

```
1 !pip install -r rec.txt
```

```

Collecting numpy==1.19.5 (from -r rec.txt (line 1))
  Using cached numpy-1.19.5.zip (7.3 MB)
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Preparing metadata (pyproject.toml) ... done
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from -r rec.txt (line 2)) (2.2.
Collecting matplotlib==3.3.4 (from -r rec.txt (line 3))
  Using cached matplotlib-3.3.4.tar.gz (37.9 MB)
  Preparing metadata (setup.py) ... done
Collecting seaborn==0.11.1 (from -r rec.txt (line 4))
  Using cached seaborn-0.11.1-py3-none-any.whl.metadata (2.3 kB)
Collecting scikit-learn==0.24.1 (from -r rec.txt (line 5))
  Using cached scikit-learn-0.24.1.tar.gz (7.4 MB)
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Preparing metadata (pyproject.toml) ... done
Collecting scipy==1.5.4 (from -r rec.txt (line 6))
  Using cached scipy-1.5.4.tar.gz (25.2 MB)
error: subprocess-exited-with-error

  × pip subprocess to install build dependencies did not run successfully.
  | exit code: 1
  |_ See above for output.

note: This error originates from a subprocess, and is likely not a problem with pip.
Installing build dependencies ... error
error: subprocess-exited-with-error

  × pip subprocess to install build dependencies did not run successfully.
  | exit code: 1
  |_ See above for output.

note: This error originates from a subprocess, and is likely not a problem with pip.

```

```
1 !pip install catboost
```

```

Collecting catboost
  Downloading catboost-1.2.7-cp310-cp310-manylinux2014_x86_64.whl.metadata (1.2 kB)
Requirement already satisfied: graphviz in /usr/local/lib/python3.10/dist-packages (from catboost) (0.20.3)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from catboost) (3.7.1)
Requirement already satisfied: numpy<2.0,>=1.16.0 in /usr/local/lib/python3.10/dist-packages (from catboost) (1.2
Requirement already satisfied: pandas>=0.24 in /usr/local/lib/python3.10/dist-packages (from catboost) (2.2.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from catboost) (1.13.1)
Requirement already satisfied: plotly in /usr/local/lib/python3.10/dist-packages (from catboost) (5.24.1)
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from catboost) (1.16.0)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24->catbo
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24->catb
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catb
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->cat
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->cat
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catbo
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboos
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catb
Requirement already satisfied: tenacity>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from plotly->catboost)
Downloading catboost-1.2.7-cp310-cp310-manylinux2014_x86_64.whl (98.7 MB)
98.7/98.7 MB 7.4 MB/s eta 0:00:00

Installing collected packages: catboost
Successfully installed catboost-1.2.7

```

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn import preprocessing
6 import scipy.stats as stats
7 from sklearn.model_selection import train_test_split
8 from collections import Counter
9 from imblearn.over_sampling import SMOTE

```

```


10 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
11 from sklearn import metrics
12 from sklearn.ensemble import RandomForestClassifier
13 from catboost import CatBoostClassifier
14 from xgboost import XGBClassifier
15 from sklearn.svm import SVC
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.naive_bayes import GaussianNB
18 from sklearn.neighbors import KNeighborsClassifier
19 import joblib




```

```

1 df = pd.read_csv("processed_weather_data.csv")
2 pd.set_option("display.max_columns", None)
3 df

```




	Date	Location	Temp	WindSpeed	Humidity	RainToday	RainTomorrow	
0	2008-12-01	Albury	18.15	22.0	46.5	No	No	
1	2008-12-02	Albury	16.25	13.0	34.5	No	No	
2	2008-12-03	Albury	19.30	22.5	34.0	No	No	
3	2008-12-04	Albury	18.60	10.0	30.5	No	No	
4	2008-12-05	Albury	24.90	13.5	57.5	No	No	
...	
145455	2017-06-21	Uluru	13.10	12.0	37.5	No	No	
145456	2017-06-22	Uluru	14.45	11.0	38.5	No	No	
145457	2017-06-23	Uluru	16.15	9.0	38.5	No	No	
145458	2017-06-24	Uluru	17.40	10.0	37.5	No	No	
145459	2017-06-25	Uluru	14.90	17.0	49.0	No	NaN	

145460 rows × 7 columns

```

1 numerical_feature = [feature for feature in df.columns if df[feature].dtypes != 'O']
2 discrete_feature = [feature for feature in numerical_feature if len(df[feature].unique()) < 25]
3 continuous_feature = [feature for feature in numerical_feature if feature not in discrete_feature]
4 categorical_feature = [feature for feature in df.columns if feature not in numerical_feature]
5 print("Numerical Features Count {}".format(len(numerical_feature)))
6 print("Discrete feature Count {}".format(len(discrete_feature)))
7 print("Continuous feature Count {}".format(len(continuous_feature)))
8 print("Categorical feature Count {}".format(len(categorical_feature)))


```

 Numerical Features Count 3
Discrete feature Count 0
Continuous feature Count 3
Categorical feature Count 4

```

1 # Handle Missing Values
2 df.isnull().sum()*100/len(df)

```



	0
Date	0.000000
Location	0.000000
Temp	0.594665
WindSpeed	0.833906
Humidity	1.297264
RainToday	2.241853
RainTomorrow	2.245978

dtype: float64

```
1 print(numerical_feature)
```

```
↳ ['Temp', 'WindSpeed', 'Humidity']
```


```
1 def randomsampleimputation(df, variable):
2     df[variable]=df[variable]
3     random_sample=df[variable].dropna().sample(df[variable].isnull().sum(),random_state=0)
4     random_sample.index=df[df[variable].isnull()].index
5     df.loc[df[variable].isnull(),variable]=random_sample
```

```
1 df
```

	Date	Location	Temp	WindSpeed	Humidity	RainToday	RainTomorrow	
0	2008-12-01	Albury	18.15	22.0	46.5	No	No	
1	2008-12-02	Albury	16.25	13.0	34.5	No	No	
2	2008-12-03	Albury	19.30	22.5	34.0	No	No	
3	2008-12-04	Albury	18.60	10.0	30.5	No	No	
4	2008-12-05	Albury	24.90	13.5	57.5	No	No	
...	
145455	2017-06-21	Uluru	13.10	12.0	37.5	No	No	
145456	2017-06-22	Uluru	14.45	11.0	38.5	No	No	
145457	2017-06-23	Uluru	16.15	9.0	38.5	No	No	
145458	2017-06-24	Uluru	17.40	10.0	37.5	No	No	
145459	2017-06-25	Uluru	14.90	17.0	49.0	No	NaN	

145460 rows × 7 columns

