

## Data Collection and Preprocessing Phase

Date	15 March 2024
Team ID	Team-739770
Project Title	Predicting the energy output of wind turbine based on weather condition
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	<p>Basic statistics, dimensions, and structure of the data.</p> <pre>&lt;class 'pandas.core.frame.DataFrame'&gt; RangeIndex: 4447 entries, 0 to 4446 Data columns (total 6 columns): #      Column              Non-Null Count  Dtype ---  - 0     Wind Speed (m/s)      4447 non-null   float64 1     Wind Direction        4447 non-null   float64 2     maxtempC              4447 non-null   int64 3     humidity              4447 non-null   float64 4     pressure              4447 non-null   float64 5     Output_Energy         4447 non-null   float64 dtypes: float64(5), int64(1)</pre>
Univariate Analysis	Exploration of individual variables (mean, median, mode, etc.).

	<div><div>Wind Speed (m/s)</div><div>Wind Direction</div><div>maxtempC</div><div>humidity</div><div>pressure</div><div>Output_Eng</div></div> <div><div>count</div><div>4447.000000</div><div>4447.000000</div><div>4447.000000</div><div>4447.000000</div><div>4447.000000</div></div> <div><div>mean</div><div>7.357389</div><div>140.667803</div><div>8.535192</div><div>78.648874</div><div>1019.49165</div></div> <div><div>std</div><div>4.361162</div><div>93.616266</div><div>3.034301</div><div>9.004574</div><div>5.15432</div></div> <div><div>min</div><div>0.000000</div><div>0.000000</div><div>4.000000</div><div>54.125000</div><div>1004.54160</div></div> <div><div>25%</div><div>3.669025</div><div>53.272396</div><div>6.000000</div><div>74.000000</div><div>1015.87500</div></div> <div><div>50%</div><div>6.717962</div><div>143.424896</div><div>8.000000</div><div>80.041667</div><div>1020.83333</div></div> <div><div>75%</div><div>10.197950</div><div>206.816154</div><div>12.000000</div><div>84.708333</div><div>1023.45833</div></div> <div><div>max</div><div>21.621000</div><div>359.942291</div><div>14.000000</div><div>93.958333</div><div>1028.20833</div></div>
Bivariate Analysis	<div>Relationships between two variables (correlation, scatter plots).</div> <div><div>Wind Speed (m/s)</div><div>Wind Direction</div><div>maxtempC</div><div>humidity</div><div>pressure</div></div> <div><div>Wind Speed (m/s)</div><div>1.000000</div><div>0.017336</div><div>0.339107</div><div>0.15185</div></div> <div><div>Wind Direction</div><div>0.017336</div><div>1.000000</div><div>0.080762</div><div>0.31354</div></div> <div><div>maxtempC</div><div>0.339107</div><div>0.080762</div><div>1.000000</div><div>0.06532</div></div>

	<div> <div>humidity</div> <div>0.151853-0.3135420.0653291.000000-0.12</div> </div> <div> <div>pressure</div> <div>0.234967-0.0209620.5973240.1292951.00</div> </div> <div> <div>Output_Energy</div> <div>0.8824570.1229130.4033820.251067-0.24</div> </div>
Multivariate Analysis	<p>Patterns and relationships involving multiple variables.</p>
Outliers and Anomalies	<pre> for col in df.columns:     q1 = df[col].quantile(0.25)     q3 = df[col].quantile(0.75)     iqr = q3 - q1     lower_bound = q1 - 1.5 * iqr     upper_bound = q3 + 1.5 * iqr     df[col]=np.where(df[col]&lt;lower_bound,lower_bound,df[col])     df[col]=np.where(df[col]&gt;upper_bound,upper_bound,df[col])  for col in df.columns:     sns.boxplot(df[col])     plt.show() </pre>
Data Preprocessing Code Screenshots	

Loading Data	<pre>data = pd.read_csv('/content/data.csv') target = pd.read_csv('/content/target.csv')</pre>
Handling Missing Data	<pre>Data columns (total 6 columns): #   Column                Non-Null Count  Dtype ---  - 0   Wind Speed (m/s)      4447 non-null   float64 1   Wind Direction        4447 non-null   float64 2   maxtempC              4447 non-null   int64 3   humidity              4447 non-null   float64 4   pressure              4447 non-null   float64 5   Output_Energy         4447 non-null   float64 dtypes: float64(5), int64(1) memory usage: 208.6 KB</pre>
Data Transformation	<pre>Scaler = StandardScaler() for col in df.columns:     if col != 'Output_Energy':         df[col] = Scaler.fit_transform(df[[col]])  df.head()</pre>
Feature Engineering	Code for creating new features or modifying existing ones.
Save Processed Data	Code to save the cleaned and processed data for future use. df = data