

## Milestone 1 & 2 Report

### Dataset Description:

<b>Transaction Data</b>	<p><b>Transaction_ID:</b> Unique identifier for each transaction.</p> <p><b>Customer_ID:</b> Unique identifier for the customer associated with the transaction.</p> <p><b>Product_ID:</b> Unique identifier for the product purchased in the transaction.</p> <p><b>Transaction_Date:</b> The date when the transaction occurred.</p>
<b>Product Information</b>	<p><b>Category:</b> The category to which the product belongs (i.e., Electronics, Clothing).</p> <p><b>Units_Sold:</b> The quantity of the product sold in the transaction.</p> <p><b>Discount_Applied:</b> The discount percentage applied to the product during the transaction.</p> <p><b>Revenue:</b> Total revenue generated from the transaction, calculated as Price x Units Sold x (1 - Discount).</p>
<b>Customer Demographics</b>	<p><b>Customer_ID:</b> Unique identifier for each customer (repeated for easier reference).</p>
<b>Advertising Metrics</b>	<p><b>Clicks:</b> Number of ad clicks associated with the product during the time of the transaction.</p> <p><b>Impressions:</b> Number of ad impressions served during the campaign.</p> <p><b>Conversion_Rate:</b> Calculated as Clicks / Impressions, representing the percentage of impressions that resulted in clicks.</p> <p><b>Ad_CTR:</b> Click-through rate (CTR) for the advertisement, representing the effectiveness of the ad campaign.</p> <p><b>Ad_CPC:</b> Cost-per-click for the advertisement.</p> <p><b>Ad_Spend:</b> Total advertising spend for the product, calculated as Ad_CTR x Ad_CPC x 1000.</p>

<b>Seasonal and Regional Information</b>	<p><b>Region:</b> The geographical region where the transaction occurred (e.g., North America, Europe, Asia).</p> <p><b>Seasonality Effects:</b> Implied through patterns in transaction dates and revenue, reflecting holiday promotions and season-based purchasing trends.</p>
--	---

### **Data Cleaning:**

We have checked the missing values, duplicates in addition to removing the trivial columns ('Transaction\_ID', 'Customer\_ID', 'Product\_ID') and converting 'Transaction\_Date' into a time column. Furthermore, we have used the interquartile range method to remove the outliers, remove the highly correlated features to avoid overfitting occurrence, and round the float columns to avoid noise effects which may sometimes lead to overfitting. The categorical (nominal) columns are dealt with through one-hot encoding then checking the categories balance, finding that almost all the categories are divided equally. We also have checked each column distribution to do the scaling based on their distribution: for the ('Ad\_Spend', 'Revenue', 'Conversion\_Rate'), they are exponentially distributed so a logarithmic scaling is done while for the ('Ad\_CPC', 'Ad\_CTR', 'Clicks', 'Impressions', 'Discount\_Applied'), they are uniformly distributed so a minmax scaling is done.

### **Data visualization:**

**Univariate analysis:** 1. We did pie chart for the 'Region' and 'Category' columns with finding that almost the distribution are balanced in both categorical features  
2. We also did a histogram for the numerical variables to check the data distribution in order to decide which scaling technique will be used

**Bivariate Analysis:** 1. We have checked the Revenue over time (aggregated Monthly), finding that November is the month with the highest revenue which is reassured through checking the units sold among months to get November also scores the highest all over the year.

2. We have checked the revenue per category, finding almost equal revenues with the electronics, the highest one but a very slight increase among the other categories with checking the revenue per region, which appears to be similar across them.

3. We have looked for the Units Sold by Region and Season, getting that the highest number of units sold is in North America in November and also checked the advertising effects on the revenue which is observed to have a real impact.

4. We have also checked the holidays versus the units sold because most of the time, the businesses earn more at holidays, so we have found that the median number of units

sold is slightly higher on holidays than on non-holidays which indicates a potential boost in sales during holidays. However, there are significantly more extreme high outliers on non-holidays than holidays, indicating occasional promotional campaigns or irregular sales spikes on non-holidays to attract customers.

5. We have done a time series plot which shows as the previous plots that November has the highest number of 'units\_sold'.

**Multivariate Analysis:** we have done a pairplot between the columns and the target feature to check how they are related to each other.

**Feature engineering:**

We have checked the correlation between the numerical features to remove the highly correlated ones with the target column to prevent or lessen the probability of overfitting occurrence, then checked the heatmap of revenue by day of week and month. We found that October and November have highest correlations as demonstrated from the other analysis in addition to having the highest correlation on Friday (which does make sense).