

Data Collection and Preprocessing Phase

Date	June 2024
Team ID	739971
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. Link: https://www.kaggle.com/datasets/aswathrao/demand-forecasting
Importing the libraries	<pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns from sklearn.model_selection import train_test_split from sklearn.linear_model import LinearRegression from sklearn.metrics import mean_squared_error, r2_score from sklearn.ensemble import RandomForestRegressor from sklearn.tree import DecisionTreeRegressor from sklearn.metrics import accuracy_score</pre>

Loading Data	<p>We use the code</p> <pre>df=pd.read_csv('/content/train_0irEZ2H.csv')</pre> <p>For reading the dataset</p>
Handling missing values	<pre>df.isnull().sum()</pre> <pre>record_ID 0 week 0 store_id 0 sku_id 0 total_price 1 base_price 0 is_featured_sku 0 is_display_sku 0 units_sold 0 dtype: int64</pre>
Finding Categorical values and numerical values	<pre>categorical_cols=df.select_dtypes(include=['object']) print(categorical_cols)</pre> <pre>Index(['week'], dtype='object')</pre> <pre>numerical_cols=df.select_dtypes(include=['int64'], print(numerical_cols)</pre> <pre>Index(['record_ID', 'store_id', 'sku_id', 'total_ 'is_featured_sku', 'is_display_sku', 'unit dtype='object')</pre>

Handling Outliers



