# Data Collection and Preprocessing Phase

| Date | 05 June 2024 |
|---|---|
| Team ID | 739975 |
| Project Title | To Predict Consumer Price Index |
| 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 | Screenshots |
|---|---|---|
| Data Overview | **Basic Statistics:** Summarize the dataset with measures such as mean, median, standard deviation, and range for each category. **Dimensions:** Describe the size of the dataset (e.g., number of records, number of features). **Structure:** Explain the format and organization of the data, including any hierarchical structure (e.g., categories and subcategories of items). |  |
| Univariate Analysis | Analyze each variable related to CPI, such as the price of vegetables, fruits, meat, housing, etc. Calculate and interpret measures of central tendency (mean, median, mode) and dispersion (variance, standard deviation) for each category. Visualize the distribution of each variable using histograms or box plots to understand their behavior over time |  |

| | | |
|---|---|---|
| Bivariate Analysis | Investigate the relationships between pairs of variables, such as the price of pulses and the price of vegetables. Calculate correlation coefficients to quantify the strength and direction of relationships between variables. Use scatter plots to visualize these relationships and identify any patterns or trends. |  |
| Multivariate Analysis | Explore the interactions between multiple variables simultaneously, such as the combined effect of prices of different food items on the overall CPI. Use techniques such as multiple regression analysis, principal component analysis (PCA), or cluster analysis to uncover underlying patterns and relationships. Visualize multivariate relationships using pair plots, heatmaps, or 3D plots |  |
| Outliers and Anomalies | - | - |
| **Data Preprocessing Code Screenshots** | | |
| Loading Data | Code to load the dataset into the preferred environment (e.g., Python, R). |  |

| Handling Missing Data | Code for identifying and handling missing values. |  |
|---|---|---|
| Feature Engineering | Enhance the accuracy and the robustness of the CPI predictions | |
| Save Processed Data | - | |



```
missing_values=df.isnull().sum()
missing_values

Sector                                  0
Year                                    0
Month                                   0
Cereals and products                    3
Meat and fish                           6
Egg                                     3
Milk and products                       3
Oils and fats                           3
Fruits                                  3
Vegetables                              3
Pulses and products                     3
Sugar and Confectionery                 3
Spices                                  3
Non-alcoholic beverages                 3
Prepared meals, snacks, sweets etc.     6
Food and beverages                      3
Pan, tobacco and intoxicants            6
Clothing                                6
Footwear                                6
Clothing and footwear                   6
Housing                               122
Fuel and light                          3
Household goods and services            6
Health                                  3
Transport and communication             6
...
Education                               6
Personal care and effects               6
Miscellaneous                           6
General index                           6
dtype: int64
```