', 'love', 'it'], tags=[0]),
    TaggedDocument(words=['poor', 'quality', 'waste', 'money'], tags=[1])
]
```

## **Example - What we are getting now:

This **merges** the original DataFrame ( `reviews_df` ) with the new DataFrame ( `doc2vec_df` ) containing **vectorized text representations**.

---

### 1 Original `reviews_df`

| review_clean |
| --- |
| "great product love it" |
| "poor quality waste money" |

### 2 `doc2vec_df` (Generated Vectors)

| doc2vec_vector_0 | doc2vec_vector_1 | doc2vec_vector_2 | doc2vec_vector_3 | doc2vec_vector_4 |
| --- | --- | --- | --- | --- |
| 0.12 | -0.45 | 0.78 | 0.34 | -0.23 |
| -0.21 | 0.67 | -0.88 | 0.13 | 0.45 |

### 3 After `pd.concat()` (Final `reviews_df` )

| review_clean | doc2vec_vector_0 | doc2vec_vector_1 | doc2vec_vector_2 | doc2vec_vector_3 | doc2vec_vector_4 |
| --- | --- | --- | --- | --- | --- |
| "great product love it" | 0.12 | -0.45 | 0.78 | 0.34 | -0.23 |

| review_clean | doc2vec_vector_0 | doc2vec_vector_1 | doc2vec_vector_2 | doc2vec_vector_3 | doc2vec_vector_4 |
|---|---|---|---|---|---|
| "poor quality waste money" | -0.21 | 0.67 | -0.88 | 0.13 | 0.45 |

Now, each text review has **5 numerical vector values** representing its semantic meaning.

---

☑ Merges the **original text data** with the **generated document vectors**
☑ Allows **further processing**, such as training Machine Learning models
☑ Keeps the **original review data intact** while adding meaningful features

## Extracts TF-IDF (Term Frequency-Inverse Document Frequency) features from the review_clean column.

In [ ]:

In [15]:
```python
from sklearn.feature_extraction.text import TfidfVectorizer
import pandas as pd

# Sample DataFrame
###reviews_df = pd.DataFrame({"review_clean": ["This is a sample review", "Another review text", "More text here"]})

# Apply TF-IDF
tfidf = TfidfVectorizer(min_df=10)
tfidf_result = tfidf.fit_transform(reviews_df["review_clean"]).toarray()

# Create DataFrame with TF-IDF values
tfidf_df = pd.DataFrame(tfidf_result, columns=tfidf.get_feature_names_out())
tfidf_df.columns = ["word_" + str(x) for x in tfidf_df.columns]
tfidf_df.index = reviews_df.index

# Concatenate with original DataFrame
reviews_df = pd.concat([reviews_df, tfidf_df], axis=1)
```

## The TF-IDF metric solves this problem:

- TF computes the classic number of times the word appears in the text
- IDF computes the relative importance of this word which depends on how many texts the word can be found

## ◆ Example: Before & After TF-IDF Transformation

### 📝 Original `reviews_df`

| review_clean |
| --- |
| "great product love it" |
| "poor quality waste money" |

### 📊 After Adding TF-IDF Features

| review_clean | word_great | word_product | word_love | word_quality | word_money |
| --- | --- | --- | --- | --- | --- |
| "great product love it" | 0.45 | 0.38 | 0.52 | 0.00 | 0.00 |
| "poor quality waste money" | 0.00 | 0.00 | 0.00 | 0.58 | 0.47 |

Each review is now represented by **TF-IDF scores**, making it suitable for **Machine Learning tasks** such as sentiment analysis and classification.