```
1 for feature in continuous_feature:
2     data=df.copy()
3     sns.distplot(df[feature])
4     plt.xlabel(feature)
5     plt.ylabel("Count")
6     plt.title(feature)
7     plt.figure(figsize=(15,15))
8     plt.show()
```

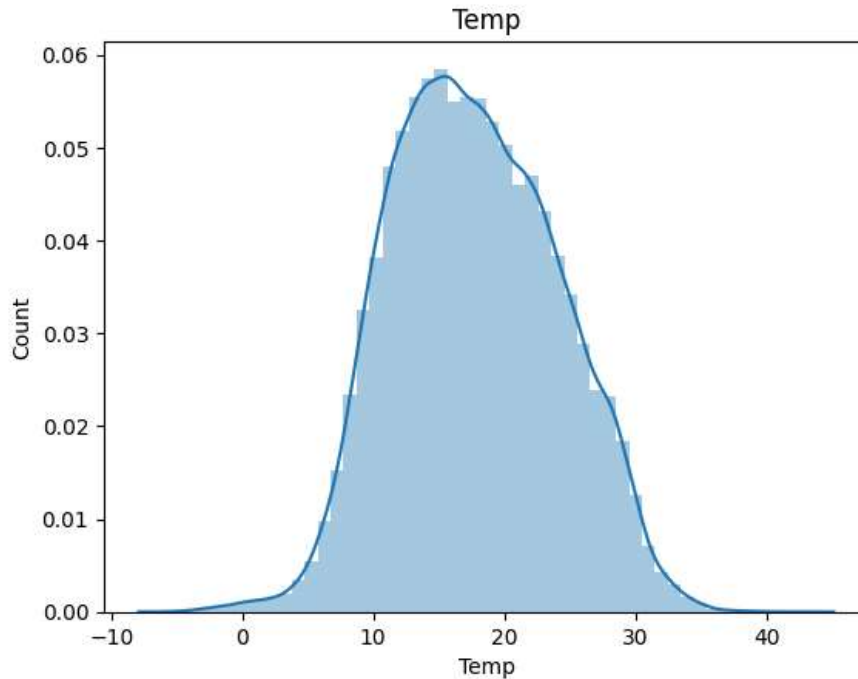
 <ipython-input-12-f3f36df57cc9>:3: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(df[feature])
```



