# Data Collection and Preprocessing Phase

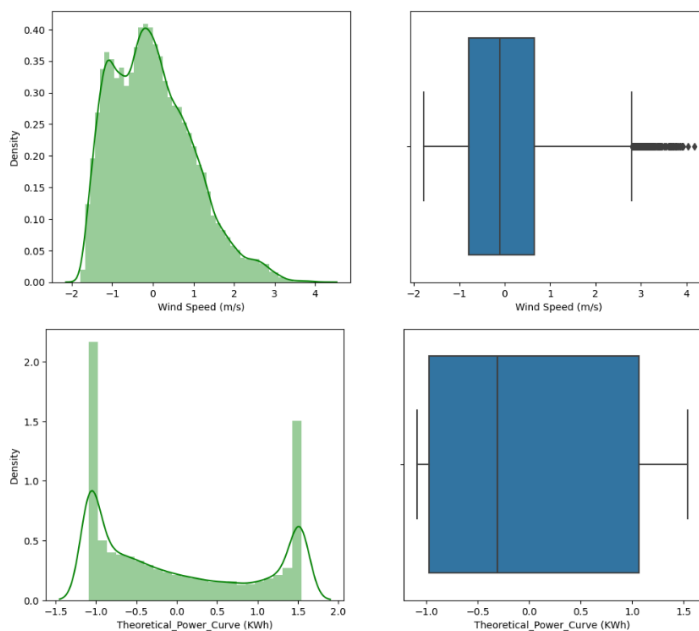| | |
|---|---|
| Date | 12 July 2024 |
| Team ID | SWTID1720108739 |
| Project Title | **Predicting The Energy Output Of Wind Turbine Based On Weather Condition** |
| Maximum Marks | 6 Marks |

## 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 | Dimension:<br>50530 rows × 5columns<br>Descriptive statistics:<br> |



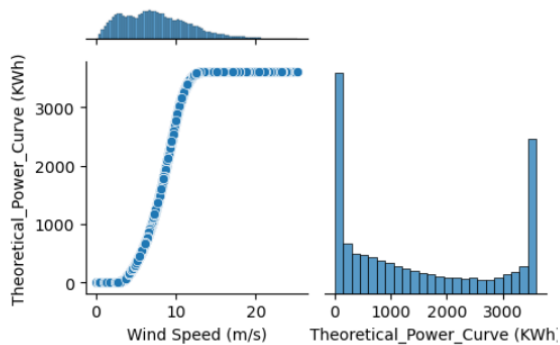| [18]: | | LV ActivePower (kW) | Wind Speed (m/s) | Theoretical_Power_Curve (KWh) | Wind Direction (°) |
|---|---|---|---|---|---|
| | count | 50530.000000 | 50530.000000 | 50530.000000 | 50530.000000 |
| | mean | 1307.684332 | 7.557952 | 1492.175463 | 123.687559 |
| | std | 1312.459242 | 4.227166 | 1368.018238 | 93.443736 |
| | min | -2.471405 | 0.000000 | 0.000000 | 0.000000 |
| | 25% | 50.677890 | 4.201395 | 161.328167 | 49.315437 |
| | 50% | 825.838074 | 7.104594 | 1063.776283 | 73.712978 |
| | 75% | 2482.507568 | 10.300020 | 2964.972462 | 201.696720 |
| | max | 3618.732910 | 25.206011 | 3600.000000 | 359.997589 |

| Univariate Analysis |  |
|---|---|
| | ```python
[136]: #univariate analysis
       for col in num_col:
           fig, ax = plt.subplots(1, 2, figsize=(12, 5))
           sns.distplot(x_train[col], ax=ax[0] ,color='green')
           sns.boxplot(x=x_train[col], ax=ax[1])
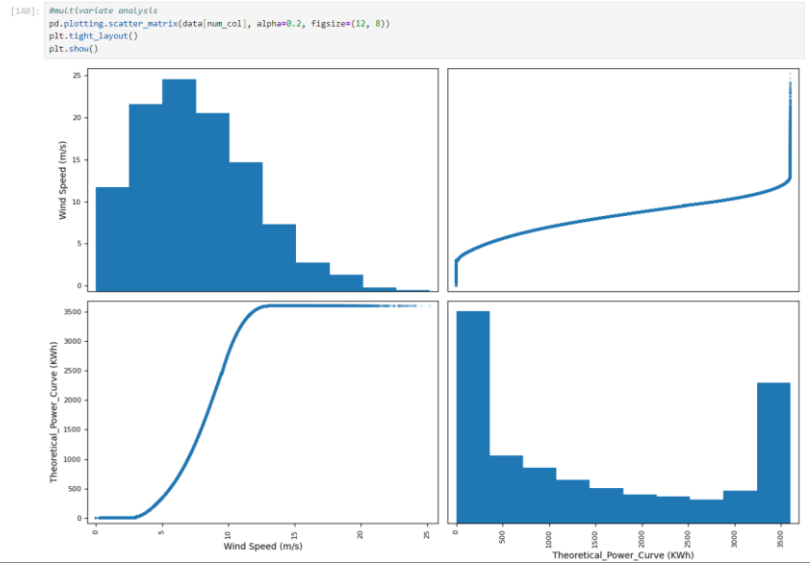           plt.show()
``` |
| Bivariate Analysis | ```python
[138]: # bivariate analysis
       sns.pairplot(data[num_col], corner=True)
       sns.boxplot(x=x_train[col], ax=ax[1])
       plt.tight_layout()
       plt.show()
```
 |

| | |
|---|---|
| Multivariate Analysis |  |
| Outliers and Anomalies | - |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data |  |

| | |
|---|---|
| Handling Missing Data | • There is no null values in the dataset provided<br><br>[14]: `data.isnull().sum()`<br><br>[14]: Date/Time     0<br>LV ActivePower (kW)     0<br>Wind Speed (m/s)     0<br>Theoretical_Power_Curve (Kwh)     0<br>Wind Direction (°)     0<br>dtype: int64<br><br>[16]: `#there is no null values in the dataset provided`<br><br>[18]: `data.describe()`<br><br>[18]:<br><table><tr><td></td><td>LV ActivePower (kW)</td><td>Wind Speed (m/s)</td><td>Theoretical_Power_Curve (KWh)</td><td>Wind Direction (°)</td></tr><tr><td>count</td><td>50530.000000</td><td>50530.000000</td><td>50530.000000</td><td>50530.000000</td></tr><tr><td>mean</td><td>1307.684332</td><td>7.557952</td><td>1492.175463</td><td>123.687559</td></tr><tr><td>std</td><td>1312.459242</td><td>4.227166</td><td>1368.018238</td><td>93.443736</td></tr><tr><td>min</td><td>-2.471405</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr><tr><td>25%</td><td>50.677890</td><td>4.201395</td><td>161.328167</td><td>49.315437</td></tr><tr><td>50%</td><td>825.838074</td><td>7.104594</td><td>1063.776283</td><td>73.712978</td></tr><tr><td>75%</td><td>2482.507568</td><td>10.300020</td><td>2964.972462</td><td>201.696720</td></tr><tr><td>max</td><td>3618.732910</td><td>25.206011</td><td>3600.000000</td><td>359.997589</td></tr></table><br>[20]: `#we see in the active power there is a negative value`<br>`count_negative_values = (data['LV ActivePower (kW)'] < 0).sum()`<br>`count_negative_values`<br><br>[20]: 57<br><br>[22]: `#so we changed all the negative values to 0 for better preprocessing`<br>`data.loc[data['LV ActivePower (kW)'] < 0, 'LV ActivePower (kW)'] = 0`<br><br>[24]: `data.describe()`<br><br>[24]:<br><table><tr><td></td><td>LV ActivePower (kW)</td><td>Wind Speed (m/s)</td><td>Theoretical_Power_Curve (KWh)</td><td>Wind Direction (°)</td></tr><tr><td>count</td><td>50530.000000</td><td>50530.000000</td><td>50530.000000</td><td>50530.000000</td></tr><tr><td>mean</td><td>1307.684699</td><td>7.557952</td><td>1492.175463</td><td>123.687559</td></tr><tr><td>std</td><td>1312.458876</td><td>4.227166</td><td>1368.018238</td><td>93.443736</td></tr><tr><td>min</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr><tr><td>25%</td><td>50.677890</td><td>4.201395</td><td>161.328167</td><td>49.315437</td></tr><tr><td>50%</td><td>825.838074</td><td>7.104594</td><td>1063.776283</td><td>73.712978</td></tr><tr><td>75%</td><td>2482.507568</td><td>10.300020</td><td>2964.972462</td><td>201.696720</td></tr><tr><td>max</td><td>3618.732910</td><td>25.206011</td><td>3600.000000</td><td>359.997589</td></tr></table> |
| Data Transformation | [32]: `#Scaling on Independent Features: to avoid biasing of results`<br>`from sklearn.preprocessing import StandardScaler`<br><br>[88]: `scale=StandardScaler()`<br><br>[90]: `x=scale.fit_transform(x)`<br><br>[92]: `x`<br><br>[92]: `array([[-0.78643484, -0.53147626, -1.61557807, ..., -1.68278034,`<br>`        -1.6608638 , -1.42360538],`<br>`       [-0.71071243, -0.44611545, -1.61557807, ..., -1.68278034,`<br>`        -1.6608638 , -1.42360538],`<br>`       [-0.80502315, -0.55402096, -1.61557807, ..., -1.68278034,`<br>`        -1.6608638 , -1.42360538],`<br>`       ...,`<br>`       [ 0.21645342,  0.20756566,  1.610911  , ...,  1.76866227,`<br>`         1.65585814,  1.71135408],`<br>`       [ 0.6770496 ,  0.44082298,  1.610911  , ...,  1.76866227,`<br>`         1.65585814,  1.71135408],`<br>`       [ 0.94079255,  0.57281963,  1.610911  , ...,  1.76866227,`<br>`         1.65585814,  1.71135408]])` |
| Feature Engineering | • Attached the codes in final submission. |
| Save Processed Data | - |