# Data Collection and Preprocessing Phase

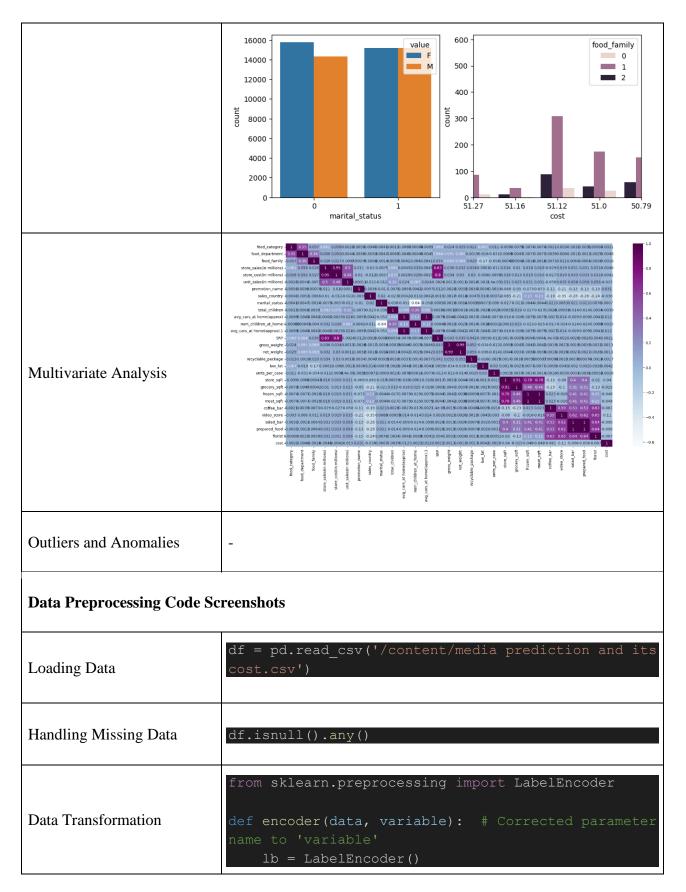| | |
|---|---|
| Date | 4th July 2024 |
| Team ID | 739804 |
| Project Title | Cost Prediction of Acquiring a Customer |
| Maximum Marks | 6 Marks |

## Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview | Dimensions: (60428, 40) |
| Univariate Analysis |  |
| Bivariate Analysis |  |

| | |
|---|---|
| |  |
| Multivariate Analysis |  |
| Outliers and Anomalies | - |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data | ```python
df = pd.read_csv('/content/media prediction and its cost.csv')
``` |
| Handling Missing Data | ```python
df.isnull().any()
``` |
| Data Transformation | ```python
from sklearn.preprocessing import LabelEncoder

def encoder(data, variable):  # Corrected parameter name to 'variable'
    lb = LabelEncoder()
``` |

| | |
|---|---|
| | ```python<br>        df[variable] = lb.fit_transform(df[variable])<br>        return lb<br>``` |
| Feature Engineering | ```python<br>food_category_le = encoder(df,'food_category')<br>brand_name_le = encoder(df,'brand_name')<br>food_department_le = encoder(df,'food_department')<br>food_family_le = encoder(df,'food_family')<br>promotion_name_le = encoder(df,'promotion_name')<br>store_city_le = encoder(df,'store_city')<br>#unit_per_case_le = encoder(df,'unit_per_case')<br>net_weight_le = encoder(df,'net_weight')<br>sales_le = encoder(df,'sales_country')<br>martial_le = encoder(df,'marital_status')<br>``` |
| Save Processed Data | ```python<br>import pickle<br>pickle.dump(rf,open('customers.pkl','wb'))<br>pickle.dump(food_category_le,open('food_category_le.pkl','wb'))<br>pickle.dump(brand_name_le,open('brand_name_le.pkl','wb'))<br>pickle.dump(promotion_name_le,open('promotion_name_le.pkl','wb'))<br>pickle.dump(store_city_le,open('store_city_le.pkl','wb'))<br>``` |