```
In [32]: reviews_df.head()
```

Out[32]:

| | review | is_bad_review | review_clean | neg | neu | pos | compound | doc2vec_vector_0 | doc2vec_vector_1 | doc2vec_v |
|---|---|---|---|---|---|---|---|---|---|---|
| 488440 | Would have appreciated a shop in the hotel th... | 0 | would appreciate shop hotel sell drinking wate... | 0.049 | 0.617 | 0.334 | 0.9924 | 0.057766 | 0.483825 | -( |
| 274649 | No tissue paper box was present at the room | 0 | tissue paper box present room | 0.216 | 0.784 | 0.000 | -0.2960 | 0.072940 | -0.038417 | -( |
| 374688 | Pillows Nice welcoming and service | 0 | pillow nice welcome service | 0.000 | 0.345 | 0.655 | 0.6908 | -0.093065 | 0.059668 | ( |
| 404352 | Everything including the nice upgrade The Hot... | 0 | everything include nice upgrade hotel revamp s... | 0.000 | 0.621 | 0.379 | 0.9153 | -0.124238 | 0.099710 | ( |
| 451596 | Lovely hotel v welcoming staff | 0 | lovely hotel welcome staff | 0.000 | 0.230 | 0.770 | 0.7717 | 0.089844 | 0.108858 | ( |

In [16]: `reviews_df.shape`

Out[16]: (51574, 3833)

# Exploratory data analysis

Let's try understading our data, what it is contains and how them concanating.

```
In [17]:   # show is_bad_review distribution
           reviews_df["is_bad_review"].value_counts(normalize = True)
```

```
Out[17]:   is_bad_review
           0    0.956761
           1    0.043239
           Name: proportion, dtype: float64
```

Our dataset is highly imbalanced because less than 5% of our reviews are considered as negative ones. This information will be very useful for the modelling part.

## Wordclouds ==> For Reviews:

```
In [18]:   # wordcloud presentation
           # total_chars ==> total characters in the  characters column
           reviews_df["total_chars"] = reviews_df["review"].apply(lambda x: len(x))

           ### total words in the words column
           reviews_df["total_words"] = reviews_df["review"].apply(lambda x: len(x.split(" ")))

           from wordcloud import WordCloud
           import matplotlib.pyplot as plt

           def show_wordcloud(data, title = None):
               wordcloud = WordCloud(
                   background_color = 'white',
                   max_words = 200,
                   max_font_size = 40,
                   scale = 3,
                   random_state = 42
               ).generate(str(data))

               fig = plt.figure(1, figsize = (20,20))
               plt.axis('off')
               if title:
                   fig.suptitle(title, fontsize = 15)
                   fig.subplots_adjust(top = 2.3)

               plt.imshow(wordcloud)
```

```
    plt.show()

# print wordcloud
show_wordcloud(reviews_df["review"])
```



```
In [19]:  # highest positive sentiment reviews (with more than 5 words)
          reviews_df[reviews_df["total_words"] >= 5].sort_values("pos", ascending = False)[["review", "pos"]].head(10)
```

|  | review | pos |
|---|---|---|
| **43101** | A perfect location comfortable great value | 0.931 |
| **211742** | Clean comfortable lovely staff | 0.907 |
| **175551** | Friendly welcome Comfortable room | 0.905 |
| **365085** | Good location great value | 0.904 |
| **109564** | Clean friendly and comfortable | 0.902 |
| **145743** | Good value amazing location | 0.901 |
| **407590** | breakfast excellent Clean comfort | 0.899 |
| **407546** | Great place I enjoyed | 0.881 |
| **218571** | Beautiful Quirky Comfortable | 0.878 |
| **436901** | Lovely comfortable rooms | 0.877 |

- **Key Finding:** The most positive reviews correspond to positive customer feedback.
- **Implication:** This suggests a direct link between review sentiment and customer satisfaction.

```python
# lowest negative sentiment reviews (with more than 5 words)
reviews_df[reviews_df["total_words"] >= 5].sort_values("neg", ascending = False)[["review", "neg"]].head(10)
```

Out[20]:

|  | review | neg |
| --- | --- | --- |
| **193086** | No dislikes LOCATION | 0.831 |
| **356368** | Nothing Great helpful wonderful staff | 0.812 |
| **318516** | A disaster Nothing | 0.804 |
| **458794** | Nothing Excellent friendly helpful staff | 0.799 |
| **29666** | A bit noisy No | 0.796 |
| **426057** | Dirty hotel Smells bad | 0.762 |
| **263187** | Very bad service No | 0.758 |
| **443796** | Nothing perfect | 0.750 |
| **181508** | Window blind was broken | 0.744 |
| **175316** | Nothing Super friendly staff | 0.743 |

- **VADER's Limitations:**
  - Phrases like "no problems" or "nothing wrong" can be misclassified.
- **Overall Trend:**
  - Most negative reviews accurately reflect negative customer sentiment.

In [21]:
```python
# plot sentiment distribution for positive and negative reviews

import seaborn as sns

for x in [0, 1]:
    subset = reviews_df[reviews_df['is_bad_review'] == x]

