# Data Collection and Preprocessing Phase

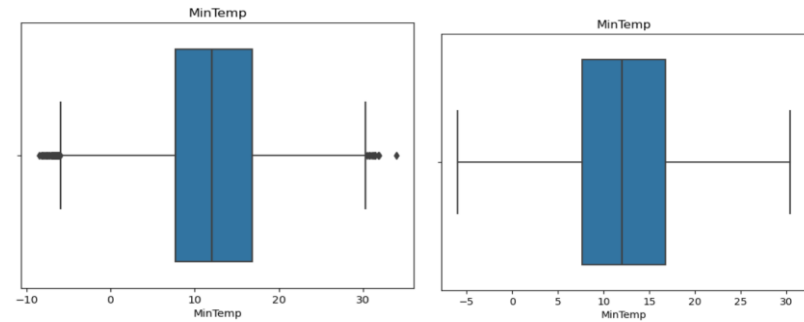| Date | 24 April 2024 |
|---|---|
| Team ID | Team-738169 |
| Project Title | Rainfall Prediction Using Machine Learning |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

This report presents the findings and steps undertaken during the exploration and preprocessing of the rainfall dataset. The primary objectives were to gain insights into the data distribution, identify potential issues, and prepare the data for further analysis and modeling.

| Section | Description |
|---|---|
| Data Overview | Dimension:<br>(145460, 23)<br>Descriptive Statistics<br><br>|  | **MinTemp** | **MaxTemp** | **Rainfall** | **Evaporation** | **Sunshine** |<br>|---|---|---|---|---|---|<br>| **count** | 145460.000000 | 145460.000000 | 145460.000000 | 145460.000000 | 145460.000000 |<br>| **mean** | 12.192053 | 23.215962 | 2.307990 | 5.464988 | 7.609641 |<br>| **std** | 6.365780 | 7.088358 | 8.389771 | 4.210586 | 3.785983 |<br>| **min** | -8.500000 | -4.800000 | 0.000000 | 0.000000 | 0.000000 |<br>| **25%** | 7.700000 | 18.000000 | 0.000000 | 2.600000 | 4.800000 |<br>| **50%** | 12.000000 | 22.600000 | 0.000000 | 4.800000 | 8.400000 |<br>| **75%** | 16.800000 | 28.200000 | 0.600000 | 7.400000 | 10.600000 |<br>| **max** | 33.900000 | 48.100000 | 371.000000 | 145.000000 | 14.500000 | |
| Univariate Analysis |  |

| | |
|---|---|
| |  |
| Bivariate Analysis |  |
| Multivariate Analysis | - |

| Outliers and Anomalies |  Handeling Outliers by IQR(Inter Quartile Range). <br><br>```python<br>IQR=df.MinTemp.quantile(0.75)-df.MinTemp.quantile(0.25)<br>lower_bridge=df.MinTemp.quantile(0.25)-(IQR*1.5)<br>upper_bridge=df.MinTemp.quantile(0.75)+(IQR*1.5)<br>print(lower_bridge, upper_bridge)<br><br>-5.950000000000002 30.450000000000003<br>```<br><br>```python<br>df.loc[df['MinTemp']>=30.45,'MinTemp']=30.45<br>df.loc[df['MinTemp']<=-5.95,'MinTemp']=-5.95<br>``` |
|---|---|

## Data Preprocessing Code Screenshots

| Loading Data | ```python<br>df = pd.read_csv("Dataset.csv")<br>pd.set_option("display.max_columns", None)<br>df.head()<br>``` <br><br>  |
|---|---|
| Handling Missing Data | **Replacing the null values of the remaining continuous features by median** <br><br>```python<br>for feature in continuous_feature:<br>    if(df[feature].isnull().sum()*100/len(df))>0:<br>        df[feature] = df[feature].fillna(df[feature].median())<br>``` <br><br>**Replacing the null values of the discrete features by mode** <br><br>```python<br>def mode_nan(df,variable):<br>    mode=df[variable].value_counts().index[0]<br>    df[variable].fillna(mode,inplace=True)<br>mode_nan(df,"Cloud9am")<br>mode_nan(df,"Cloud3pm")<br>``` |

| Data Transformation | **Handling categorical features using One Hot Encoding**<br><br>```python
df["RainToday"] = pd.get_dummies(df["RainToday"], drop_first = True, dtype = np.int64)
df["RainTomorrow"] = pd.get_dummies(df["RainTomorrow"], drop_first = True, dtype = np.int64)
df
```<br>**Performing Label Encoding on "Location"**<br><br>```python
df1 = df.groupby(["Location"])["RainTomorrow"].value_counts().sort_values().unstack()
``` |
|---|---|
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | - |