**Mining Massive Datasets: Hotel Recommendation System**

**Project Report**

**Group Members:**

Anzla Aslam                    20L-2109

Shahbano Waqar                 20L-2162

National University of Computer and Emerging Sciences
Department of Computer Science
Lahore, Pakistan

## 1. Problem statement

The hospitality industry faces the challenge of efficiently recommending hotels to customers based on their preferences and past experiences. Traditional recommendation systems often lack personalization and fail to capture the nuances of individual preferences, leading to suboptimal user experiences. The goal of this project is to develop a hotel recommendation system that leverages machine learning techniques to provide personalized recommendations to users.

## 2. Solution

The proposed solution involves the development of a hotel recommendation system using PySpark, a powerful framework for large-scale data processing. The system integrates various machine learning algorithms and natural language processing techniques to analyze user reviews, extract meaningful features, and generate personalized recommendations.

## 3. Implementation

- Sentiment Analysis: TextBlob library is used for sentiment analysis to categorize user reviews into positive, neutral, or negative sentiments.
- K-means Clustering: PySpark's K-means clustering algorithm is applied to group hotels based on their features extracted from reviews and ratings.
- ALS for Recommendation: Alternating Least Squares (ALS) algorithm is employed for collaborative filtering-based recommendation, utilizing user ratings to generate personalized hotel recommendations.
- Text Processing: Text data from user reviews is processed using a pipeline consisting of tokenization, hashing term frequency (TF), and inverse document frequency (IDF) to extract meaningful features.
- Pairwise Similarity: Cosine similarity is calculated between hotel feature vectors to identify similar hotels and generate recommendations.

## 4. Problems in implementation

- Scalability: Depending on the size of the dataset, the scalability of the implemented algorithms might be a concern. Further optimization may be needed for large-scale deployment.
- Model Evaluation: While cross-validation is used for hyperparameter tuning, additional evaluation metrics and techniques could be explored to ensure the robustness and accuracy of the recommendation model.

## 5. Solution

Scalability:

- Partitioning and Caching: Spark's DataFrame API inherently handles data partitioning and caching. However, you can optimize partitioning strategies based on the data distribution and cluster configuration to improve parallelism and reduce shuffle operations during computations.
- Optimized Algorithms: PySpark leverages distributed computing, allowing algorithms like ALS for recommendation and K-means clustering to scale horizontally across multiple nodes in a Spark cluster. By utilizing PySpark's MLlib library, which is built for scalable machine learning, the implemented algorithms inherently handle large-scale datasets efficiently.
- Cluster Sizing and Resource Management: You can adjust the Spark cluster configuration, such as the number of worker nodes, executor memory, and executor cores, to optimize resource utilization based on workload characteristics and available hardware resources.

Model Evaluation:

- Additional Evaluation Metrics: While the code snippet primarily focuses on RMSE as the evaluation metric for the ALS recommendation model, you can extend the evaluation by calculating additional metrics such as precision, recall, or AUC for the recommendation system. These metrics provide deeper insights into model performance and user satisfaction.
- Online Evaluation: The provided code doesn't include real-time or online evaluation techniques. However, you can implement feedback loops and online evaluation pipelines to continuously monitor model performance, collect user feedback, and adapt the recommendation model dynamically.
- A/B Testing: A/B testing is not explicitly implemented in the provided code. However, you can conduct A/B tests by deploying multiple versions of the recommendation model with different configurations (e.g., hyperparameters) and measuring key metrics (e.g., user engagement) to identify the most effective model variant.

In summary, while the provided code demonstrates the implementation of key components for building a hotel recommendation system using PySpark, additional steps can be taken to further optimize scalability, enhance model evaluation, and ensure the effectiveness of the recommendation system in real-world scenarios.

## 6. Findings

Root-mean-square Error (RMSE):

- Before tuning: 1.6397937826251077
- After tuning: 1.2597498087540997

    The decrease in RMSE after parameter tuning indicates improved model performance.

TF-IDF Features:

Sample TF-IDF features are displayed, demonstrating the transformation of text data into numerical features.

Pairwise Similarity:

Pairwise similarity between hotels is calculated, showing the similarity scores between different hotels based on their features extracted from reviews.

Overall, the findings suggest that the implemented recommendation system shows promising results in generating personalized hotel recommendations based on user preferences and past experiences. Further refinement and optimization could enhance its effectiveness and scalability for real-world deployment.