    # Draw the density plot
    if x == 0:
        label = "Good reviews"
    else:
        label = "Bad reviews"
    sns.kdeplot(subset['compound'], label=label)
```

## 📊 Sentiment Distribution Analysis

The above graph illustrates the **distribution of sentiment scores** among **good and bad reviews**.

- **Good reviews** are mostly classified as **very positive** by the Vader sentiment analysis tool.
- **Bad reviews** tend to have **lower compound sentiment scores**, indicating more negative sentiment.

◆ **Key Insight:**

This analysis confirms that the **computed sentiment features** will play a crucial role in the **modeling phase**, helping to distinguish between positive and negative reviews effectively. 🚀

## 📌 Feature Selection & Train-Test Split

Let's performs **feature selection and dataset splitting** for training a machine learning model.

- ◆ **Steps Explained:**

  1. 🔍 **Feature Selection:**

     - The **target label** ( `is_bad_review` ), **raw text** ( `review` ), and **cleaned text** ( `review_clean` ) are **excluded** from the feature set.
     - Remaining columns in `reviews_df` are selected as **features** for model training.
  2. 📊 **Train-Test Split:**

     - The dataset is **split into training (80%) and testing (20%) subsets**.
     - `train_test_split()` ensures **randomized sampling** with a fixed seed ( `random_state=42` ) for reproducibility.

🚀 **Purpose:**

This step **prepares the dataset** for training a **Random Forest Classifier** by selecting meaningful features and ensuring a proper training/testing distribution.

```python
In [22]:   # feature selection
           label = "is_bad_review"
           ignore_cols = [label, "review", "review_clean"]
           features = [c for c in reviews_df.columns if c not in ignore_cols]

           # split the data into train and test
           from sklearn.ensemble import RandomForestClassifier
           from sklearn.model_selection import train_test_split

           X_train, X_test, y_train, y_test = train_test_split(reviews_df[features], reviews_df[label], test_size = 0.20, randon
```

## Model Training for Classification & Duplicate Column Handling

- 🔍 **Identify Duplicates:** Scrutinize column names to detect any duplications.
- 📝 **Rename Duplicates:** Apply a systematic method to rename duplicate columns, ensuring uniqueness.
- 🚀 **Prepare Data:** Conduct this preprocessing *before* splitting the data for consistency.
- 📊 **Guarantee Compatibility:** Ensure data is in a clean, compatible format for optimal classifier performance.
- ✂️ **Split Dataset:** Divide the preprocessed dataset into training and testing subsets using `train_test_split`.
- 🌲 **Train RandomForest:** Train a `RandomForestClassifier` on the training data, utilizing the features identified.
- 📈 **Feature Importance:** Extract and analyze feature importance scores from the trained RandomForest model to understand feature relevance.

- 📋 **Evaluate Model:** Evaluate model performance using appropriate metrics (e.g., AUC-ROC, Precision-Recall) on the test dataset.

In [23]:
```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier

duplicate_columns = reviews_df.columns[reviews_df.columns.duplicated()]

if len(duplicate_columns) > 0:
    print("Duplicate columns found: ", duplicate_columns)
    new_columns = []
    count_dict = {}
    for col in reviews_df.columns:
        if col in count_dict:
            count_dict[col] += 1
            new_columns.append(f"{col}_{count_dict[col]}")
        else:
            count_dict[col] = 0
            new_columns.append(col)
    reviews_df.columns = new_columns

# feature selection
label = "is_bad_review"
ignore_cols = [label, "review", "review_clean"]
features = [c for c in reviews_df.columns if c not in ignore_cols]

# split the data into train and test
X_train, X_test, y_train, y_test = train_test_split(
    reviews_df[features], reviews_df[label], test_size=0.20, random_state=42
)

# Train the Random Forest Classifier
rf = RandomForestClassifier(random_state=42)
rf.fit(X_train, y_train)

# Check lengths and create the DataFrame
print(f"Number of features: {len(X_train.columns)}")
print(f"Number of feature importances: {len(rf.feature_importances_)}")

if len(X_train.columns) == len(rf.feature_importances_):
```

```
    feature_importances_df = pd.DataFrame(
        {"feature": X_train.columns, "importance": rf.feature_importances_}
    ).sort_values("importance", ascending=False)
    print(feature_importances_df.head(20))
else:
    print("Error: Feature list and feature importances have different lengths.")
    print("Debugging information:")
    print("X_train columns:", X_train.columns)
```

```
Number of features: 3832
Number of feature importances: 3832
            feature  importance
3          compound    0.044380
2               pos    0.025604
0               neg    0.024508
3830     total_chars   0.023911
3831     total_words   0.018288
1               neu    0.017590
2846      word_room    0.011517
2232   word_nothing    0.010200
277         word_bad    0.009894
942       word_dirty    0.008723
1937   word_location   0.008266
3195      word_staff    0.007653
3209       word_star    0.007624
1631      word_hotel    0.007175
2277        word_old    0.006407
2196      word_never    0.006181
3081      word_small    0.005772
2510       word_poor    0.005631
424    word_breakfast   0.005502
2860       word_rude    0.005249
```

In [ ]:

## Receiver operating characteristic example
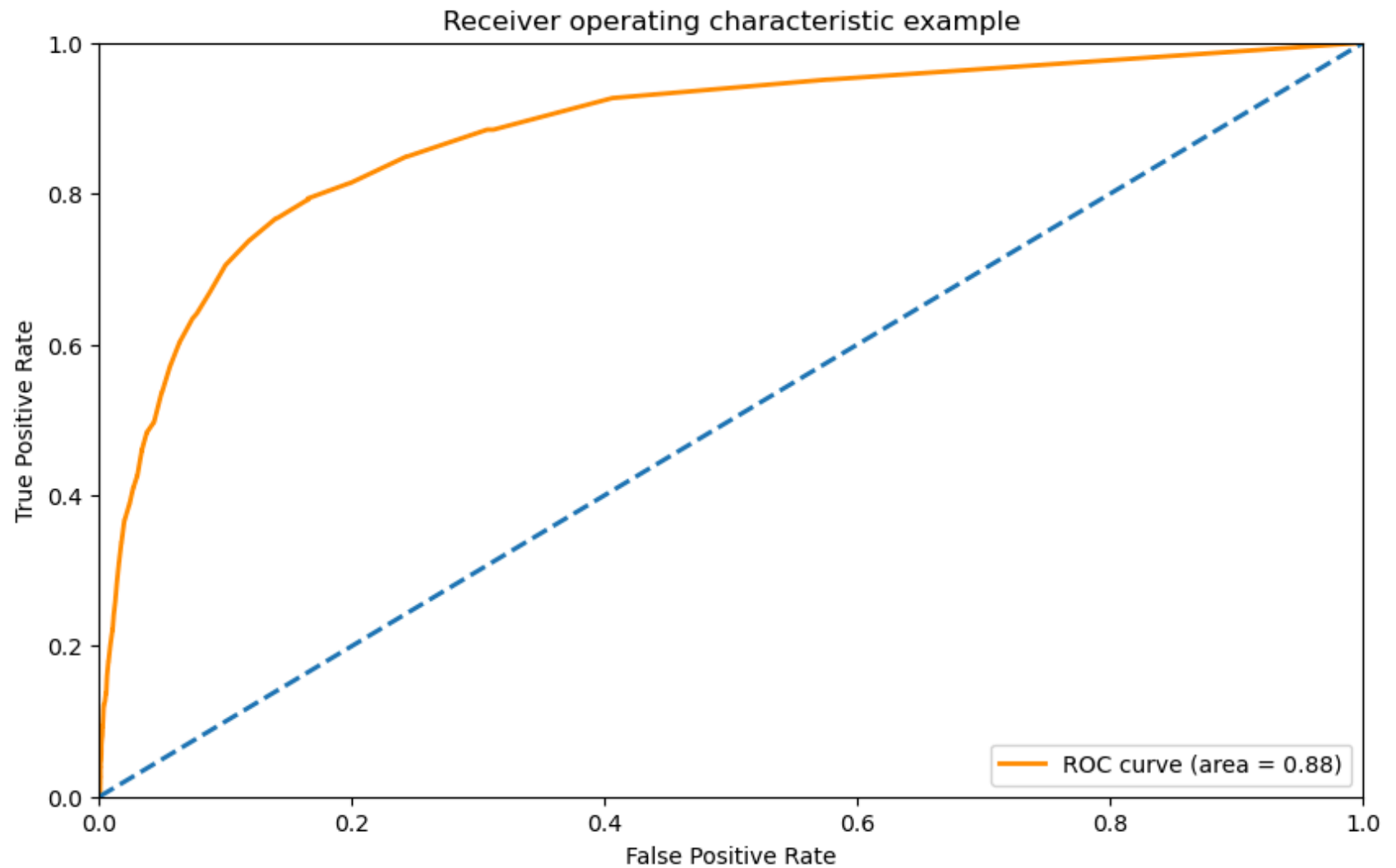
In [24]:
```python
# ROC curve

