

Objective:

The goal of this code is to build a hotel recommendations system using Natural Language Processing (NLP) techniques. The system is based on the similarity of hotel reviews, with the assumption that hotels with similar reviews are likely to be preferred by the same type of guests.

1. Libraries Used:

pandas: For data manipulation and analysis.

numpy: For numerical operations.

nltk: For natural language processing tasks like tokenization and stop word removal.

sklearn: For various machine learning tasks such as text vectorization, cosine similarity, and k-Nearest Neighbors.

matplotlib: For data visualization.

2. Data Loading and Preprocessing:

The dataset is loaded using Pandas from a CSV file.

Unnecessary columns are dropped, and positive and negative reviews are combined into a single column.

NLTK is used for text preprocessing, including tokenization and stop word removal.

The distribution of reviewer scores is visualized using a histogram.

3. Text Vectorization:

The combined reviews are transformed into a TF-IDF (Term Frequency-Inverse Document Frequency) matrix using the TfidfVectorizer from scikit-learn.

4. k-Nearest Neighbors (k-NN) Model:

A k-NN model is created using the cosine similarity metric, as it is suitable for text data.

The model is trained on the TF-IDF matrix.

5. Recommender System Function:

A function `get_recommendations` is defined to provide personalized hotel recommendations based on the input hotel name.

The function calculates the cosine similarity between the input hotel and others, excluding the queried hotel itself.

It returns a DataFrame with similar hotels and their reviewer scores.

6. Model Evaluation:

For evaluation, a hypothetical query hotel is selected ('Hotel Arena').

Recommendations for this hotel are obtained and displayed.
Mean Squared Error (MSE) is calculated by comparing actual and predicted reviewer scores for each hotel. The choice of MSE is motivated by the nature of the problem and the goal of assessing the accuracy of the predicted reviewer scores.

7. Visualization:

The distribution of reviewer scores is visualized using a histogram to understand the overall pattern.

8. Comments:

Inline comments are provided to explain key steps and decisions throughout the code.
Variable names are chosen to be descriptive, aiding in code readability.

9. Results:

Recommendations for the hypothetical query hotel are printed.
The Mean Squared Error (MSE) between actual and predicted scores is displayed, providing a measure of model performance