

## Data Collection and Preprocessing Phase

|               |  |
|---------------|--|
| Date          | 7th July 2024                            |
| Team ID       | 739719                                   |
| Project Title | Garment Workers Productivity Predictions |
| Maximum Marks | 6 Marks                                  |

### Data Exploration and Preprocessing Template

Data collection for gather data from internal and external sources, assesses quality issues like missing values and duplicates, normalize features, engineer new features, data merging and splitting then split into training and testing.

| Section       | Description   |
|---------------|---|
| Data Overview | <p><b>Basic statistics:</b> Summarize the central tendency, dispersion, and shape of the dataset's distribution.</p> <p><b>dimensions:</b> Display the number of rows and columns in each dataset to understand the size of the data.</p> <p><b>structure of the data:</b> Provide information about the data types and non-null values in each column to understand the structure.</p> |

|                     |  |
|---------------------|--|
| Univariate Analysis | <p><b>1. Summary Statistics</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Calculate and display basic statistics such as mean, median, mode, standard deviation, and range for individual variables.</li> </ul> <p><b>2.. Distribution Plots</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Visualize the distribution of individual variables using histograms and density plots.</li> </ul> <p><b>3. Box Plots</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Use box plots to visualize the spread and identify potential outliers in individual variables</li> </ul> <p>This code will provide a comprehensive univariate analysis for individual variables, helping to understand their distributions, central tendencies, and variations.</p> |
| Bivariate Analysis  | <p><b>1.Correlation</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Calculate the correlation coefficient to measure the strength and direction of the relationship between two variables.</li> </ul> <p><b>2. Scatter Plot</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Create scatter plots to visually inspect the relationship between two variables.</li> </ul>  |

|  |   |
|--|---|
| Multivariate Analysis                      | <p><b>1. Pair Plot</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Visualize pairwise relationships between multiple variables to identify potential interactions.</li> </ul> <p><b>2. Heatmap for Correlation Matrix</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Visualize the correlation matrix to see the strength of relationships between multiple variables.</li> </ul> <p><b>3. Multiple Regression Analysis</b></p> <ul style="list-style-type: none"> <li><b>Description:</b> Fit a multiple regression model to understand the combined effect of multiple predictors on a target variable.</li> </ul>   |
| Outliers and Anomalies                     | <p><b>Treatment of Outliers:</b></p> <p><b>1.Strategies</b></p> <ul style="list-style-type: none"> <li><b>Remove outliers:</b> Exclude them from further analysis if they are likely to be errors or irrelevant extreme values.</li> <li><b>Transform variables:</b> Apply mathematical transformations like log or square root to reduce skewness caused by outliers.</li> <li><b>Winsorization:</b> Cap extreme values by replacing them with values at a specified percentile (e.g., 95th percentile).</li> </ul> <p><b>2.Considerations:</b></p> <ul style="list-style-type: none"> <li><b>Impact assessment:</b> Evaluate how removing or transforming outliers affects the overall dataset and model performance.</li> <li><b>Documentation:</b> Document all decisions and actions taken regarding outlier handling for transparency and reproducibility.</li> </ul> |
| <b>Data Preprocessing Code Screenshots</b> |   |
| Loading Data                               | <pre>df = pd.read_csv(r'C:\Users\srira\Downloads\miniProject\garments_worker_productivity.csv') df.head()</pre>   |

|                       |   |
|-----------------------|---|
| Handling Missing Data | <pre>df2.isnull().sum()  quarter      0 department   0 day          0 team         0 targeted_productivity  0 smv          0 wip         506 over_time    0 incentive    0 idle_time    0 idle_man     0 no_of_style_change  0 no_of_workers  0 actual_productivity  0 dtype: int64</pre> |
| Data Transformation   | <pre># Read CSV file df = pd.read_csv('data.csv')  # Print the first 5 rows df.head()  # Print the last 5 rows df.tail()  # Print the shape of the DataFrame df.shape  # Print the data types of the columns df.dtypes</pre>  |
| Feature Engineering   | <pre>df7 = df6[(df6.over_time &gt; lower_limit) &amp; (df6.over_time &lt; upper_limit)]  df7.shape  (1176, 13)</pre>  |
| Save Processed Data   | <pre># Create DataFrame df = pd.DataFrame(data)  # Save to CSV df.to_csv('data.csv', index=False)  print("Data saved to data.csv")</pre>  |