

# Online Grocery Recommendation System

## Final Project Report

### Team Members :

[ Zeina Ayman 202200351 - Habiba Khalil 202200720 - Nour Helmy 202202012 - Youssef Abohendia 202202699 - Youssef Fathy 202202193 ]

### Introduction :

This project presents a comprehensive analysis of a grocery recommendation system based on users' past purchases. The primary goal is to predict whether a user will reorder a particular product in their next purchase. Three machine learning models were evaluated: K-Nearest Neighbors (KNN), Decision Tree, and Logistic Regression.

### Methodology :

#### 1. Data Preparation :

- **Dataset :** The merged dataset consists of user orders and product information with the following columns : [ order\_id, product\_id, add\_to\_cart\_order, reordered, product\_name, aisle\_id, department\_id, aisle, department, user\_id, eval\_set, order\_number, order\_dow, order\_hour\_of\_day, days\_since\_prior\_order, user\_total\_orders, user\_total\_products, user\_reorder\_ratio, product\_total\_orders, product\_reorder\_ratio ].
- **Features :** Selected features for model training : [ user\_total\_orders, user\_total\_products, user\_reorder\_ratio, product\_total\_orders, product\_reorder\_ratio, add\_to\_cart\_order, order\_dow, order\_hour\_of\_day, days\_since\_prior\_order ].
- **Target :** reordered (binary variable indicating whether a product was reordered).

#### 2. Data Splitting :

- **Train -Test Split :** The data was split into training and testing sets with an 80-20 ratio using train\_test\_split from Scikit-learn.

### Model Training and Evaluation :

#### 1. K-Nearest Neighbors (KNN) :

n\_neighbors = 3

Suitable for small to medium-sized datasets.

## 2. Decision Tree :

random\_state = 42

Easy to interpret but can overfit.

## 3. Logistic Regression :

random\_state = 42

max\_iter = 1000

### Evaluation Metrics :

- **Accuracy** : Proportion of correct predictions.
- **Precision** : Proportion of true positive predictions among all positive predictions.
- **Recall** : Proportion of true positive predictions among all actual positives.
- **AUC-ROC** : Area under the ROC curve, indicating the model's ability to distinguish between classes.

### Results :

#### 1. K-Nearest Neighbors (KNN) :

- **Accuracy** : 0.66
- **Precision** : 0.70
- **Recall** : 0.76
- **AUC-ROC** : 0.69

#### 2. Decision Tree :

- **Accuracy** : 0.66
- **Precision** : 0.71
- **Recall** : 0.71
- **AUC-ROC** : 0.65

#### 3. Logistic Regression :

- **Accuracy** : 0.66
- **Precision** : 0.70
- **Recall** : 0.75
- **AUC-ROC** : 0.70

**Model Selection :**

Based on the evaluation metrics, the Logistic Regression model demonstrates the best overall performance with balanced accuracy, precision, recall, and AUC-ROC scores. It is also efficient in terms of prediction speed.

**Conclusion :**

The Logistic Regression model is recommended for predicting grocery reorders based on users' past purchases due to its high performance and efficiency. Future work could involve further optimization and exploring additional features to improve the model's predictive power.