# Data Collection and Preprocessing Phase

| Date | June 2025 |
|---|---|
| Team ID | Team-739774 |
| Project Title | Amazon Kindle Store Review Analysis |
| 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 | Data overview enables a thorough analysis of Kindle store reviews based on customer ratings and review texts, helping to identify reader preferences, sentiment trends, and product feedback. |
| Data Preparation | Gather attributes such as review-id, product-id, review-text, star-rating, and verified-purchase. Handle missing values, outliers, and inconsistencies. Clean the text data by removing special characters, HTML tags, and stop words. Encode categorical fields if necessary, and prepare for ML models. |
| Handling missing values | After loading the dataset, check for null values in each column. If found, take steps such as imputing missing text with placeholder tokens or removing rows if critical fields like review-text or star-rating are missing. Use fill (), median/mode imputation, and visualize with heatmaps for missing data. |
| Data Visualisation | Visualizing data helps detect trends and patterns in reviews. For Kindle reviews, this includes bar plots of rating distribution, word clouds for frequent terms, and sentiment heatmaps. Use libraries like Matplotlib, Seaborn, and Word Cloud to generate visual insights from the data. |
| Splitting The Dataset Into Dependent And Independent Variable | For sentiment modeling, split the dataset into independent variables (X) like review-text, and dependent variable (y) such as star-rating or a sentiment label. This prepares the data for supervised learning tasks such as sentiment classification or rating prediction. |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Import the libraries | ✓ Importing the libraries<br><br>```python
#import pandas library
import pandas as pd
#import numpy
import numpy as np
#import requests
import requests
#import io
import io
``` |
| Importing The Dataset | ```python
#import the dataset in the data variable
data=pd.read_csv('kindle_reviews.csv',on_bad_lines='skip',quoting=3)
data.head()
``` |

| | |
|---|---|
| Analyse The Data |  |
| Checking for Null Values or Taking Care of Missing Data |  |
| Joining review description and summary into one column |  |
| Splitting The Dataset Into Dependent And Independent Variable |  |