<Figure size 1500x1500 with 0 Axes>

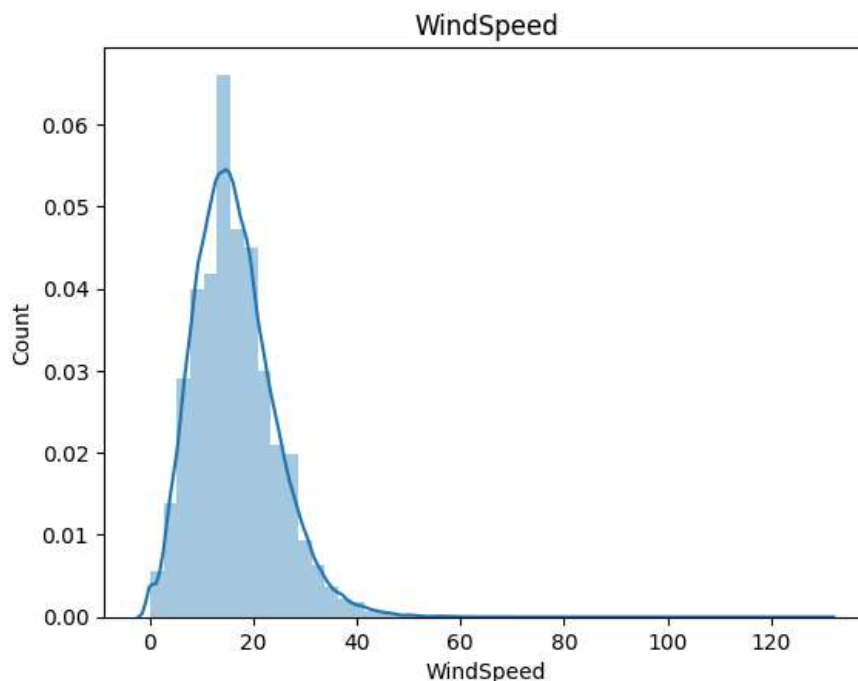
<ipython-input-12-f3f36df57cc9>:3: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(df[feature])
```



<Figure size 1500x1500 with 0 Axes>

<ipython-input-12-f3f36df57cc9>:3: UserWarning:

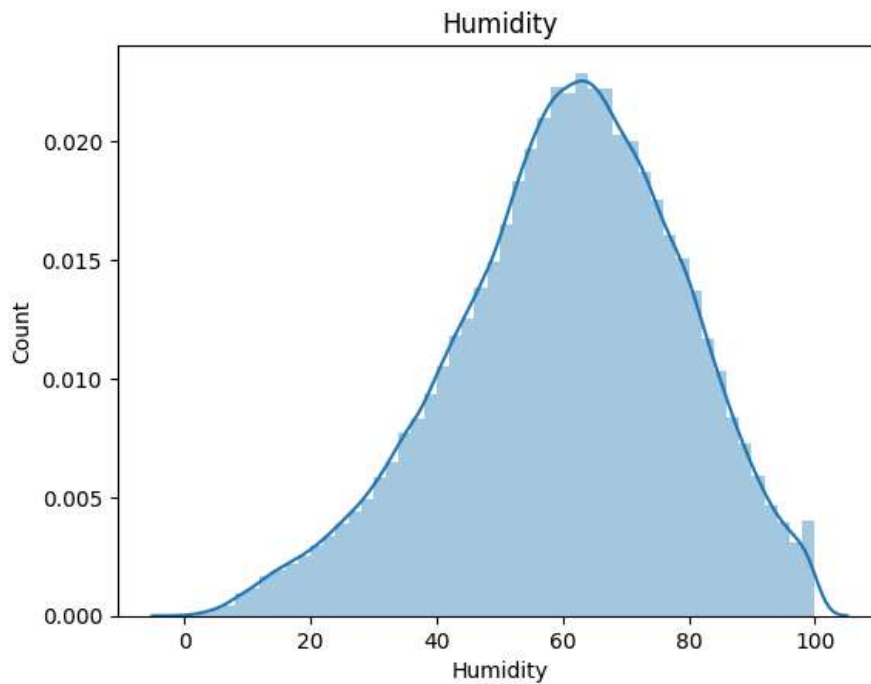
``distplot`` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see

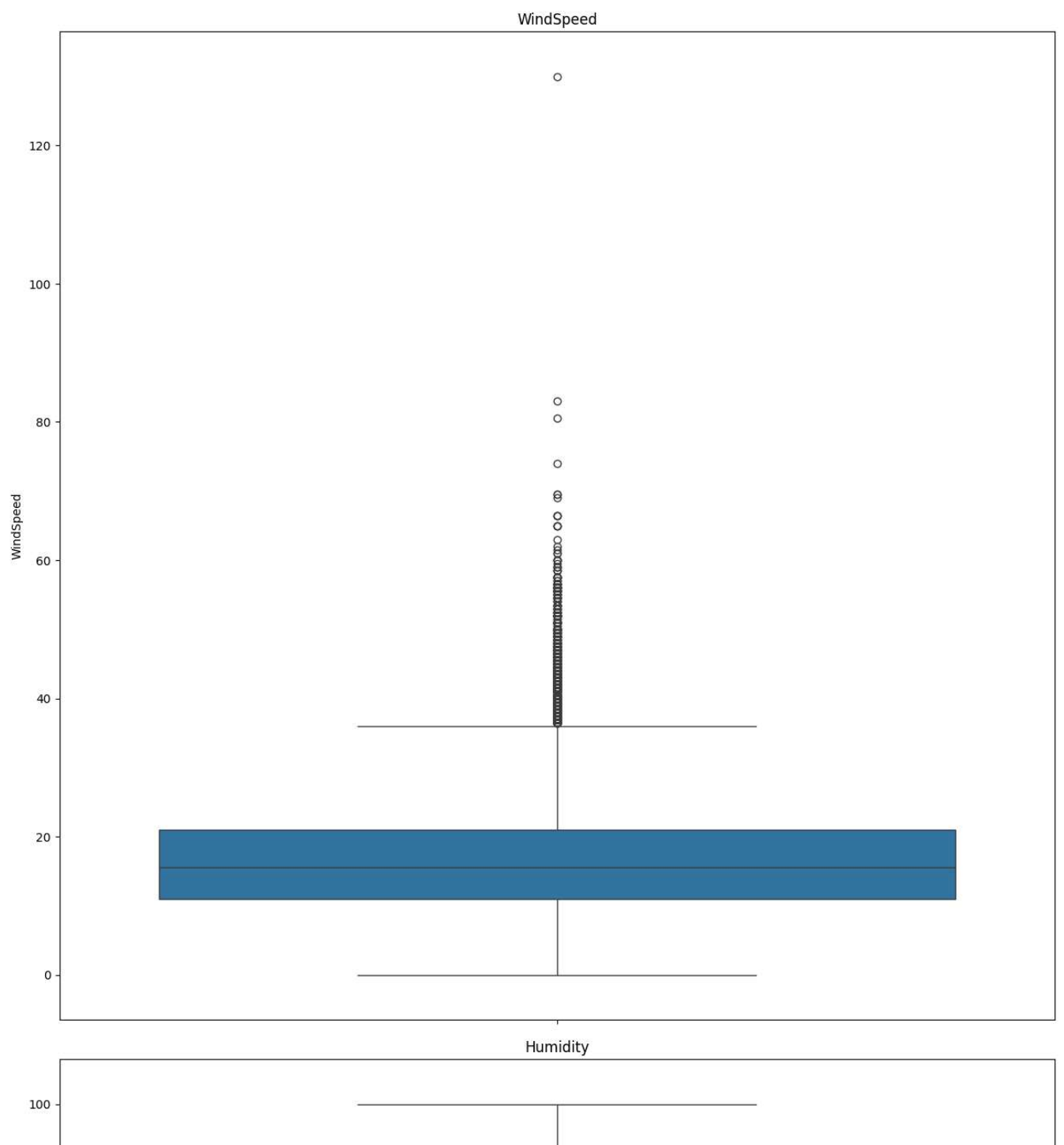
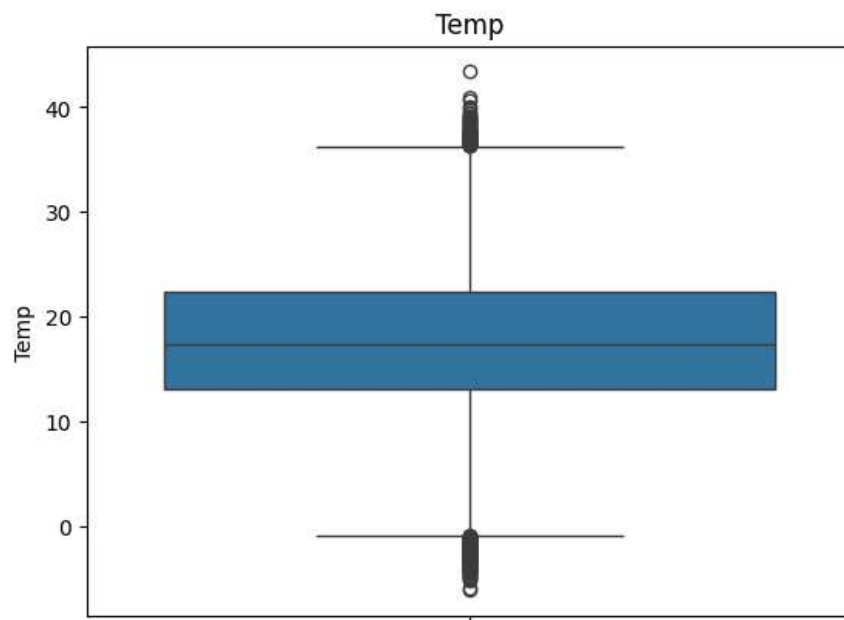
<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

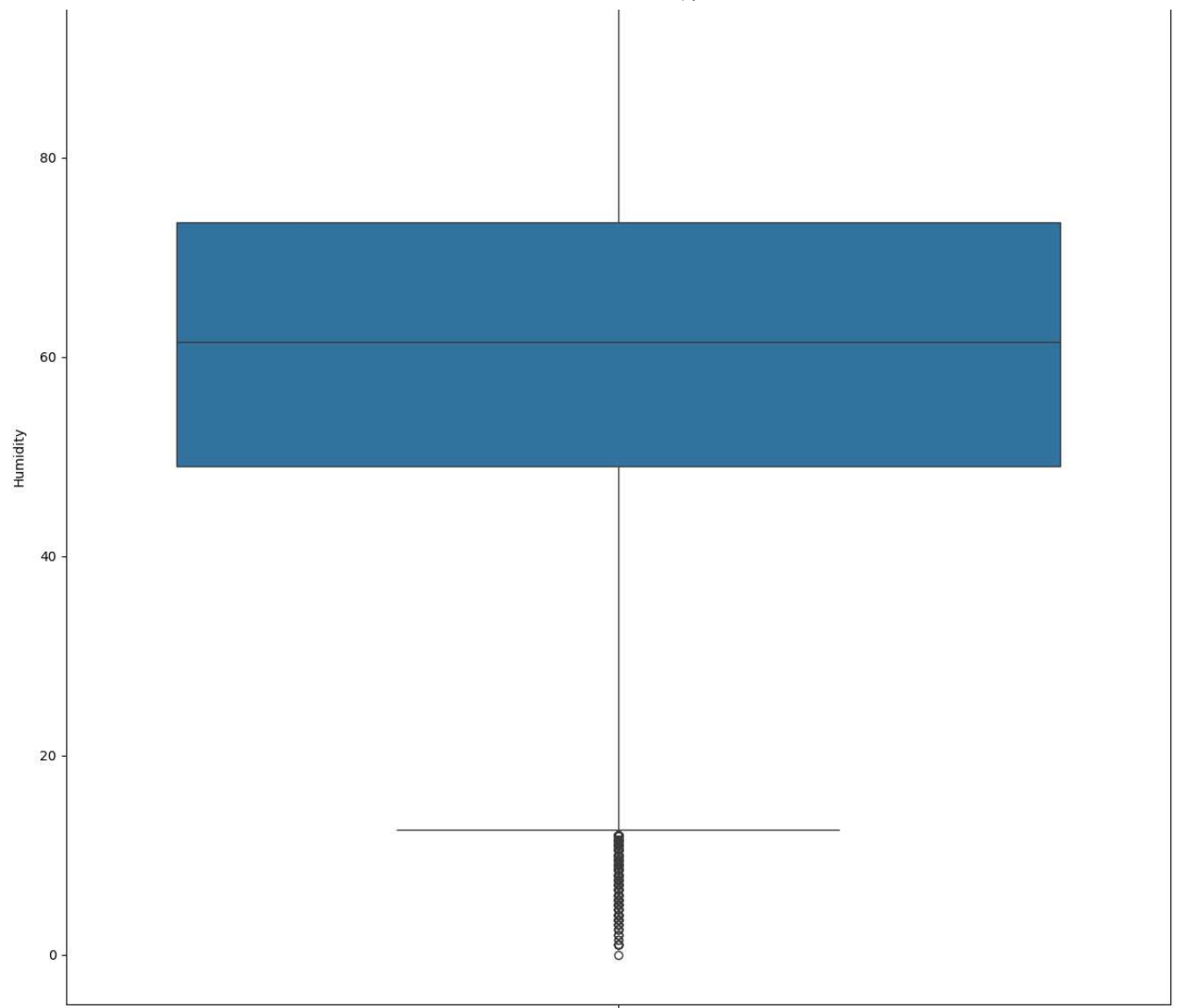
```
sns.distplot(df[feature])
```



<Figure size 1500x1500 with 0 Axes>

