# Data Collection and Preprocessing Phase

| Date | June 2024 |
|---|---|
| Team ID | 739663 |
| Project Title | Estimating the stock keeping units using Machine Learning |
| Maximum Marks | 6 Marks |

**Preparation Template**

The images will be preprocessed by resizing, normalizing, augmenting, denoising, adjusting contrast, detecting edges, converting color space, cropping, batch normalizing, and whitening data. These steps will enhance data quality, promote model generalization, and improve convergence during neural network training, ensuring robust and efficient performance across various computer vision tasks.

| Section | Description |
|---|---|
| Data Overview | There are many popular open sources for collecting the data. Eg: kaggle.com, UCI repository, etc. In this project we have used .csv data. |
| Data Preparation | These are the general steps of pre-processing the data before using it for machine learning |
| Handling missing values | We use Handling missing values For checking the null values |
| Handling categorical data | As we can see our dataset has categorical data we must convert the categorical data to integer encoding or binary encoding |
| Handling Outliers in Data | With the help of boxplot, outliers are visualized. And here we are going to find upper bound and lower bound of numerical features with some mathematical formula. |

# Data Preparation

| | |
|---|---|
| Collect the dataset | Please refer to the link given below to download the dataset.<br>Link: https://www.kaggle.com/datasets/prachi13/customer-analytics?select=Train.csv |
| Importing the libraries | ```python<br>import pandas as pd<br>import numpy as np<br>import seaborn as sns<br>import matplotlib.pyplot as plt<br>from imblearn.over_sampling import SMOTE<br>from sklearn.model_selection import train_test_split<br>from sklearn.metrics import classification_report<br>from sklearn.metrics import confusion_matrix, precision_sc<br>from sklearn.ensemble import RandomForestClassifier<br>import warnings<br>warnings.filterwarnings("ignore")<br>``` |
| Loading Data | We use the code<br>`df=pd.read_csv("/content/Train.csv")`<br>For reading the dataset |
| Handling missing values | ```<br>df.isnull().sum()<br><br>ID                    0<br>Warehouse_block       0<br>Mode_of_Shipment      0<br>Customer_care_calls   0<br>Customer_rating       0<br>Cost_of_the_Product   0<br>Prior_purchases       0<br>Product_importance    0<br>Gender                0<br>Discount_offered      0<br>Weight_in_gms         0<br>Reached.on.Time_Y.N   0<br>dtype: int64<br>``` |
| Handling Categorical values | ```python<br>label_map={}<br>for i in df.columns:<br>  if str(df[i].dtype) == 'object':<br>    temp={}<br>    cats=df[i].unique()<br>    for index in range(len(cats)):<br>      temp[cats [index]]=index<br>    label_map[i]=temp<br>    #Labeling<br>    df[i]=df[i].map(temp)<br><br>print(label_map)<br><br>{'Warehouse_block': {'D': 0, 'F': 1, 'A': 2, 'B': 3, 'C': 4}<br>``` |

| Handling Outliers | ```
c=0
plt.figure(figsize=(18, 10))
for i in df.drop(columns=['Warehouse_block', 'Mode_of_Shipment','Product_importan
  if str(df[i].dtype)=='object':
    continue
  plt.subplot(2, 3, c+1)
  plt.boxplot(df[i])
  plt.title(i)
  c+=1
plt.show()
```  |