# Data Collection and Preprocessing Phase

| | |
|---|---|
| Date | 15 March 2025 |
| Team ID | LTVIP2025TMID25168 |
| Project Title | Cosmatic Insights |
| Maximum Marks | 10 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 | **Short description about dataset**:<br>• **Shape:** (1472 rows, 16 columns)<br>• **Columns:**<br>• Label, Brand, Name, Price, Rank, Ingredients, Combination, Dry, Normal, Oily, Sensitive, Sensitive Skin Suitability, Dry Skin Suitability, Normal Skin Suitability, No. of Records, Oily Skin Suitability.<br>• **Missing Values:** None<br>• **Duplicate Rows:** None |
| Data Cleaning | **Data Cleaning Summary:**<br>• **Duplicates Removed:** 0 (No duplicates found)<br>• **Invalid Price Entries:** 0 (All prices are valid)<br>• **Invalid Rank Entries:** 0 (All ranks are within the 0-5 range)<br>• **Final Shape After Cleaning:** (1472 rows, 16 columns) |
| Data Transformation | **Data Transformation Summary**<br>• **Sorting:** Sorted products by price in descending order.<br>• **Filtering:** Extracted **666 products** suitable for sensitive skin.<br>• **Calculated Field:** Added a new column **"Price per** |

|  | Rating" (Price divided by Rank). |
| --- | --- |
|  | • **Pivot Table:** Shows the **average price** of different product types based on suitability for sensitive skin. |
|  | **Pivot Table - Average Price by Product Type & Skin Suitability** |

| Product Type | Suitable ($) | Not Suitable ($) |
| --- | --- | --- |
| Cleanser | 34.28 | 31.64 |
| Eye Cream | 67.55 | 59.04 |
| Face Mask | 45.93 | 40.36 |
| Moisturizer | 74.55 | 63.32 |
| Sun Protect | 47.03 | 45.02 |
| Treatment | 78.96 | 79.37 |

| Data Type Conversion | **Data Type Rectification Summary:**<br>• **No changes were needed** for most columns since they already had appropriate data types.<br>• **Key numeric columns verified:**<br> o **Price:** int64<br> o **Rank:** float64<br> o **No. of Records:** int64<br> o **Price per Rating:** float64 |
| --- | --- |
| Column Splitting and Merging | **Column Splitting & Merging Summary**<br>• **Splitting:** Extracted the **first ingredient** into a new column **"Primary Ingredient"**.<br>• **Merging:** Created a **"Full Product Name"** column by combining **Brand** and **Name**. |
| Data Modeling | **Data Modelling Overview:**<br>Since we have a **single dataset**, we can define relationships for a structured database if we plan to integrate it with other tables. Below are possible relationships:<br>**Potential Tables & Relationships**<br>1. **Products Table** (Main Table)<br> o **Primary Key:** Product ID (can be created if needed)<br> o Contains product details: Label, Brand, Name, Price, Rank, Ingredients, Full Product Name<br>2. **Skin Suitability Table** (One-to-Many Relationship)<br> o **Foreign Key:** Product ID<br> o Columns: Sensitive Skin Suitability, Dry Skin |

| | Suitability, Normal Skin Suitability, Oily Skin Suitability<br>3. **Ingredients Table** (One-to-Many Relationship)<br>   o **Foreign Key:** Product ID<br>   o Columns: Primary Ingredient, Full Ingredient List<br>4. **Reviews Table** (If data is available)<br>   o **Foreign Key:** Product ID<br>   o Columns: No. of Records, Rank, Price per Rating |
|---|---|
| Save Processed Data | Save the cleaned and processed data for future use.<br>cosmetics.csv |