```
1 #A for loop is used to plot a boxplot for all the continuous features to see the outliers
2 for feature in continuous_feature:
3     data=df.copy()
4     sns.boxplot(data[feature])
5     plt.title(feature)
6     plt.figure(figsize=(15,15))
```






<Figure size 1500x1500 with 0 Axes>

```
1 for feature in continuous_feature:
2     if(df[feature].isnull().sum()*100/len(df))>0:
3         df[feature] = df[feature].fillna(df[feature].median())
```

```
1 df.isnull().sum()*100/len(df)
```



	0
Date	0.000000
Location	0.000000
Temp	0.000000
WindSpeed	0.000000
Humidity	0.000000
RainToday	2.241853
RainTomorrow	2.245978


dtype: float64

```
1 discrete_feature
```



[]


```
1 def mode_nan(df,variable):
2     mode=df[variable].value_counts().index[0]
3     df[variable].fillna(mode,inplace=True)
4
1 df["RainToday"] = pd.get_dummies(df["RainToday"], drop_first = True)
2 df["RainTomorrow"] = pd.get_dummies(df["RainTomorrow"], drop_first = True)
3 df
```



	Date	Location	Temp	WindSpeed	Humidity	RainToday	RainTomorrow
0	2008-12-01	Albury	18.15	22.0	46.5	False	False
1	2008-12-02	Albury	16.25	13.0	34.5	False	False
2	2008-12-03	Albury	19.30	22.5	34.0	False	False
3	2008-12-04	Albury	18.60	10.0	30.5	False	False
4	2008-12-05	Albury	24.90	13.5	57.5	False	False
...
145455	2017-06-21	Uluru	13.10	12.0	37.5	False	False
145456	2017-06-22	Uluru	14.45	11.0	38.5	False	False
145457	2017-06-23	Uluru	16.15	9.0	38.5	False	False
145458	2017-06-24	Uluru	17.40	10.0	37.5	False	False
145459	2017-06-25	Uluru	14.90	17.0	49.0	False	False

145460 rows x 7 columns

```
1 for feature in categorical_feature:
2     print(feature, (df.groupby([feature])["RainTomorrow"].mean().sort_values(ascending = False)).index)
```



Date	Index(['2007-11-01', '2007-12-15', '2008-02-03', '2008-01-31', '2008-01-30', '2008-01-19', '2008-01-18', '2008-01-16', '2008-01-12', '2007-11-30', ... '2008-05-08', '2008-05-09', '2008-01-03', '2008-01-02', '2008-05-10', '2008-05-11', '2008-05-12', '2008-05-13', '2008-05-14', '2007-12-25'], dtype='object', name='Date', length=3436)
Location	Index(['Portland', 'Walpole', 'Cairns', 'Dartmoor', 'NorfolkIsland',


```
'MountGambier', 'Albany', 'Witchcliffe', 'CoffsHarbour', 'MountGinini',  
'NorahHead', 'Darwin', 'Sydney', 'SydneyAirport', 'Ballarat',  
'GoldCoast', 'Watsonia', 'Newcastle', 'Hobart', 'Wollongong',  
'Williamstown', 'Launceston', 'Brisbane', 'MelbourneAirport', 'Adelaide',  
'Sale', 'Albury', 'Perth', 'Melbourne', 'Nuriootpa', 'Penrith',  
'BadgerysCreek', 'PerthAirport', 'Tuggeranong', 'Richmond', 'Bendigo',  
'Canberra', 'WaggaWagga', 'Townsville', 'Katherine', 'PearceRAAF',  
'SalmonGums', 'Nhil', 'Moree', 'Cobar', 'Mildura', 'AliceSprings',  
'Uluru', 'Woomera'],  
dtype='object', name='Location')  
RainToday Index([True, False], dtype='bool', name='RainToday')  
RainTomorrow Index([True, False], dtype='bool', name='RainTomorrow')
```

```
1 df1 = df.groupby(["Location"])[ "RainTomorrow"].value_counts().sort_values().unstack()
```

```
1 df1
```



RainTomorrow	False	True
Location		
Adelaide	2505	688
Albany	2138	902
Albury	2422	618
AliceSprings	2796	244
BadgerysCreek	2426	583
Ballarat	2259	781
Bendigo	2478	562
Brisbane	2484	709
Cairns	2090	950
Canberra	2807	629
Cobar	2623	386
CoffsHarbour	2140	869
Dartmoor	2087	922
Darwin	2341	852
GoldCoast	2265	775
Hobart	2432	761
Katherine	1313	265
Launceston	2341	699
Melbourne	2557	636
MelbourneAirport	2356	653
Mildura	2682	327
Moree	2615	394
MountGambier	2120	920
MountGinini	2221	819
Newcastle	2308	731
Nhil	1336	242
NorahHead	2196	808
NorfolkIsland	2090	919
Nuriootpa	2417	592
PearceRAAF	2504	505
Penrith	2444	595
Perth	2548	645
PerthAirport	2442	567
Portland	1914	1095
Richmond	2449	560
Sale	2366	643
SalmonGums	2529	472
Sydney	2479	865
SydneyAirport	2235	774
Townsville	2521	519
Tuggeranong	2471	568



Uluru	1462	116
WaggaWagga	2473	536
Walpole	2057	949
Watsonia	2271	738
Williamtown	2309	700
Witchcliffe	2130	879
Wollongong	2327	713
Woomera	2807	202

Next steps:

Generate code with df1

View recommended plots

New interactive sheet

```
1 location = {'Portland':1, 'Cairns':2, 'Walpole':3, 'Dartmoor':4, 'MountGambier':5,
2             'NorfolkIsland':6, 'Albany':7, 'Witchcliffe':8, 'CoffsHarbour':9, 'Sydney':10,
3             'Darwin':11, 'MountGinini':12, 'NorahHead':13, 'Ballarat':14, 'GoldCoast':15,
4             'SydneyAirport':16, 'Hobart':17, 'Watsonia':18, 'Newcastle':19, 'Wollongong':20,
5             'Brisbane':21, 'Williamtown':22, 'Launceston':23, 'Adelaide':24, 'MelbourneAirport':25,
6             'Perth':26, 'Sale':27, 'Melbourne':28, 'Canberra':29, 'Albury':30, 'Penrith':31,
7             'Nuriootpa':32, 'BadgerysCreek':33, 'Tuggeranong':34, 'PerthAirport':35, 'Bendigo':36,
8             'Richmond':37, 'WaggaWagga':38, 'Townsville':39, 'PearceRAAF':40, 'SalmonGums':41,
9             'Moree':42, 'Cobar':43, 'Mildura':44, 'Katherine':45, 'AliceSprings':46, 'Nhil':47,
10            'Woomera':48, 'Uluru':49}
11 df["Location"] = df["Location"].map(location)

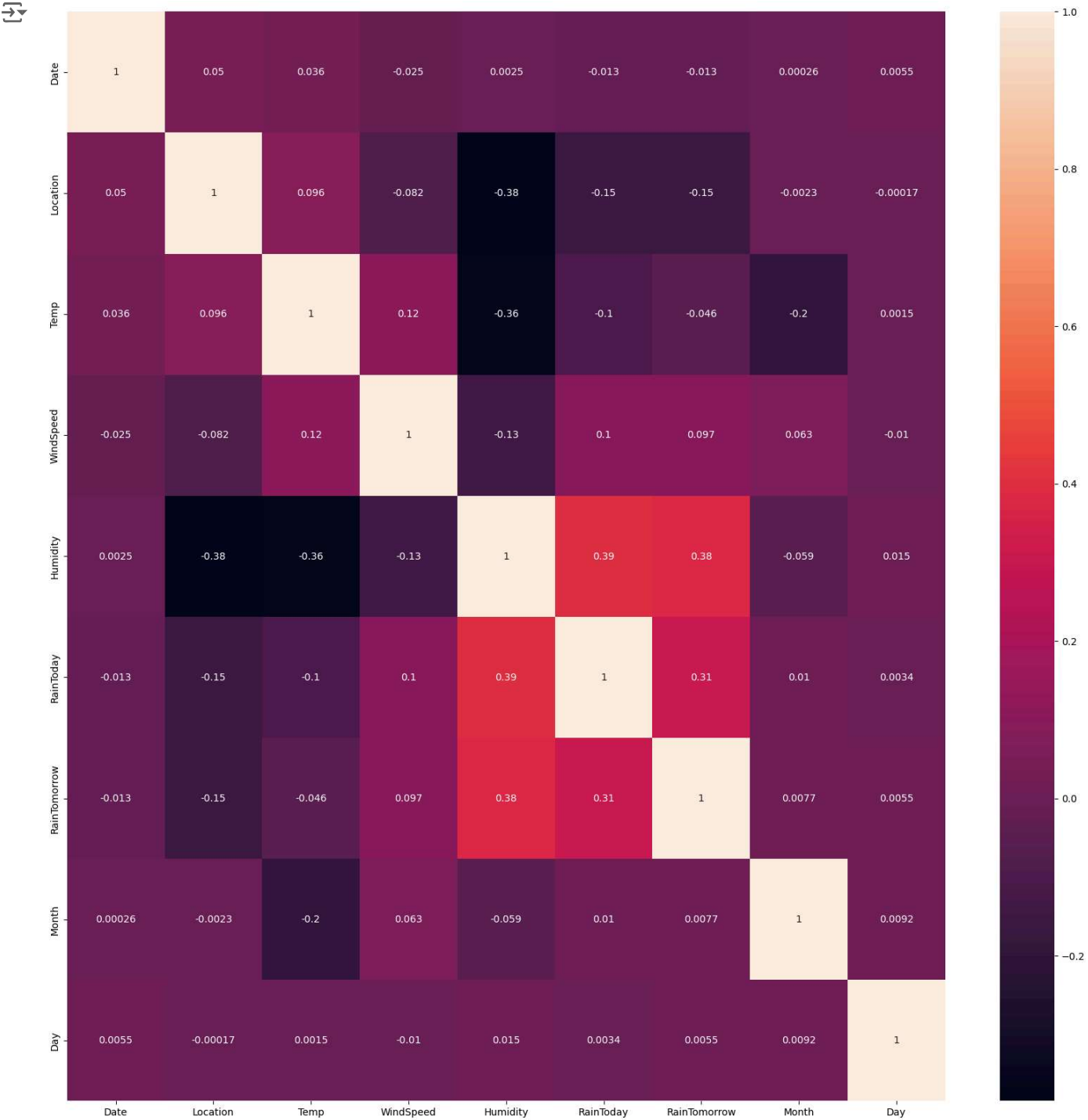
1 df['Date'] = pd.to_datetime(df['Date'])
2 df['Month'] = df['Date'].dt.month
3 df['Day'] = df['Date'].dt.day
4 df
```

	Date	Location	Temp	WindSpeed	Humidity	RainToday	RainTomorrow	Month	Day
0	2008-12-01	30	18.15	22.0	46.5	False	False	12	1
1	2008-12-02	30	16.25	13.0	34.5	False	False	12	2
2	2008-12-03	30	19.30	22.5	34.0	False	False	12	3
3	2008-12-04	30	18.60	10.0	30.5	False	False	12	4
4	2008-12-05	30	24.90	13.5	57.5	False	False	12	5
...
145455	2017-06-21	49	13.10	12.0	37.5	False	False	6	21
145456	2017-06-22	49	14.45	11.0	38.5	False	False	6	22
145457	2017-06-23	49	16.15	9.0	38.5	False	False	6	23
145458	2017-06-24	49	17.40	10.0	37.5	False	False	6	24
145459	2017-06-25	49	14.90	17.0	49.0	False	False	6	25

145460 rows × 9 columns



```
1 corrmat = df.corr()
2 plt.figure(figsize=(20,20))
3 #plot heat map
4 g=sns.heatmap(corrmat,annot=True)
```



1 df



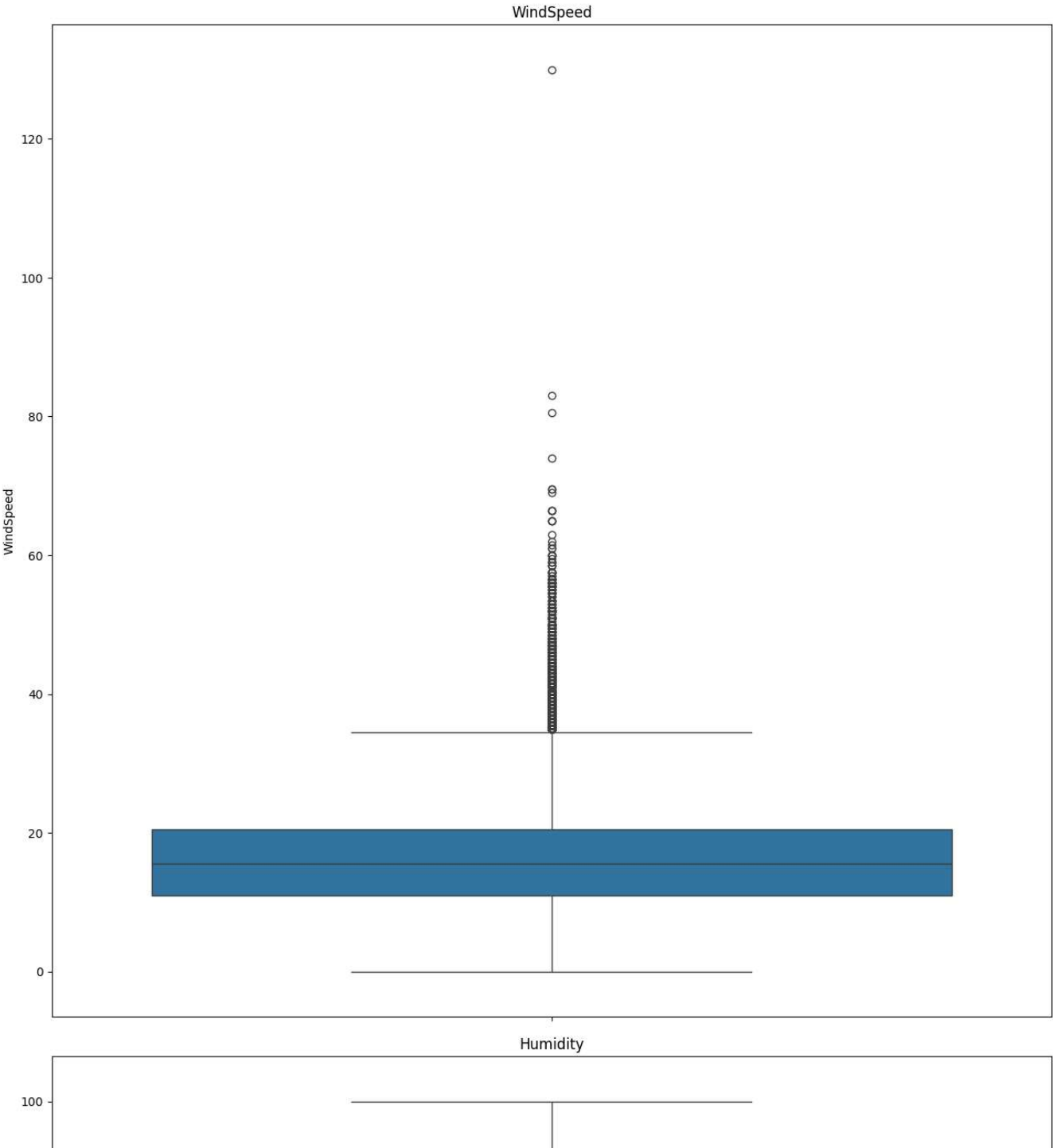
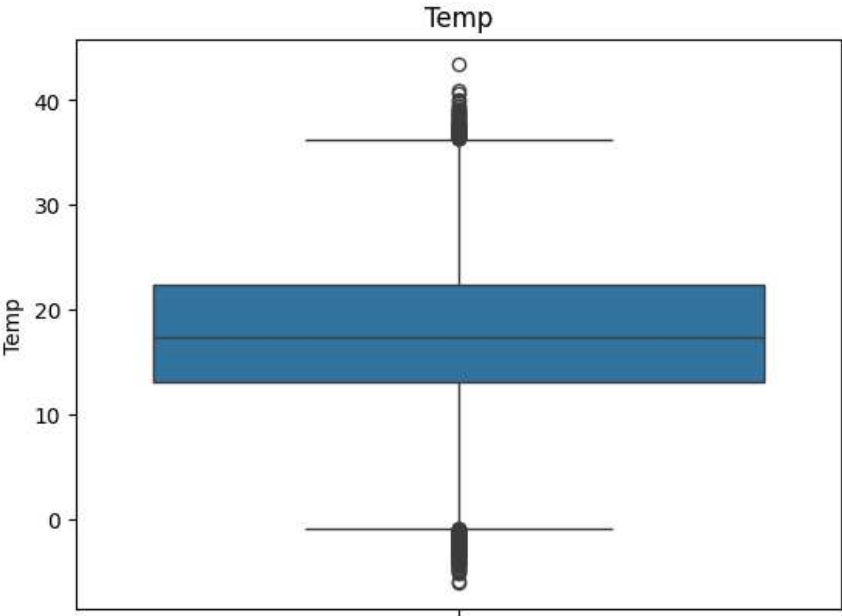
	Date	Location	Temp	WindSpeed	Humidity	RainToday	RainTomorrow	Month	Day	
0	2008-12-01	30	18.15	22.0	46.5	False	False	12	1	
1	2008-12-02	30	16.25	13.0	34.5	False	False	12	2	
2	2008-12-03	30	19.30	22.5	34.0	False	False	12	3	
3	2008-12-04	30	18.60	10.0	30.5	False	False	12	4	
4	2008-12-05	30	24.90	13.5	57.5	False	False	12	5	
...	
145455	2017-06-21	49	13.10	12.0	37.5	False	False	6	21	
145456	2017-06-22	49	14.45	11.0	38.5	False	False	6	22	
145457	2017-06-23	49	16.15	9.0	38.5	False	False	6	23	
145458	2017-06-24	49	17.40	10.0	37.5	False	False	6	24	
145459	2017-06-25	49	14.90	17.0	49.0	False	False	6	25	

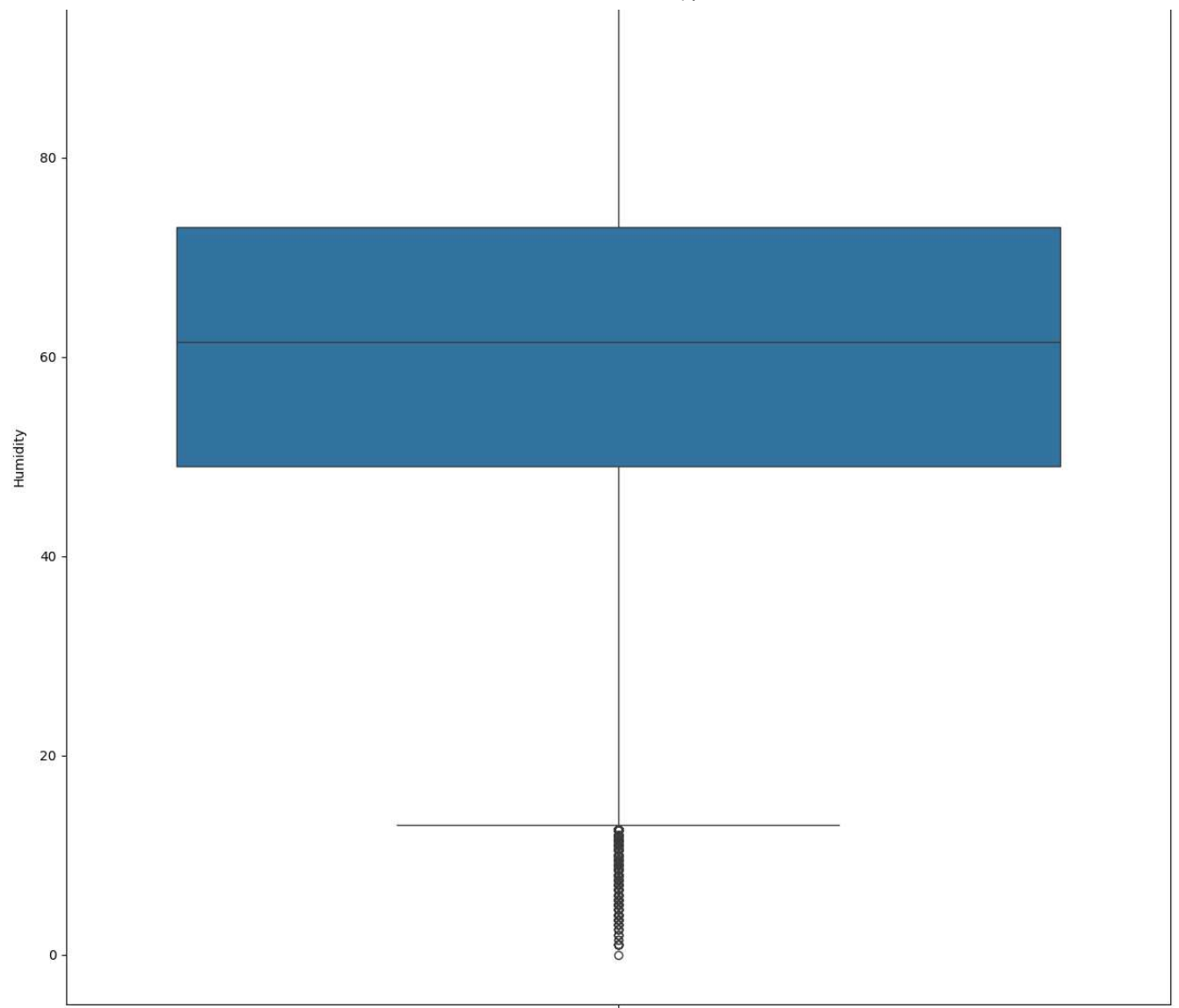
145460 rows × 9 columns

```

1 for feature in continuous_feature:
2     data=df.copy()
3     sns.boxplot(data[feature])
4     plt.title(feature)
5     plt.figure(figsize=(15,15))

```





<Figure size 1500x1500 with 0 Axes>

```
1 for feature in continuous_feature:
2     print(feature)
```

Temp
WindSpeed
Humidity

```
1 IQR=df.Temp.quantile(0.75)-df.Temp.quantile(0.25)
2 lower_bridge=df.Temp.quantile(0.25)-(IQR*1.5)
3 upper_bridge=df.Temp.quantile(0.75)+(IQR*1.5)
4 print(lower_bridge, upper_bridge)
```

-0.9500000000000011 36.25

```
1 df.loc[df['Temp']>=36.975, 'Temp']=36.975
2 df.loc[df['Temp']<=2.7, 'Temp']=2.7
```

```
1 IQR=df.WindSpeed.quantile(0.75)-df.WindSpeed.quantile(0.25)
2 lower_bridge=df.WindSpeed.quantile(0.25)-(IQR*1.5)
3 upper_bridge=df.WindSpeed.quantile(0.75)+(IQR*1.5)
4 print(lower_bridge, upper_bridge)
```

-3.25 34.75

```
1 df.loc[df['WindSpeed']>=38.75, 'WindSpeed']=38.75
2 df.loc[df['WindSpeed']<=-7.25, 'WindSpeed']=-7.25
```

```
1 IQR=df.Humidity.quantile(0.75)-df.Humidity.quantile(0.25)
2 lower_bridge=df.Humidity.quantile(0.25)-(IQR*1.5)
3 upper_bridge=df.Humidity.quantile(0.75)+(IQR*1.5)
4 print(lower_bridge, upper_bridge)
```

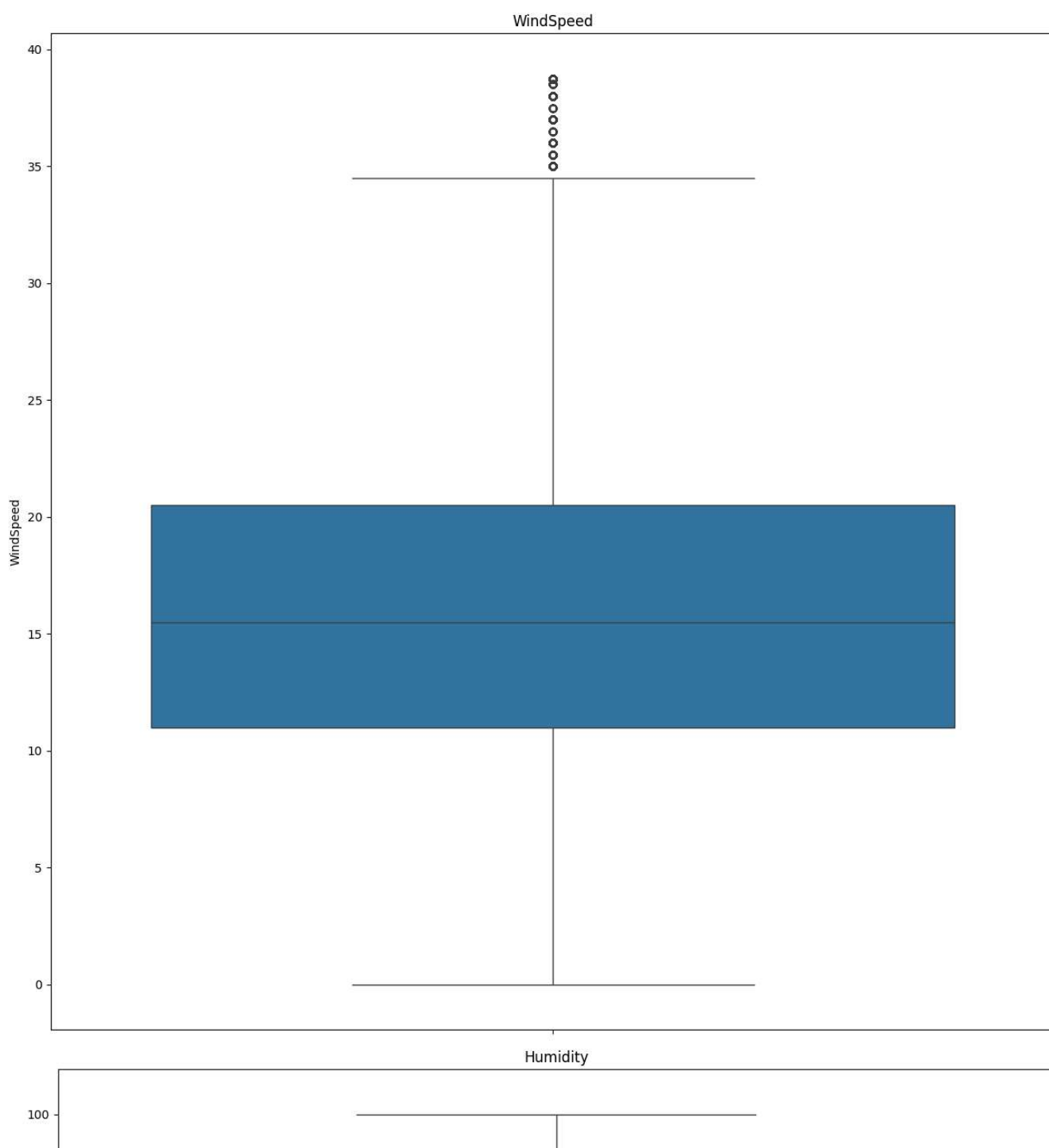
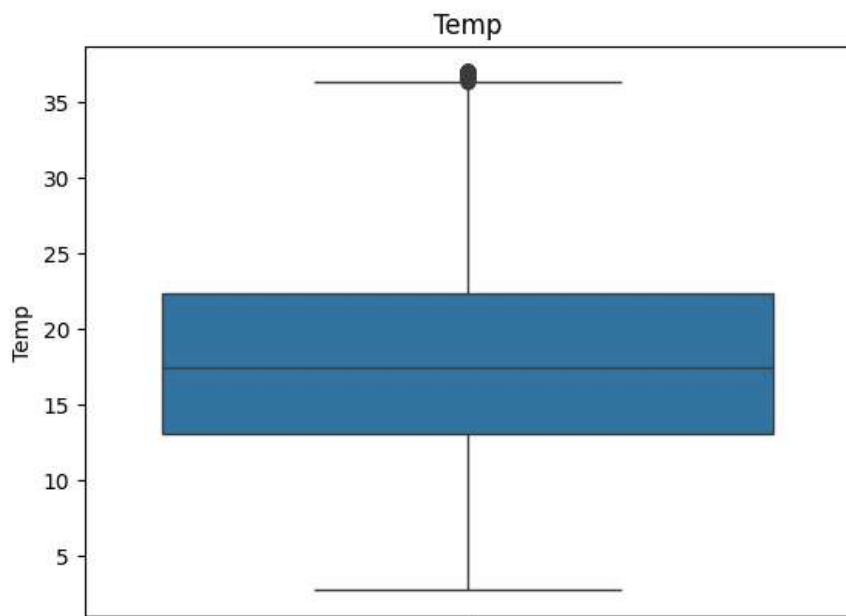
13.0 109.0

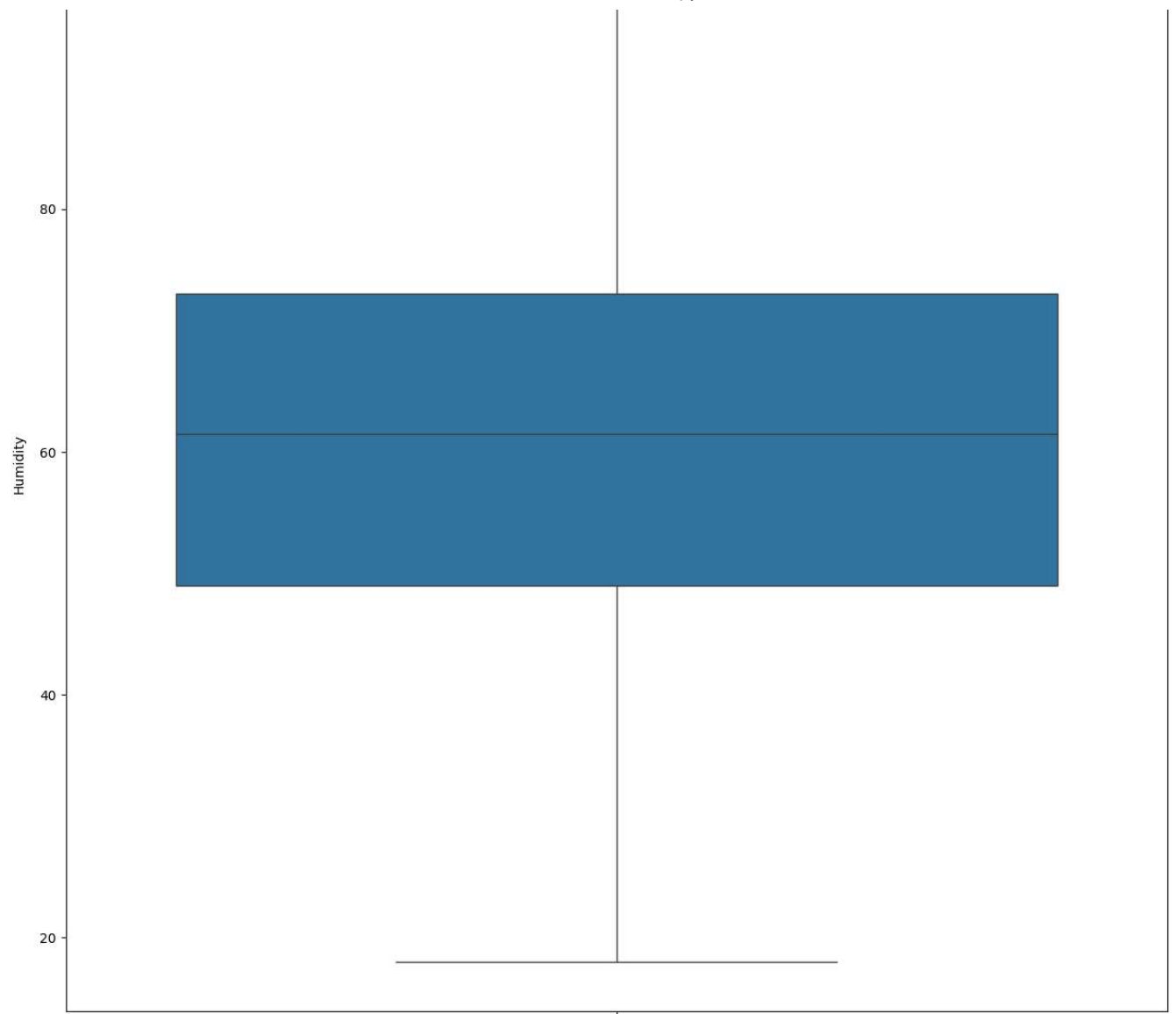
```
1 df.loc[df['Humidity']>=122, 'Humidity']=122
2 df.loc[df['Humidity']<=18, 'Humidity']=18
```

```
1 '''df.loc[df['Temp3pm']>=40.45, 'Temp3pm']=40.45
2 df.loc[df['Temp3pm']<=2.45, 'Temp3pm']=2.45'''
```

'df.loc[df['Temp3pm']>=40.45, 'Temp3pm']=40.45\ndf.loc[df['Temp3pm']<=2.45, 'Temp3pm']=2.45'

```
1 for feature in continuous_feature:
2     data=df.copy()
3     sns.boxplot(data[feature])
4     plt.title(feature)
5     plt.figure(figsize=(15,15))
```



<Figure size 1500x1500 with 0 Axes>

```
1 def qq_plots(df, variable):
2     plt.figure(figsize=(15,6))
3     plt.subplot(1, 2, 1)
4     df[variable].hist()
5     plt.subplot(1, 2, 2)
6     stats.probplot(df[variable], dist="norm", plot=plt)
7     plt.show()

1 for feature in continuous_feature:
2     print(feature)
3     plt.figure(figsize=(15,6))
4     plt.subplot(1, 2, 1)
5     df[feature].hist()
6     plt.subplot(1, 2, 2)
7     stats.probplot(df[feature], dist="norm", plot=plt)
8     plt.show()
```