from sklearn.metrics import roc_curve, auc, roc_auc_score
import matplotlib.pyplot as plt
```

```python
y_pred = [x[1] for x in rf.predict_proba(X_test)]
fpr, tpr, thresholds = roc_curve(y_test, y_pred, pos_label = 1)

roc_auc = auc(fpr, tpr)

plt.figure(1, figsize = (10, 6))
lw = 2
plt.plot(fpr, tpr, color='darkorange',
         lw=lw, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], lw=lw, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic example')
plt.legend(loc="lower right")
plt.show()
```

Receiver operating characteristic example

- **ROC/AUC Evaluation: 0.88**
  - The Receiver Operating Characteristic (ROC) curve and Area Under the Curve (AUC) are standard metrics for binary classifier evaluation.
  - An AUC-ROC of 0.88 indicates strong discriminatory power.
- **Class Imbalance Impact:**
  - Significant class imbalance (high prevalence of negative instances) compromises the ROC curve's reliability in this scenario.

- **False Positive Rate (FPR) Suppression:**
  - The False Positive Rate (FPR = False Positives / Negatives) is suppressed due to the large number of negative instances.
  - This suppression allows the model to generate numerous false positives while maintaining a low FPR.
- **Artificial AUC Inflation:**
  - The artificially low FPR leads to an inflated AUC-ROC, misrepresenting the model's true performance.
- **Alternative Metrics Recommendation:**
  - Alternative evaluation metrics, robust to class imbalance, are recommended for accurate model assessment.

In [ ]:

## Precision-Recall (PR) Curve Analysis

- 📈 **Performance Visualization:** The PR curve visually represents a classifier's performance, especially valuable in imbalanced datasets.
- 🎯 **Precision Focus:** Measures the proportion of correctly predicted positive cases out of all predicted positives (minimizes false positives).
- Recall Focus: Measures the proportion of correctly predicted positive cases out of all actual positives (minimizes false negatives).
- ⚖️ **Imbalanced Data Relevance:** Provides a more informative assessment than ROC curves when dealing with significant class imbalances.
- 📊 **Average Precision (AP):** Summarizes the PR curve into a single value, indicating the average precision across different recall thresholds.

In [ ]:

### *Visualize 2-Class Precision-Reall Curve*

In [27]:
```python
import matplotlib.pyplot as plt
from sklearn.metrics import average_precision_score, precision_recall_curve

