

Data Collection and Preprocessing Phase

Date	15 March 2024
Team ID	SWTID1720184497
Project Title	Cereal Analysis Based on Ratings by using Machine Learning Techniques
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	Basic statistics, dimensions, and structure of the data.
Univariate Analysis	Exploration of individual variables (mean, median, mode, etc.).
Bivariate Analysis	Relationships between two variables (correlation, scatter plots).
Multivariate Analysis	Patterns and relationships involving multiple variables.
Outliers and Anomalies	Identification and treatment of outliers.
Data Preprocessing Code Screenshots	
Loading Data	<pre># Load the dataset into the preferred environment data = pd.read_csv('cereal.csv')</pre>

Handling Missing Data	<pre># Check for missing values print(data.isnull().sum()) # Handle missing values (example: fill with mean) data.fillna(data.mean(), inplace=True)</pre>
Data Transformation	<pre>from sklearn.preprocessing import StandardScaler # Normalize numerical features scaler = StandardScaler() numerical_features = data.select_dtypes(include=['int64', 'float64']).columns data[numerical_features] = scaler.fit_transform(data[numerical_features])</pre>
Feature Engineering	<pre># One-hot encode categorical variables data = pd.get_dummies(data, columns=['mfr', 'type'], drop_first=True)</pre>
Save Processed Data	<pre># Save the cleaned and processed data for future use data.to_csv('processed_cereal.csv', index=False)</pre>