# Compute average precision score
average_precision = average_precision_score(y_test, y_pred)

# Compute Precision-Recall curve
```
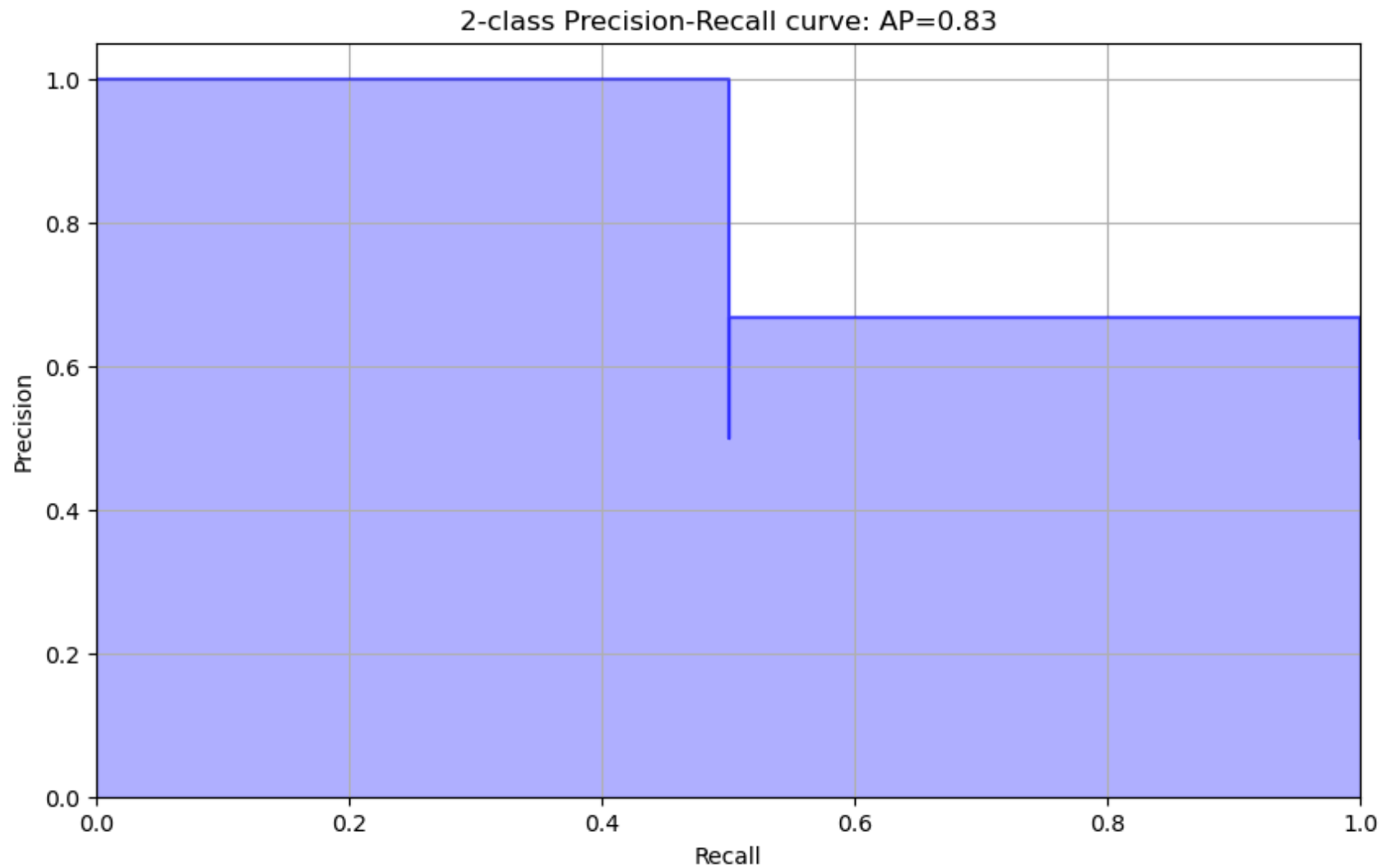
```python
precision, recall, _ = precision_recall_curve(y_test, y_pred)

# Plot Precision-Recall curve
plt.figure(figsize=(10, 6))
plt.step(recall, precision, color='b', alpha=0.6, where='post')
plt.fill_between(recall, precision, alpha=0.3, color='b', step='post')

plt.xlabel('Recall')
plt.ylabel('Precision')
plt.ylim([0.0, 1.05])
plt.xlim([0.0, 1.0])
plt.title(f'2-class Precision-Recall curve: AP={average_precision:.2f}')
plt.grid(True)
plt.show()
```
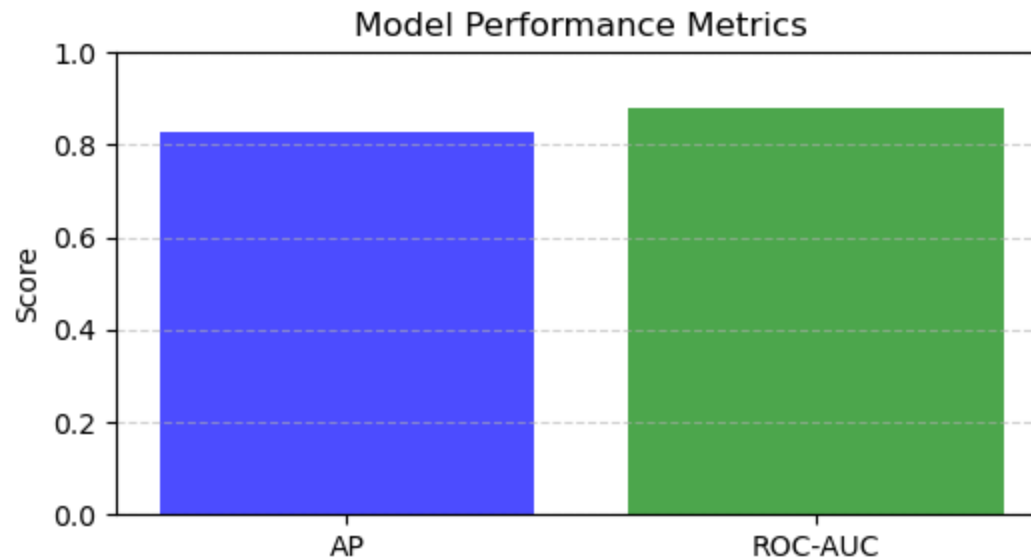
2-class Precision-Recall curve: AP=0.83

## Model Performance Metrics

```python
In [28]: import matplotlib.pyplot as plt

# Create a Visual Representation of AP and ROC-AUC
fig, ax = plt.subplots(figsize=(6, 3))
ax.bar(['AP', 'ROC-AUC'], [0.83, 0.88], color=['blue', 'green'], alpha=0.7)
ax.set_ylim([0, 1])
```

```
ax.set_ylabel('Score')
ax.set_title('Model Performance Metrics')
plt.grid(axis='y', linestyle='--', alpha=0.6)
plt.show()
```



Model Performance Metrics

# 📊 Model Performance Evaluation

## 🚀 Key Metrics

### ◆ Precision-Recall Score (PR) = 0.83

- PR summarizes the Precision-Recall curve.
- A higher PR score indicates better performance in handling class imbalance.

### ◆ ROC-AUC = 0.88

- Measures how well the model distinguishes between classes.
- A value of 0.88 suggests strong classification ability.

---

## 📈 Performance Interpretation

- With **PR = 0.83**, the model balances **Precision** and **Recall** effectively.
- A **ROC-AUC of 0.88** indicates strong **discrimination power** between classes.
- Generally, a PR score > **0.80** and ROC-AUC > **0.85** represent a well-performing model.

---

## ✅ Conclusion

- This model demonstrates **high accuracy and reliability**.
- It is a **strong candidate for deployment** in real-world scenarios.
- Further tuning may improve results, but current performance is already **robust**.

---

## 🎯 Final Verdict: ⭐ ⭐ ⭐ ⭐ ⭐ Good Model!

In [ ]: