

INFOSYS SPRINGBOARD VIRTUAL INTERNSHIP

FWI PREDICTOR

By

Srinanda C S

MILESTONE 1

Dataset Source : Kaggle

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder
```

These libraries are used for data analysis, preprocessing, visualization, and encoding.

```
df = pd.read_csv("FWI Dataset.csv")

if 'Region' in df.columns:

    df['Region'] = df['Region'].astype('category').cat.codes
```

Dataset is loaded using pandas.

The following information is displayed:

- Entire dataset
- Dataset structure (df.info())
- Statistical summary (df.describe())
- First & last 5 rows

Converts the *Region* column from strings to numerical category codes. Useful for machine learning algorithms that accept numeric inputs.

```
numeric_df = df.select_dtypes(include=['int64', 'float64'])
```

Extracts only numerical features for later analysis.

```
print(df.shape)

print(df.columns)
```

To ensure the dataset structure is correct.

```
df.isnull().sum()

df[df.isnull().any(axis=1)]

df.columns = df.columns.str.strip()
```

Identifies missing values in each column.

Displays rows containing incomplete data.

Removes extra spaces in column names (common in raw CSV files).

```
for col in df.columns:

    if df[col].dtype == 'object':

        df[col] = df[col].astype(str).str.strip()
```

Removes unnecessary spaces in string values.

Ensures uniform data format.

```
for col in df.columns:

    if df[col].dtype == 'object':

        df[col] = df[col].str.replace(" ", " ")

        if df[col].str.contains(" ").any():

            df[col] = df[col].str.split(" ").str[0]
```

This ensures numeric columns convert cleanly.

```
numeric_cols =  
['Temperature', 'RH', 'Ws', 'Rain', 'FFMC', 'DMC', 'DC', 'ISI', 'BUI',  
'FWI']
```

```
for col in numeric_cols:
```

```
    df[col] = pd.to_numeric(df[col], errors='coerce')
```

Converts corrupted strings into numeric format.

Non-convertible values become Nan.

```
df['Region'] = df['Region'].fillna(df['Region'].mode()[0])
```

```
df['Classes'] = df['Classes'].fillna(df['Classes'].mode()[0])
```

Uses mode (most frequent value) for categorical features.

Prevents ML models from failing due to null values.

Label Encoding Categorical Columns

```
le_region = LabelEncoder()
```

```
df['Region_encoded'] = le_region.fit_transform(df['Region'])
```

```
le_class = LabelEncoder()
```

```
df['Classes_encoded'] = le_class.fit_transform(df['Classes'])
```

Convert string labels to numeric classes for ML model training.

Encoding All Remaining Categorical Columns

```
df_encoded = df.copy()
```

```
label_encoders = {}
```

```
for col in df_encoded.columns:

    if df_encoded[col].dtype == 'object':

        le = LabelEncoder()

        df_encoded[col] =
le.fit_transform(df_encoded[col].astype(str))

        label_encoders[col] = le
```

This ensures all non-numeric features are usable in correlation analysis.

Correlation Heatmap

```
plt.figure(figsize=(10,8))

sns.heatmap(numeric_df.corr(), annot=True)

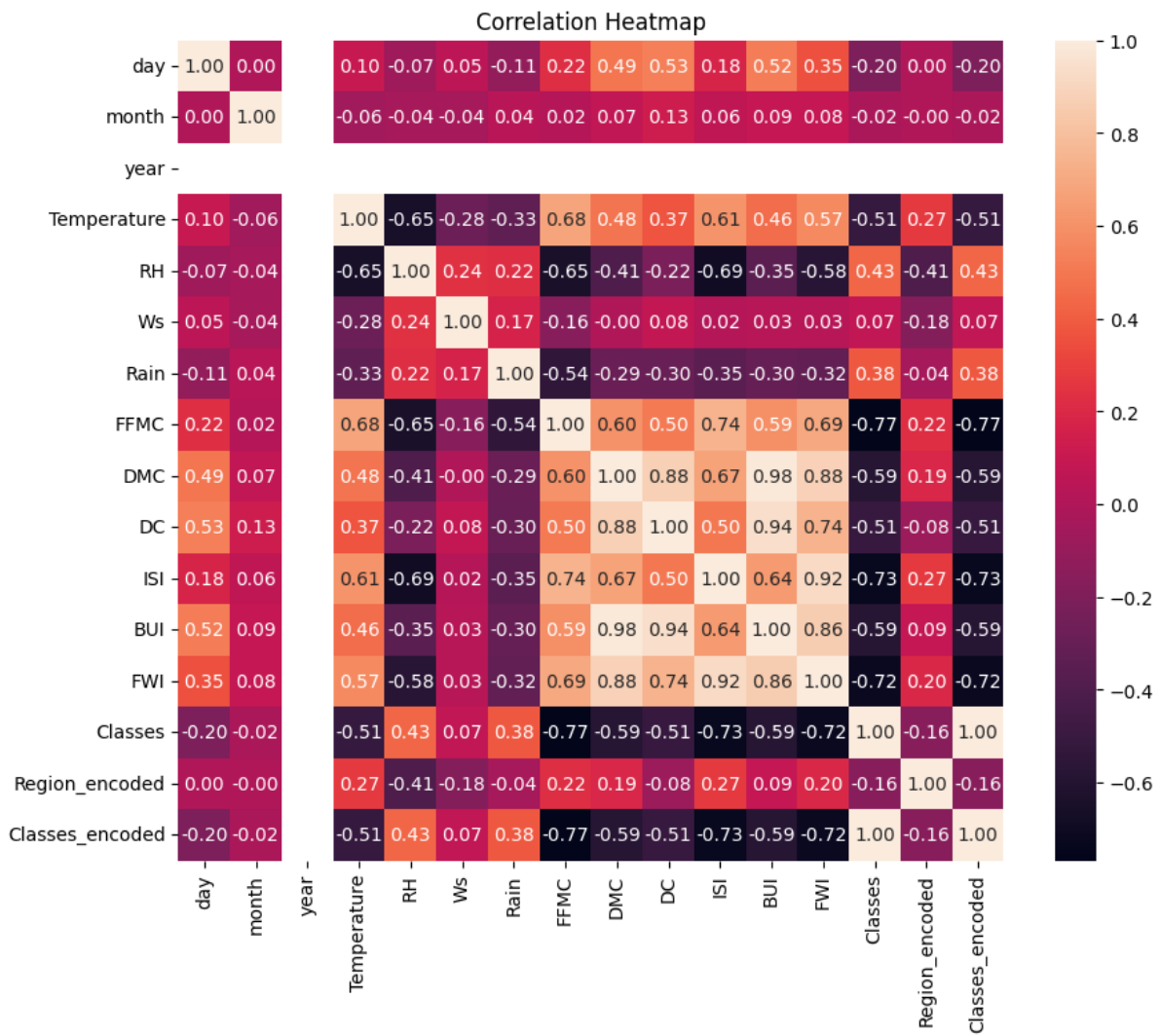
plt.show()
```

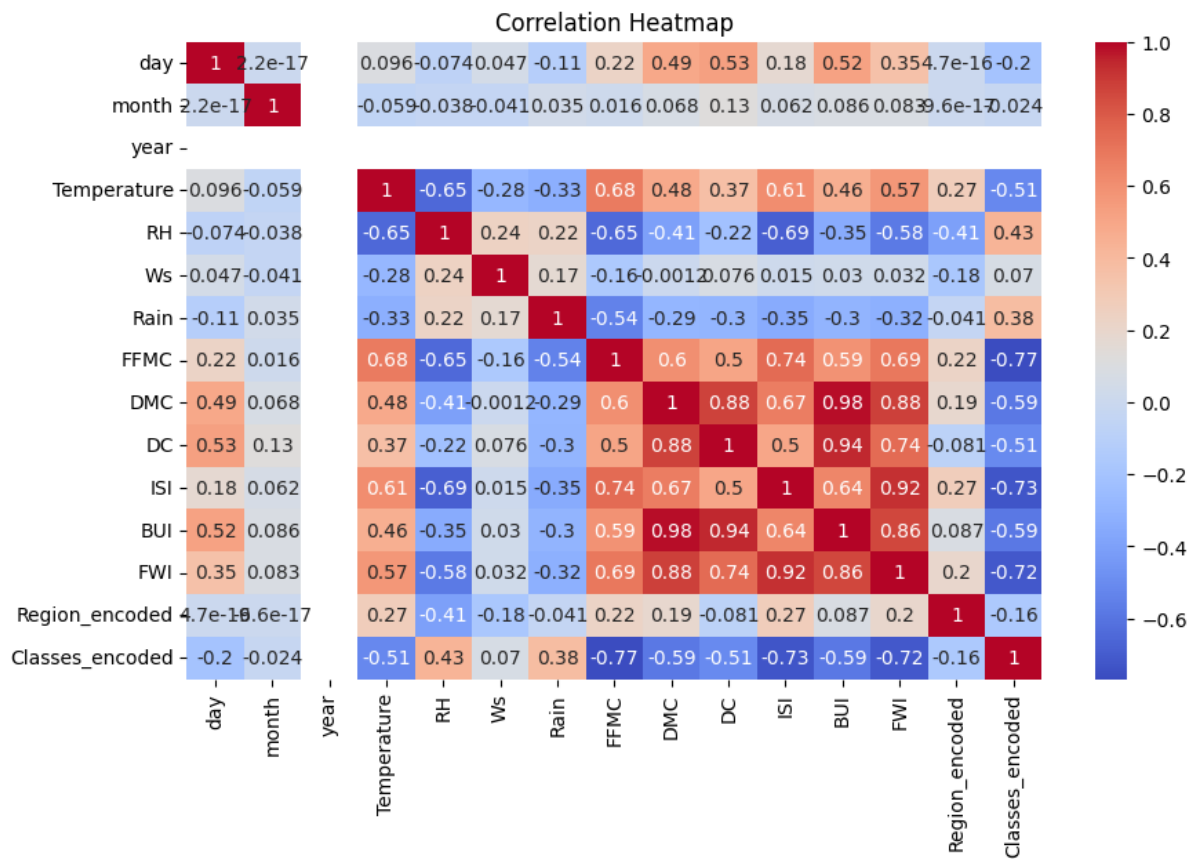
Used to understand:

Feature relationships

Which variables strongly influence FWI

Multicollinearity issues





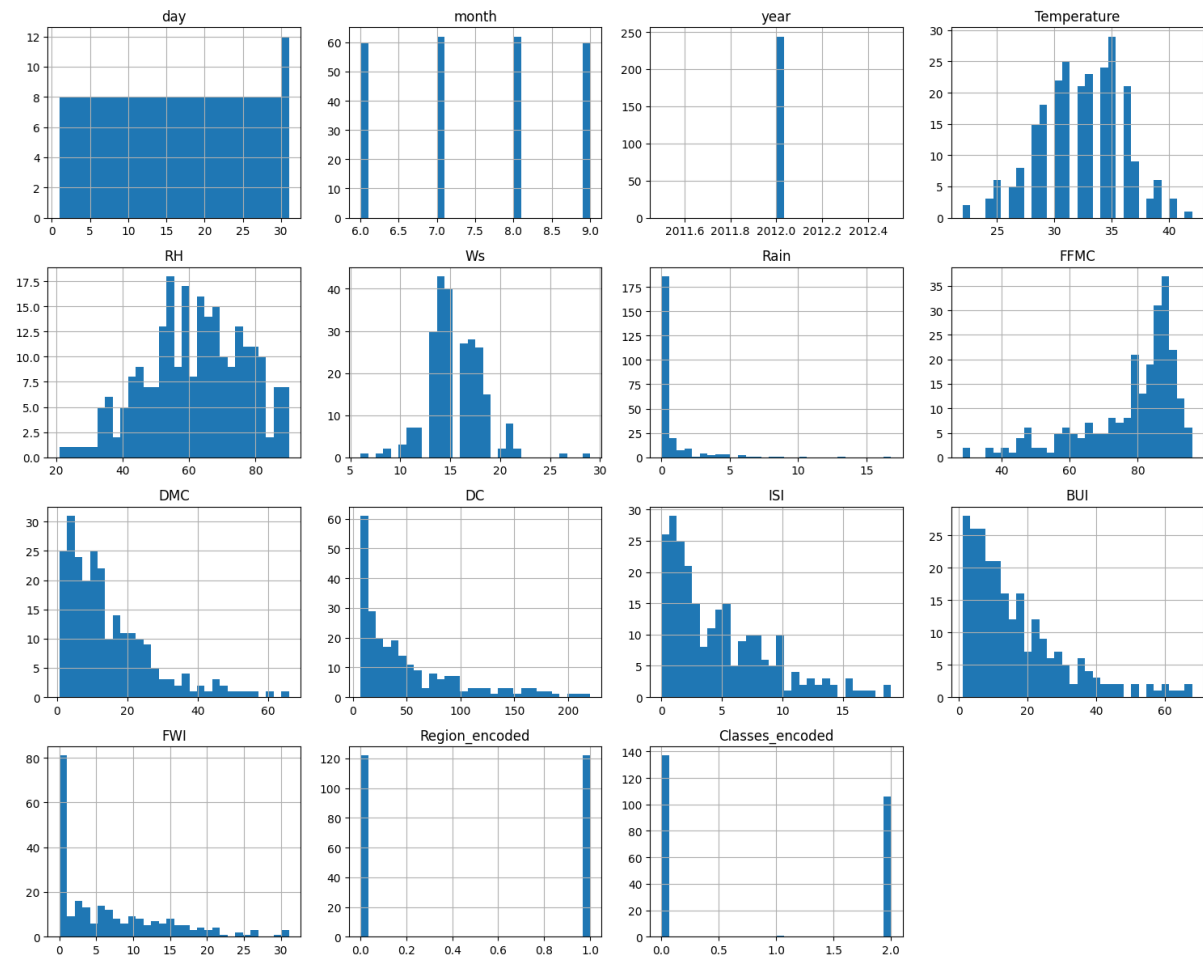
Histogram Plots

```
numeric_df.hist(figsize=(15, 12), bins=30)
```

```
plt.show()
```

Shows distribution of each numeric feature:

Normal, Skewed, Outliers



Density (KDE) Plots

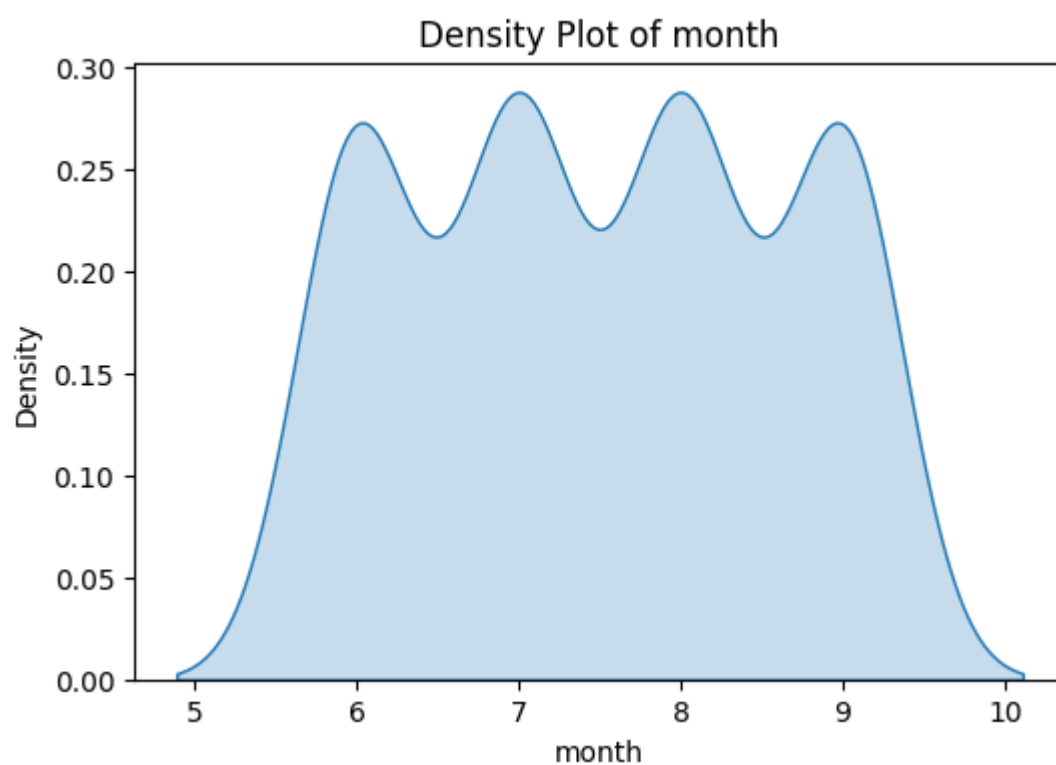
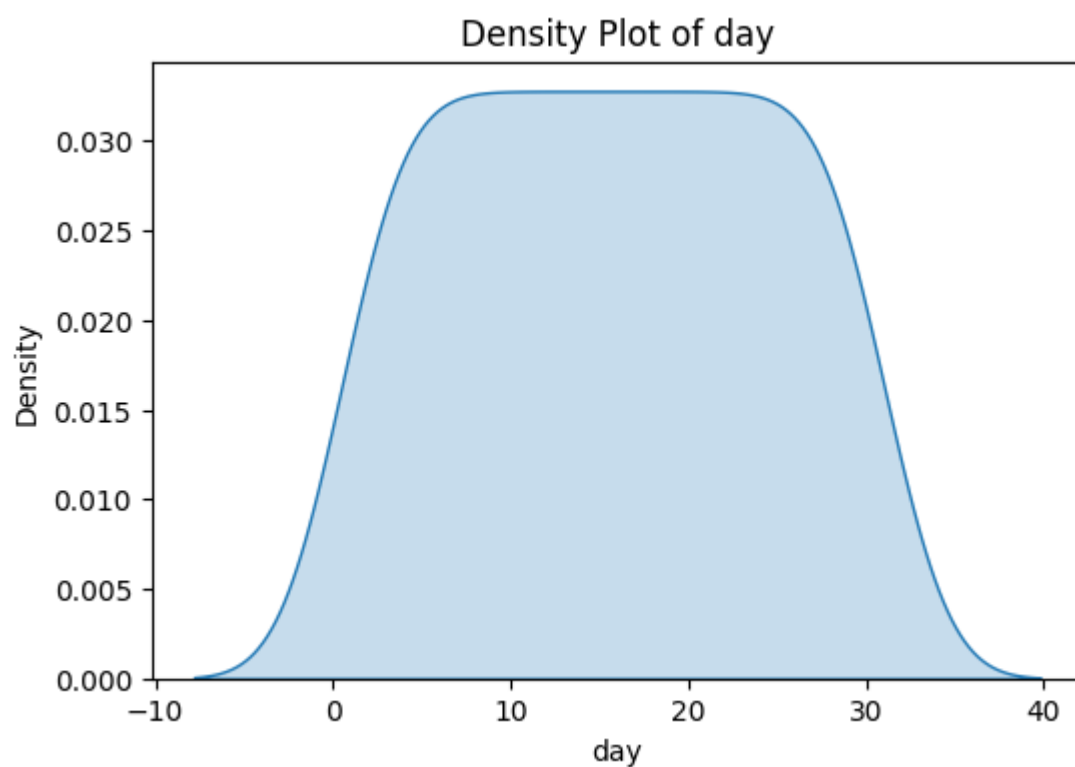
```
sns.kdeplot(numeric_df[col], fill=True)
```

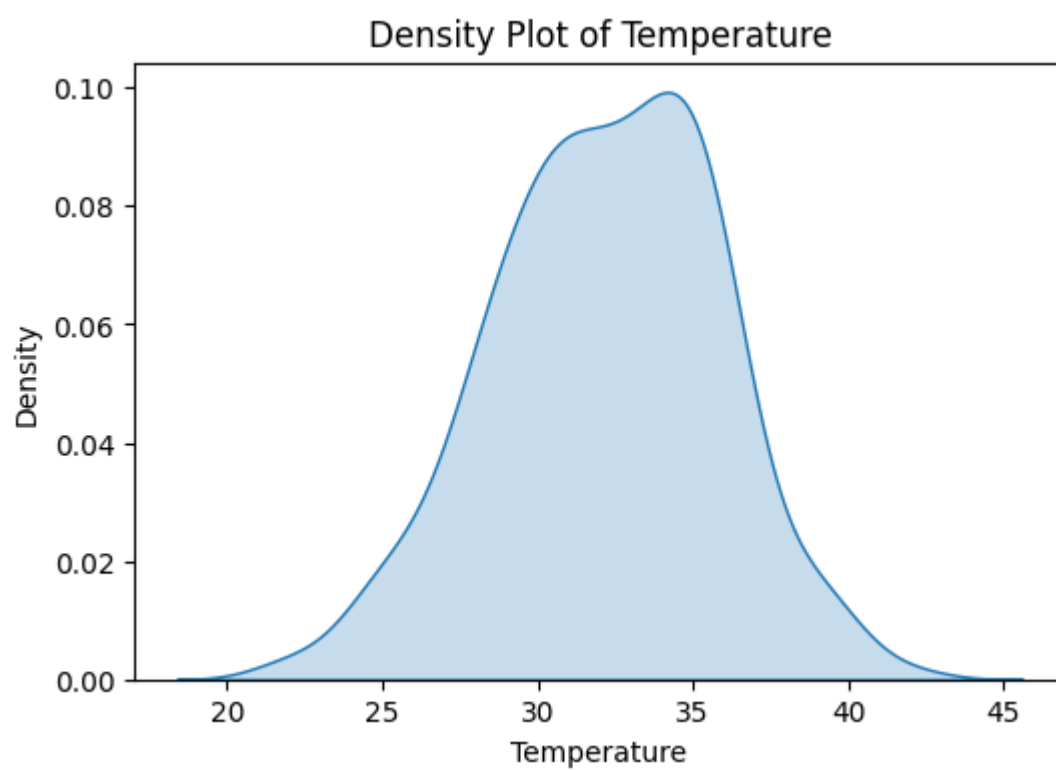
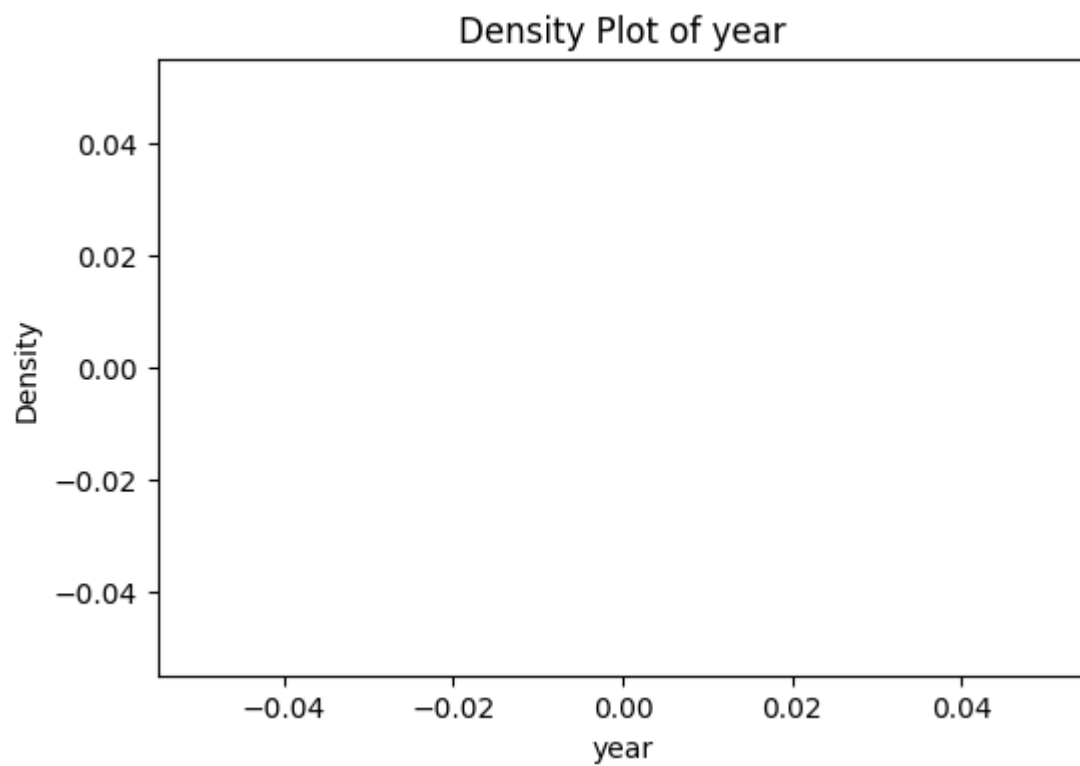
These help understand:

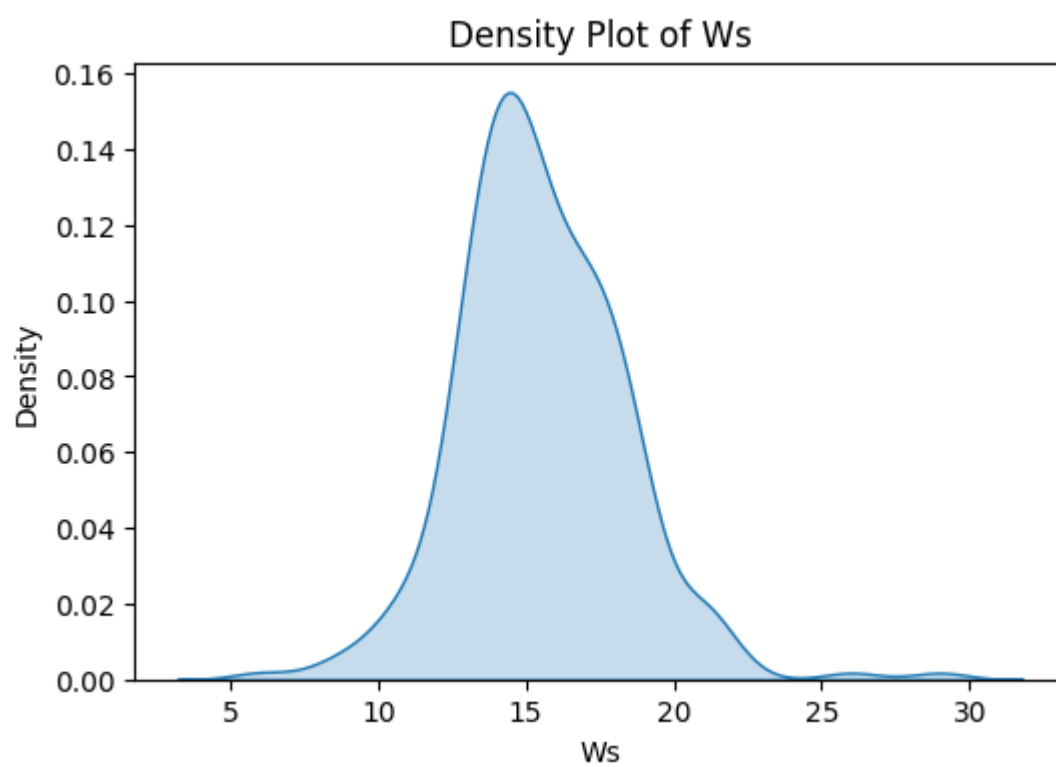
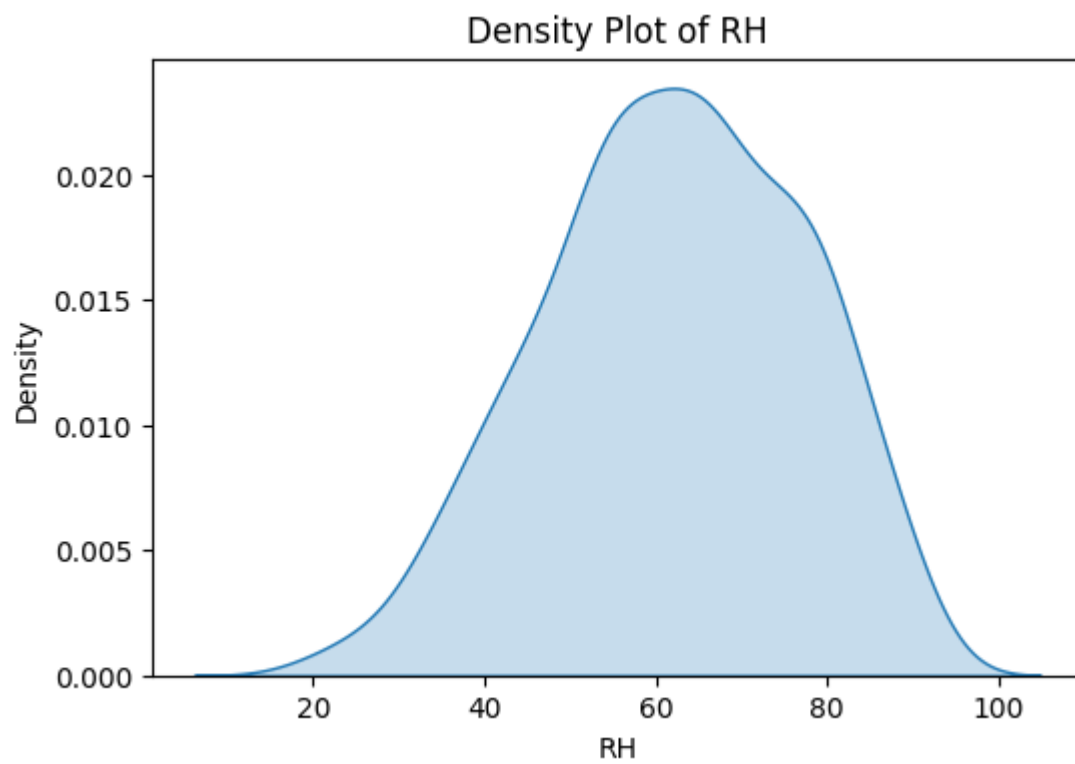
Probability distribution

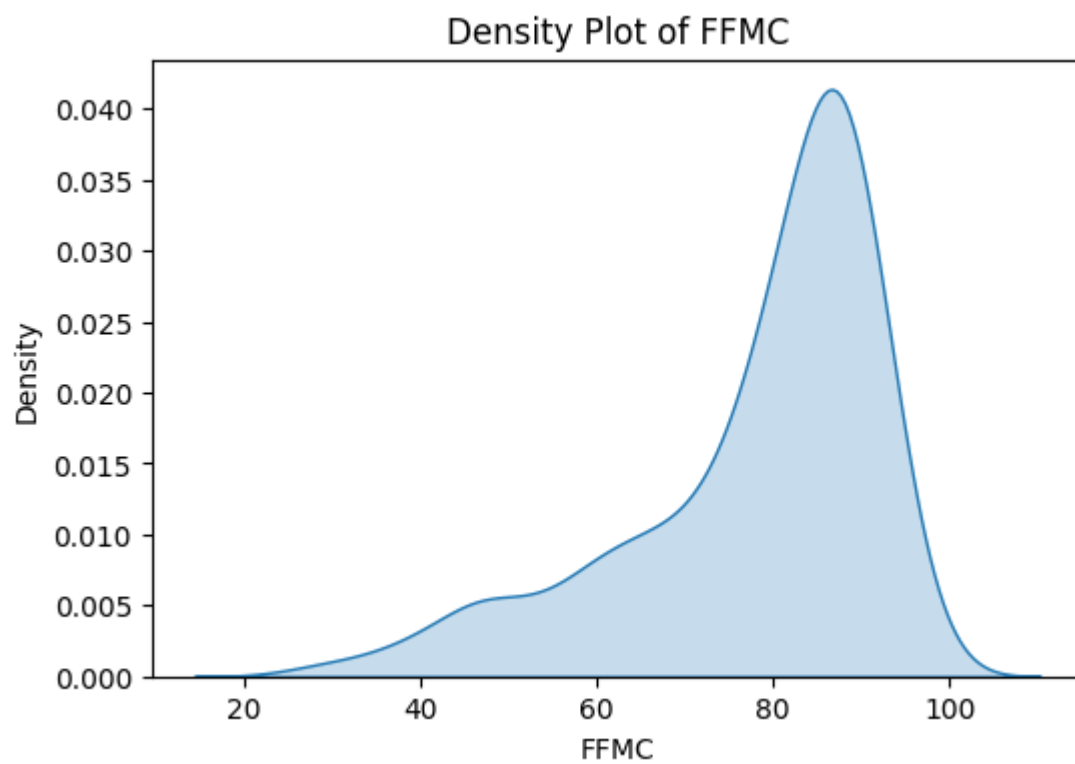
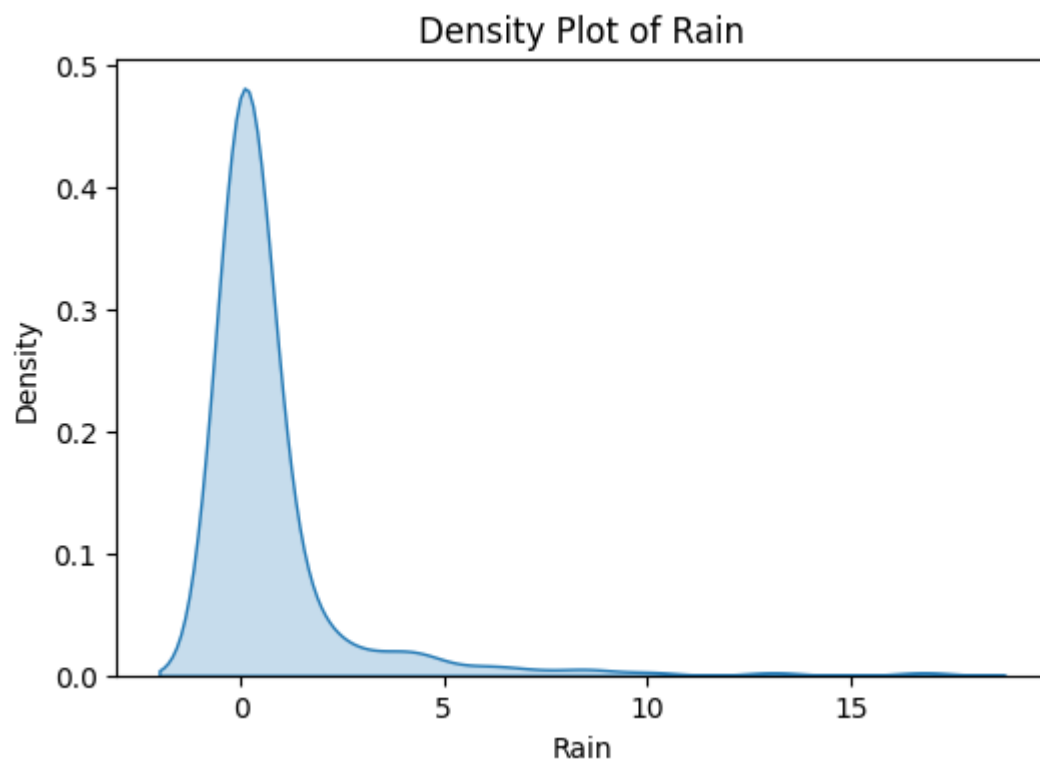
Spread of data

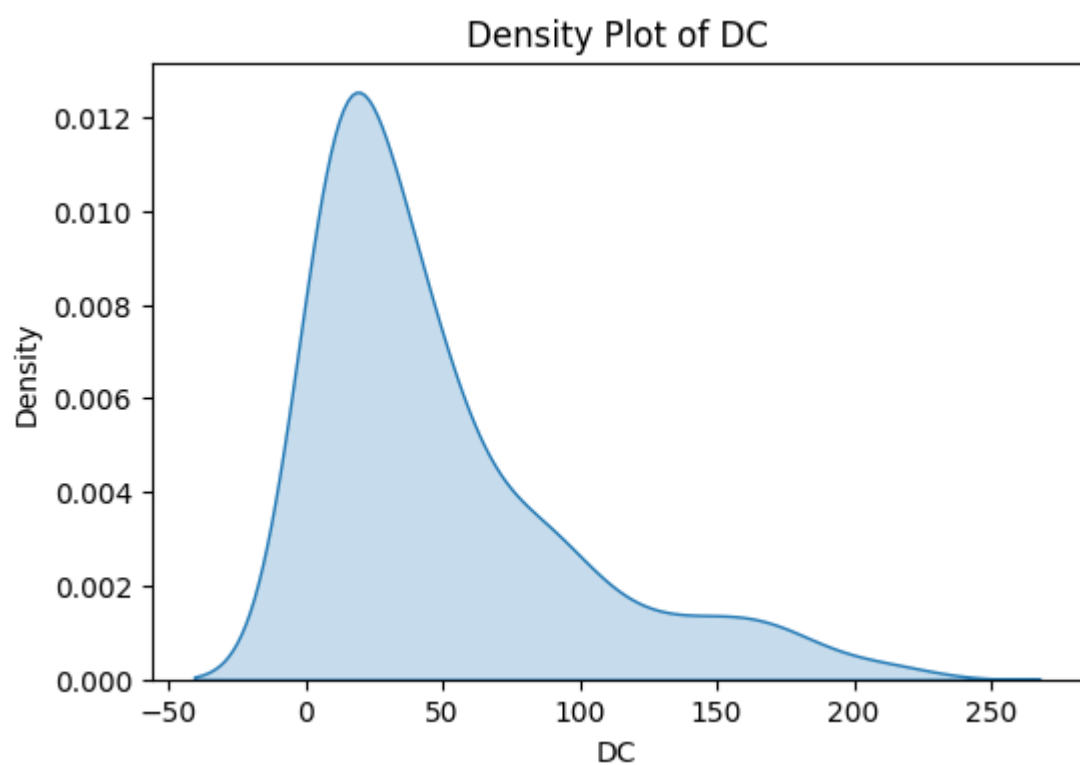
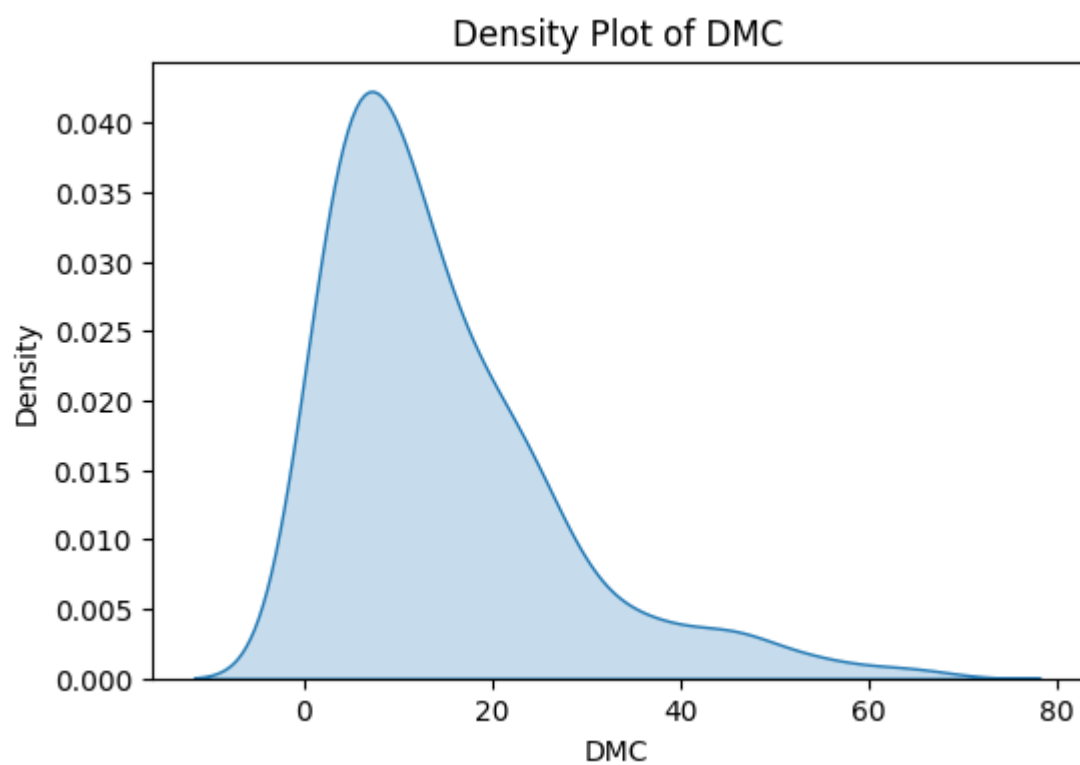
Detecting skewness

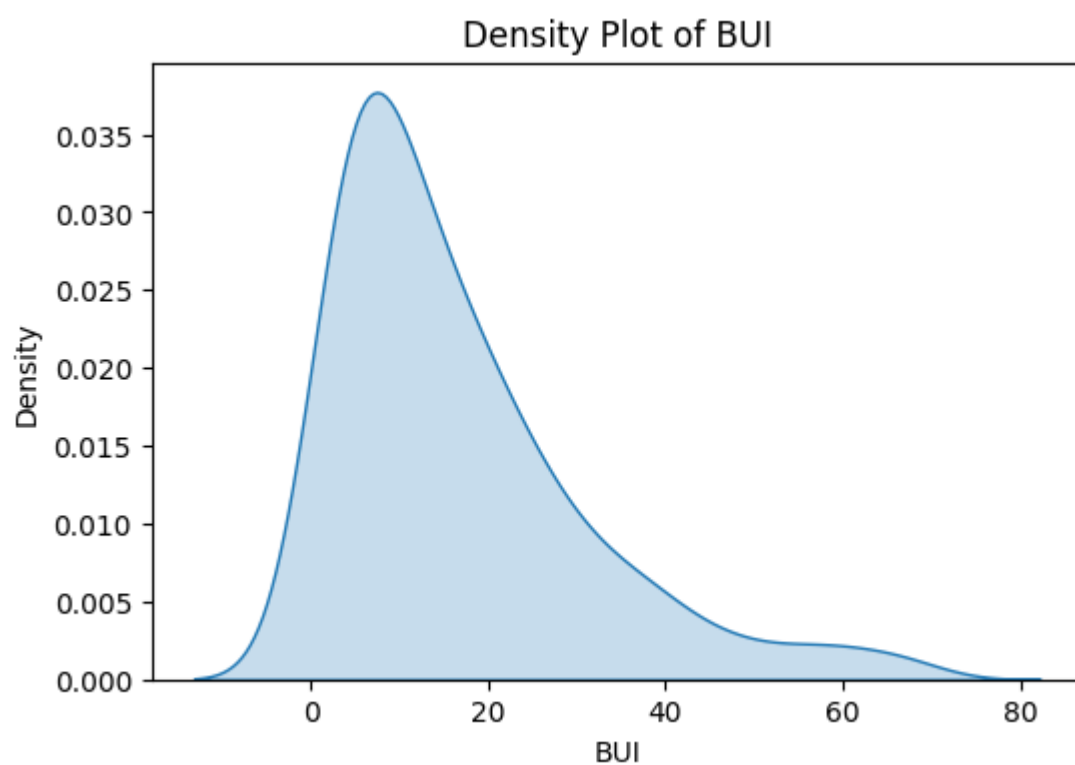
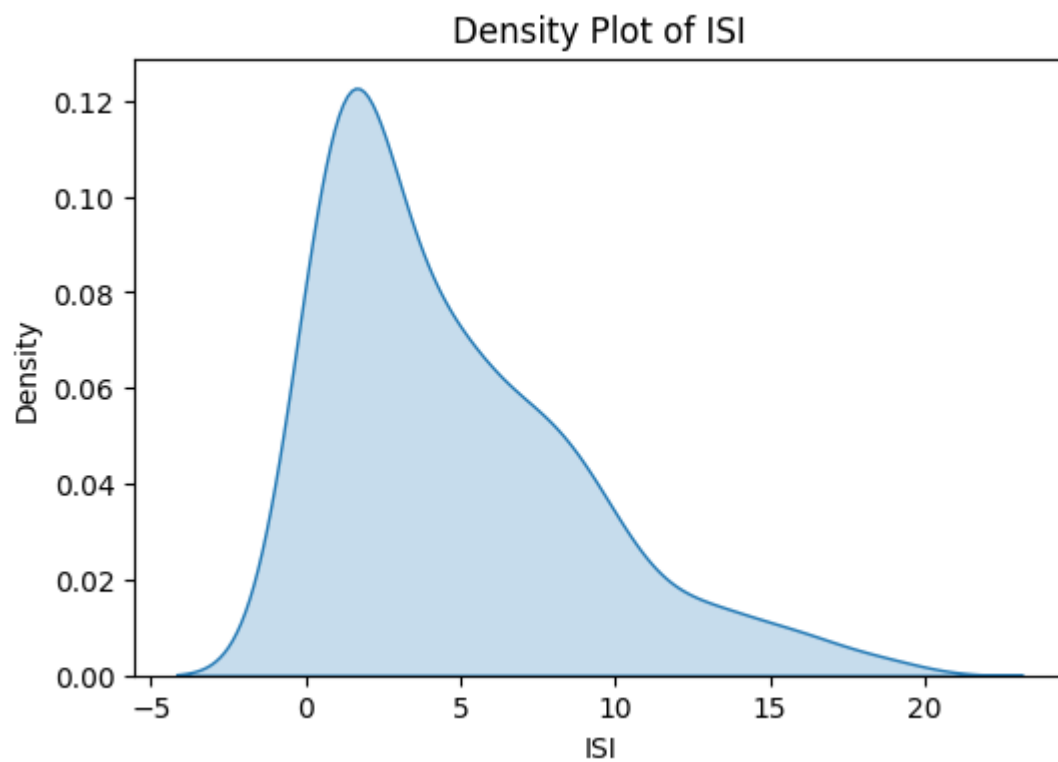


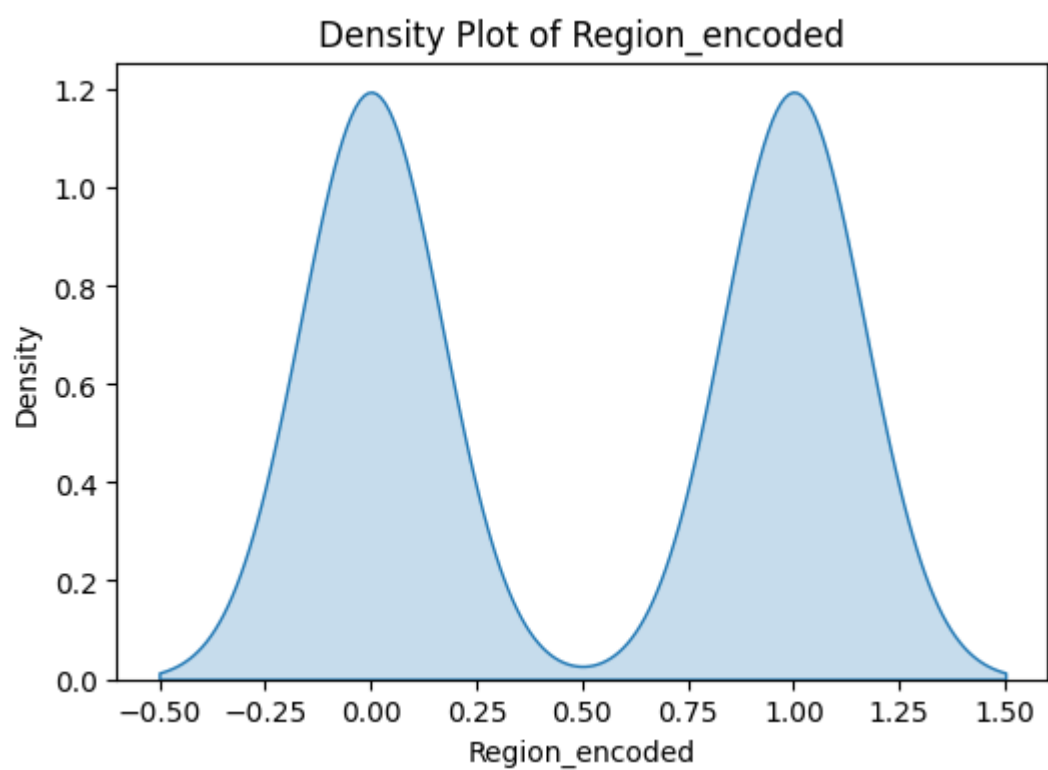
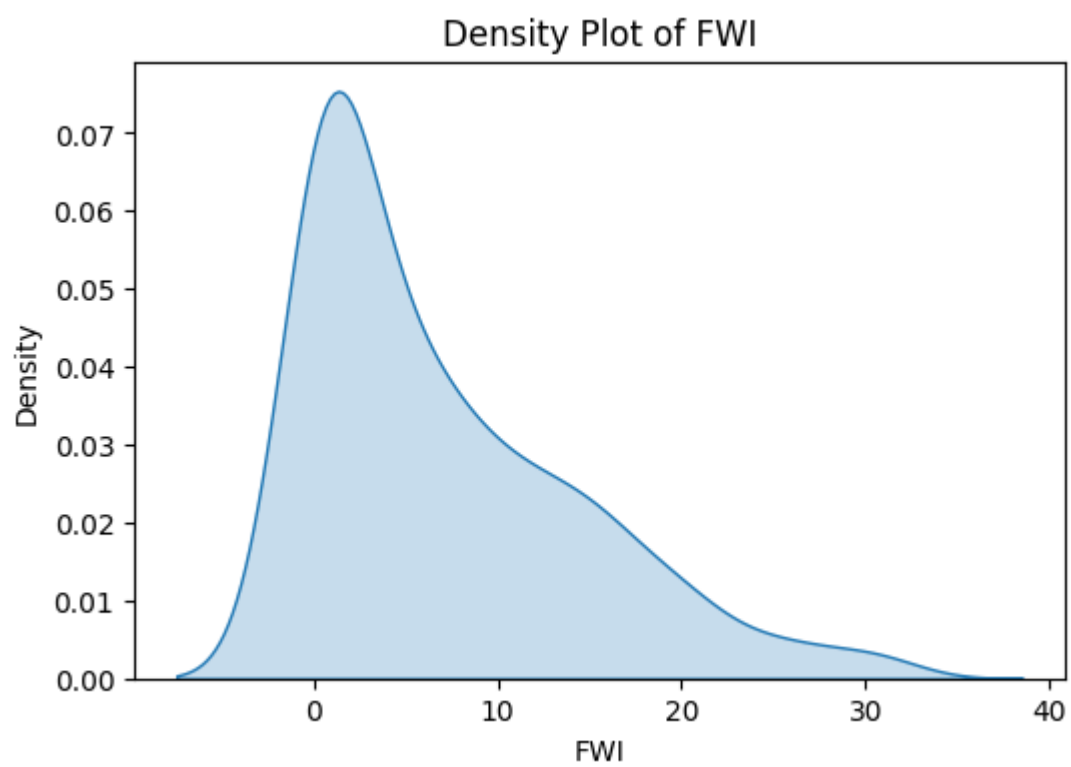


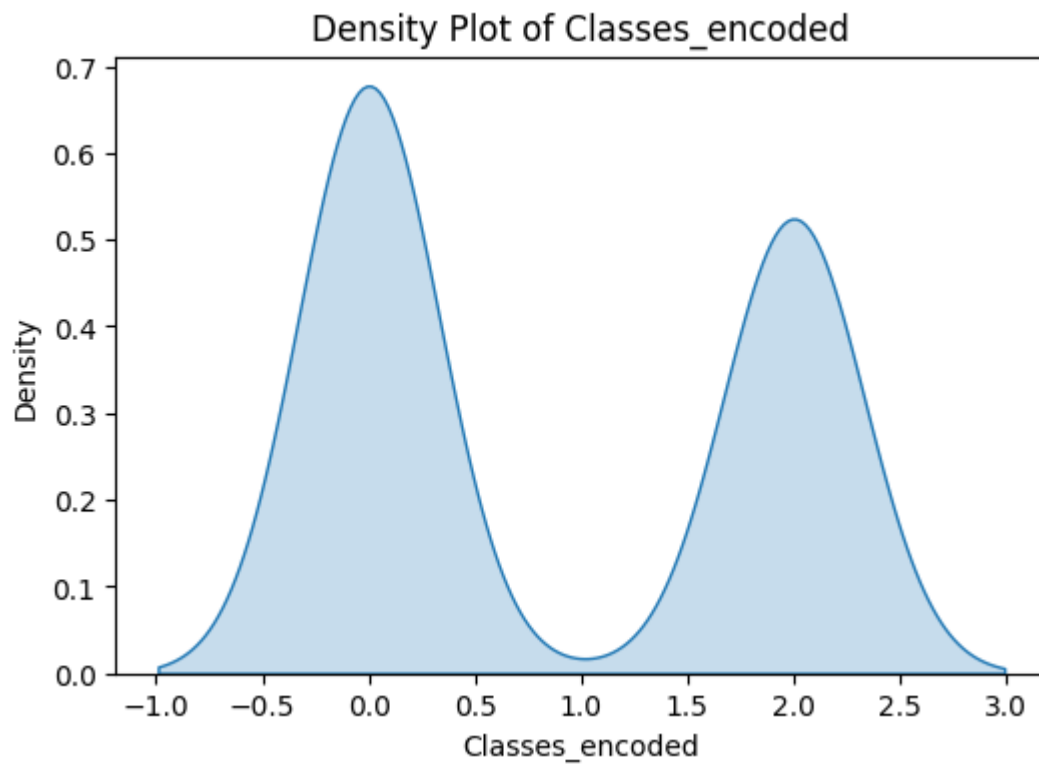












Boxplots for Outlier Detection

`sns.boxplot(x=df[col])`

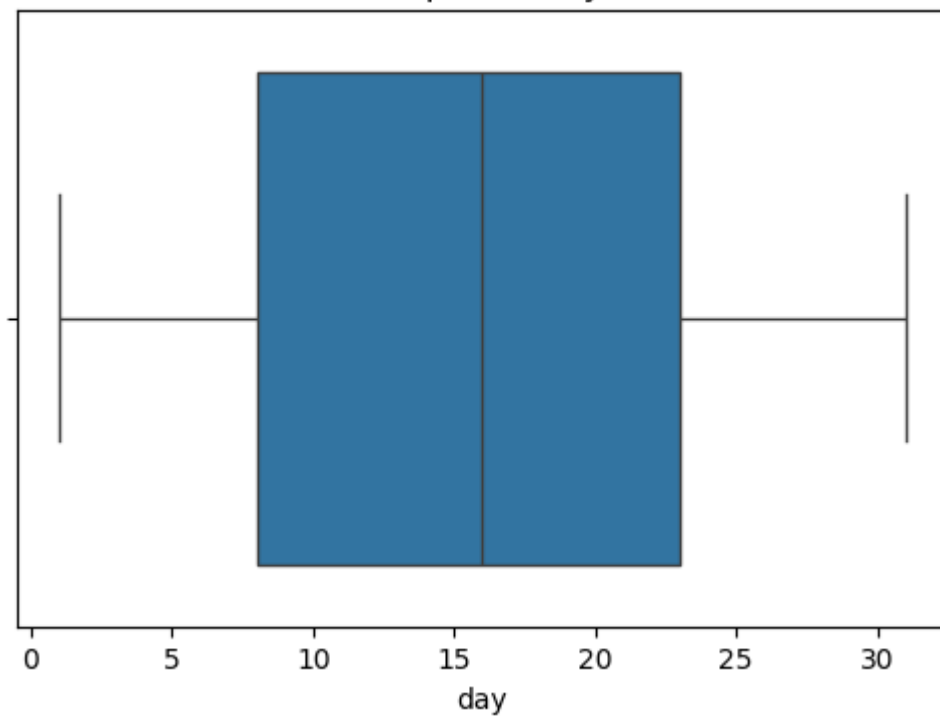
Used to visually identify:

Outliers

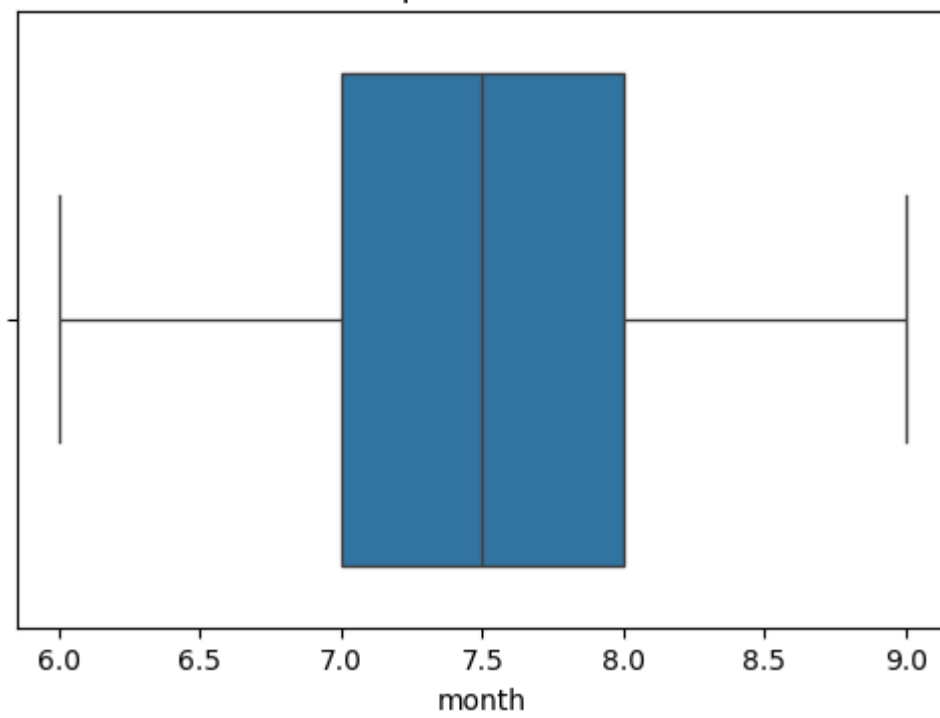
Data spread

Extreme values

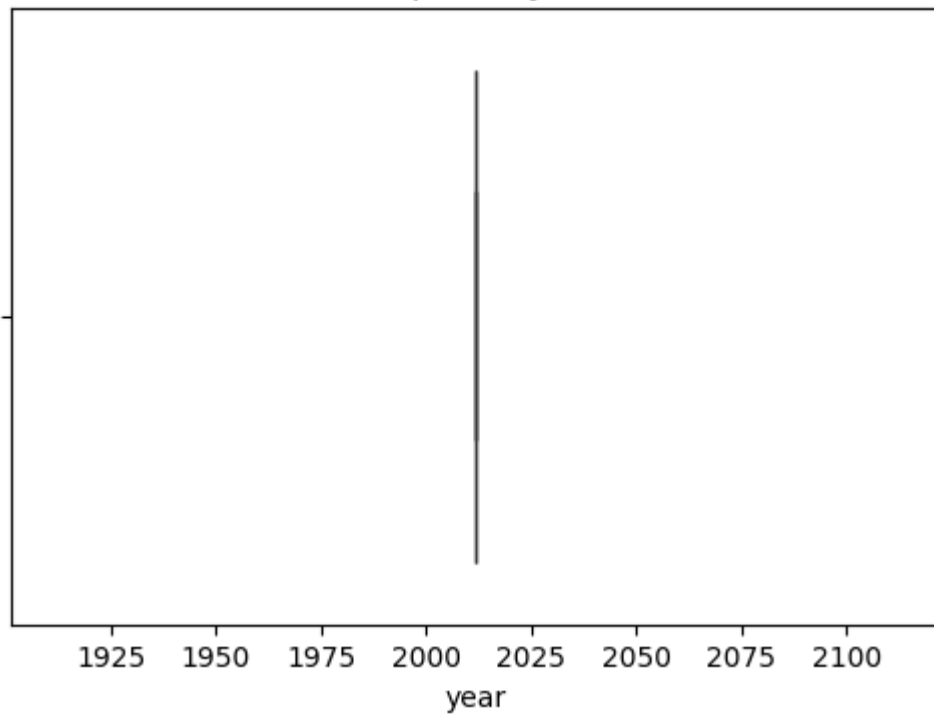
Boxplot of day



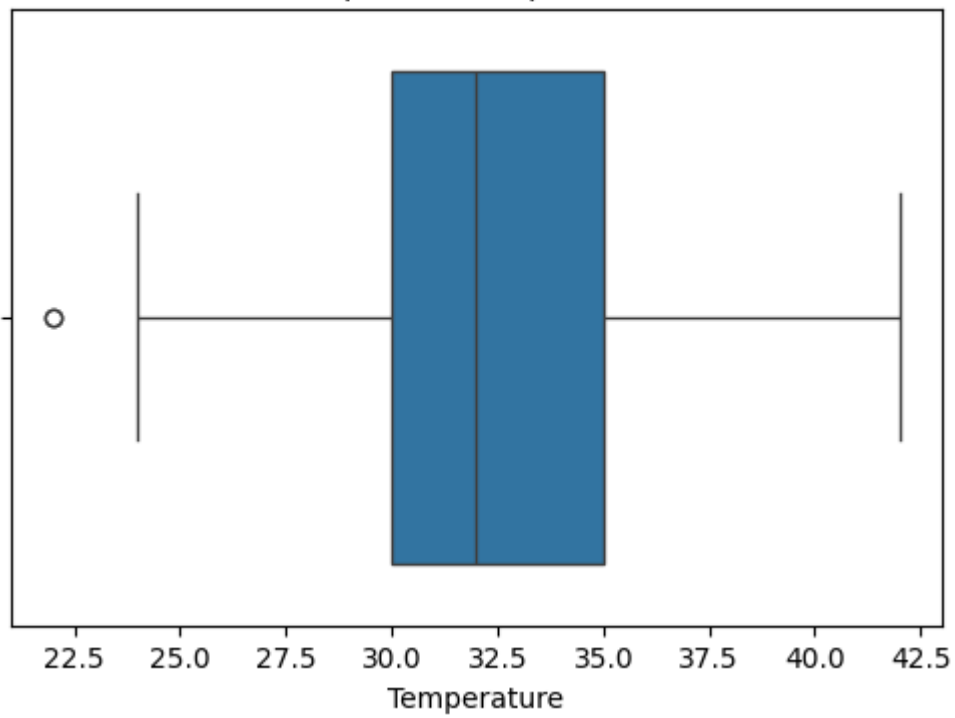
Boxplot of month



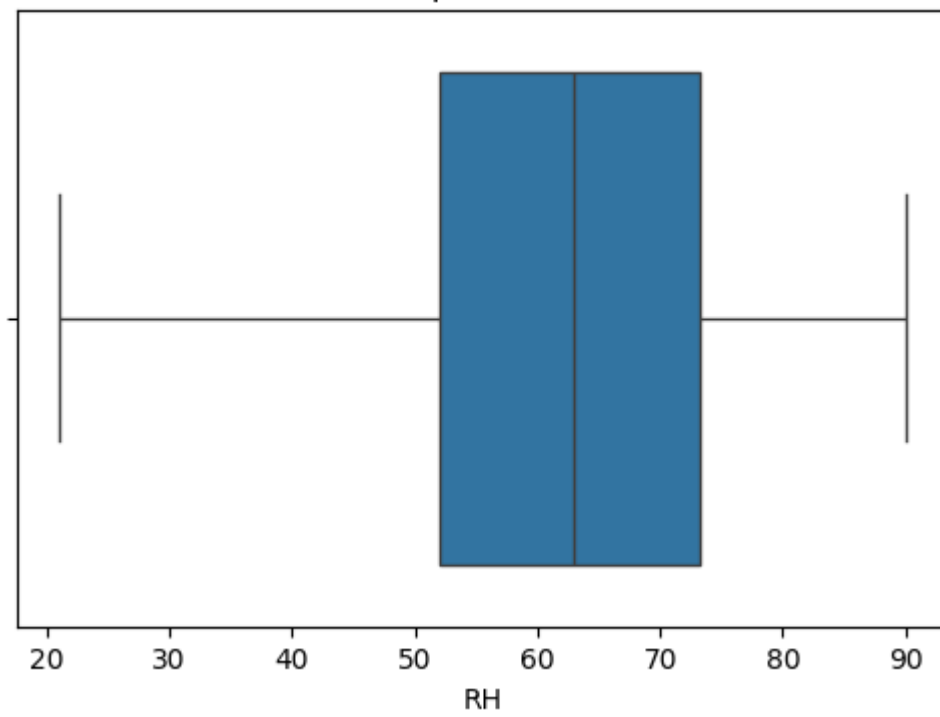
Boxplot of year



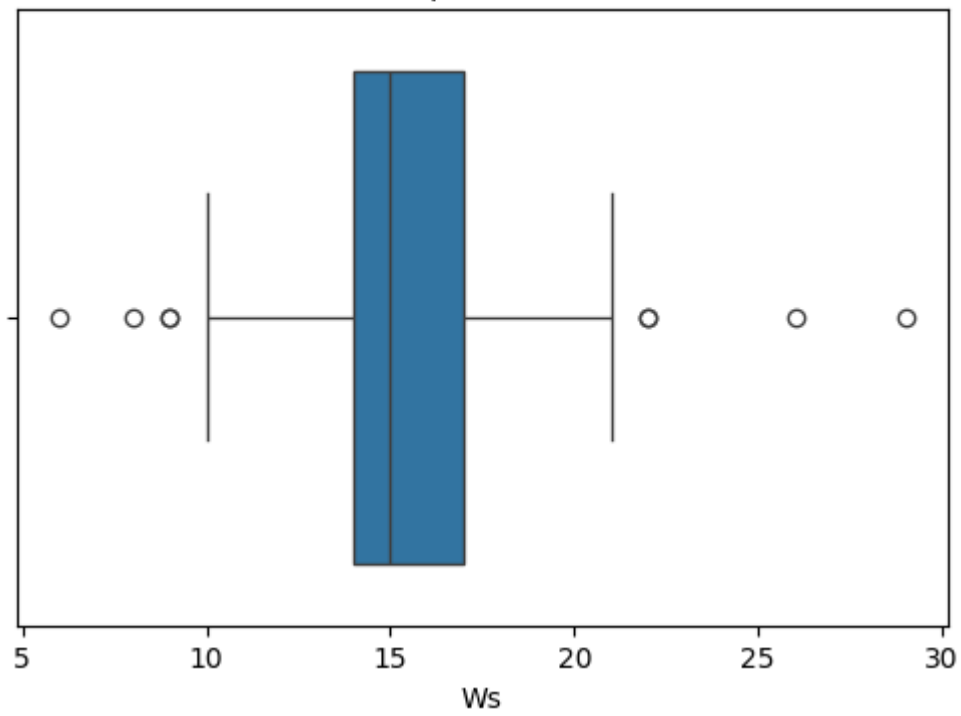
Boxplot of Temperature



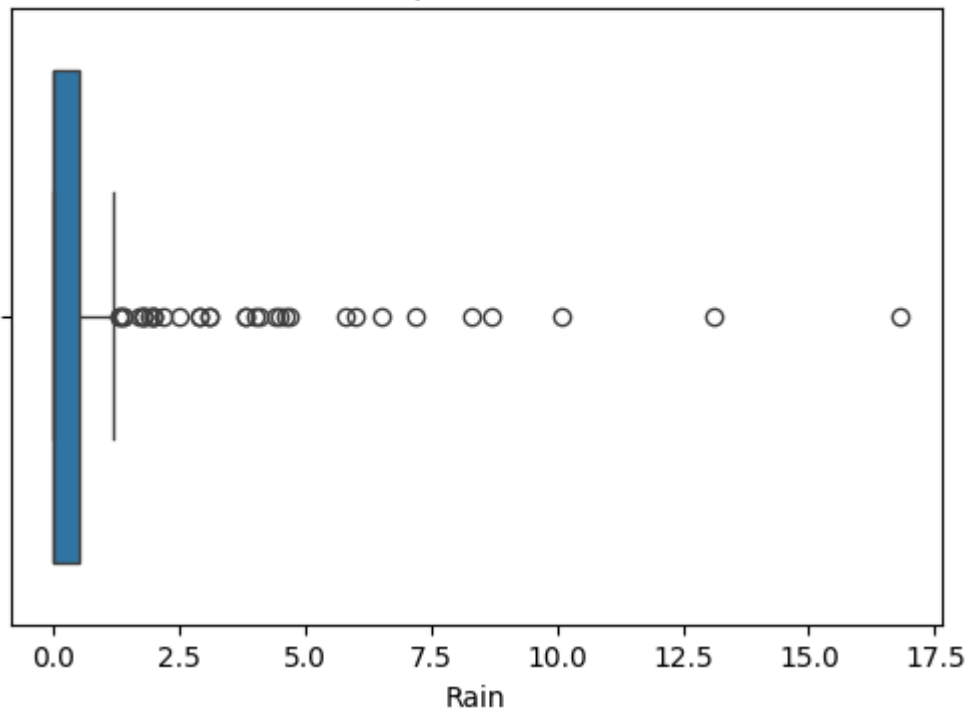
Boxplot of RH



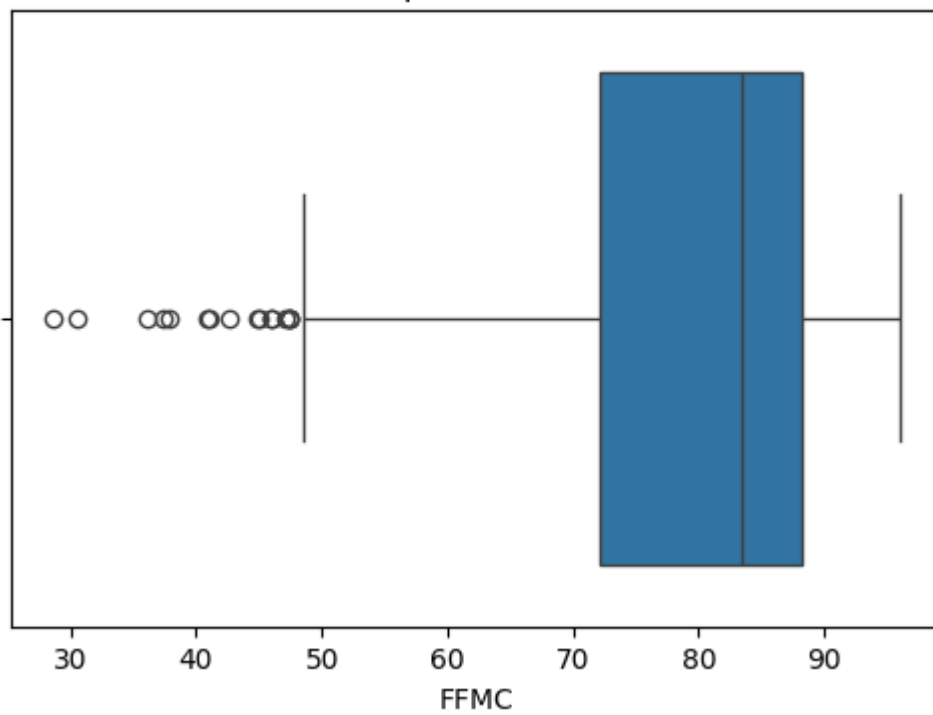
Boxplot of Ws



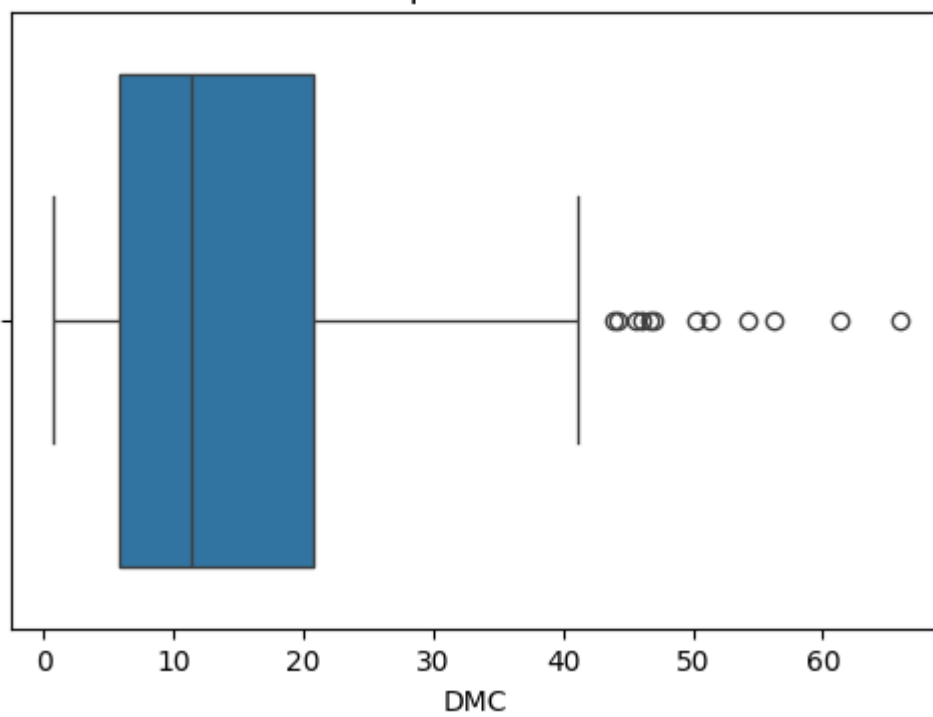
Boxplot of Rain



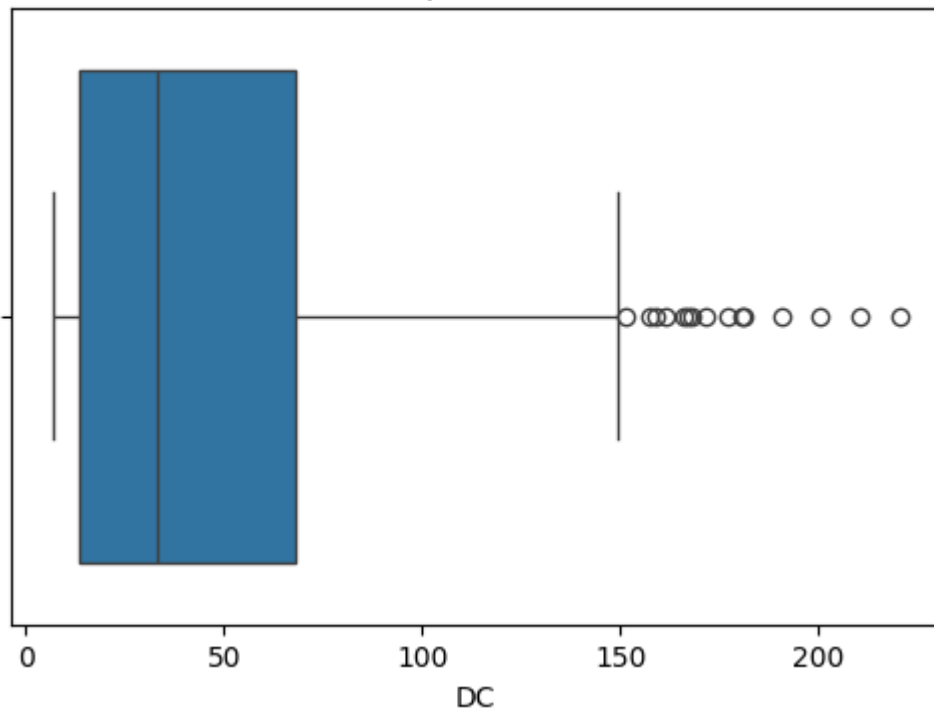
Boxplot of FFMC



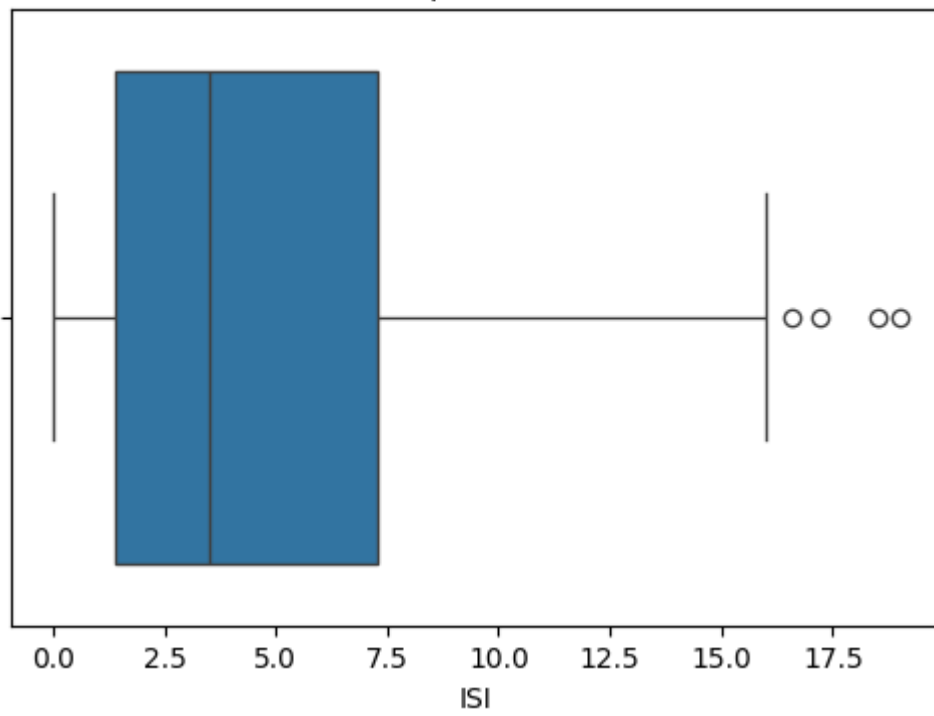
Boxplot of DMC



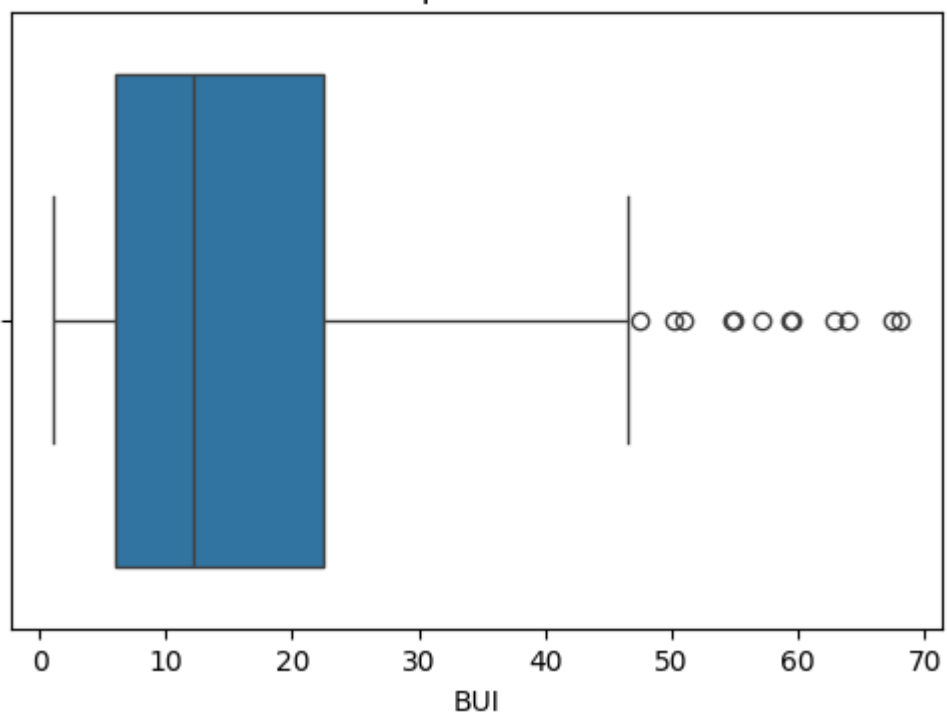
Boxplot of DC



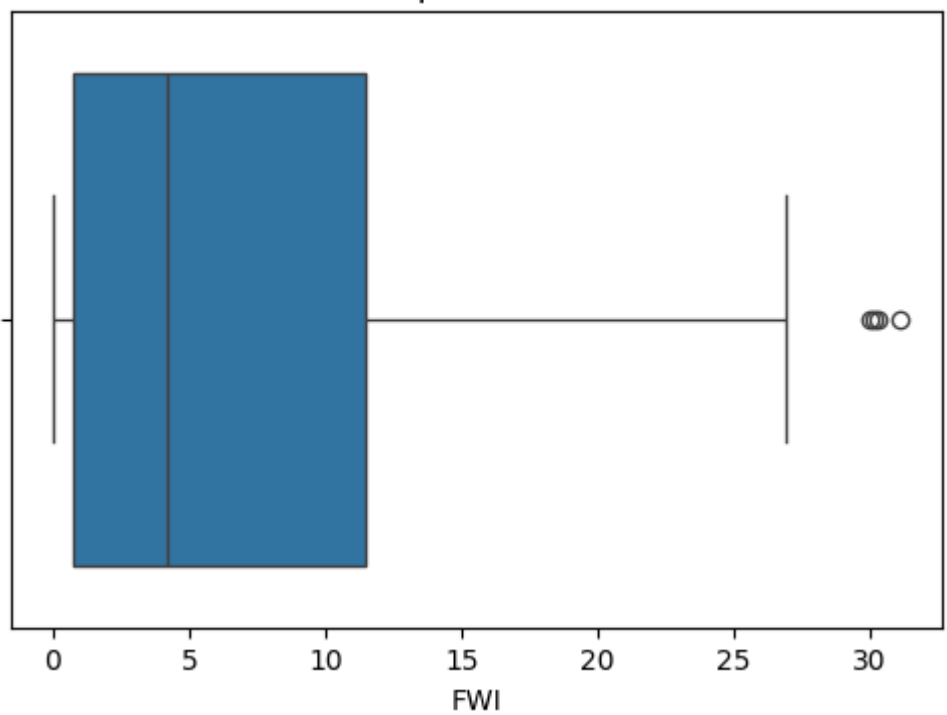
Boxplot of ISI

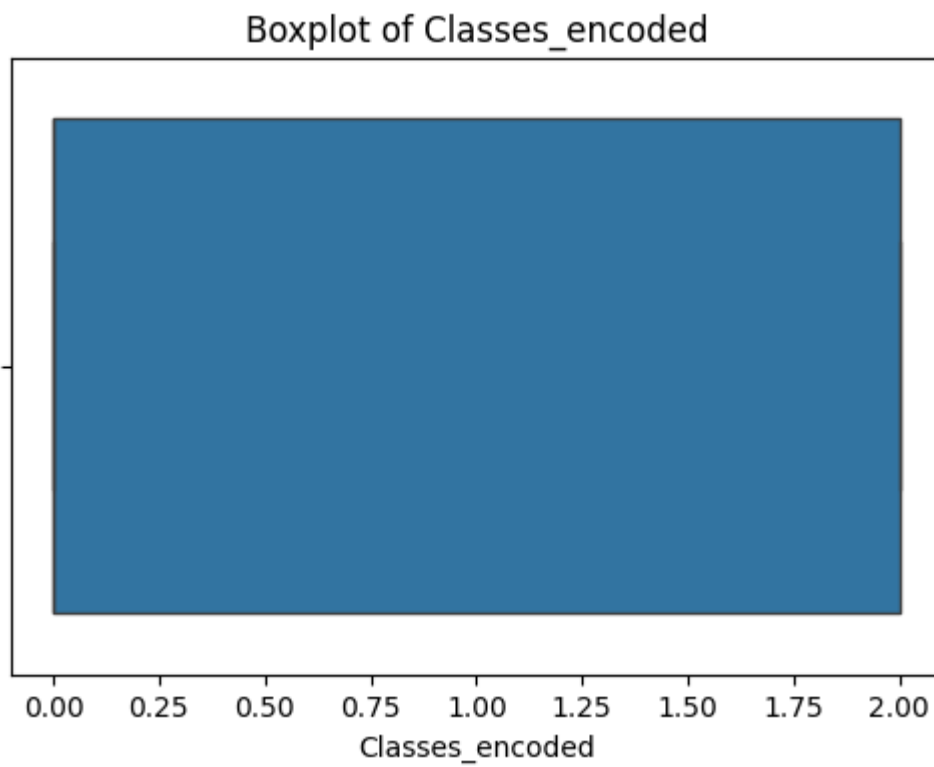
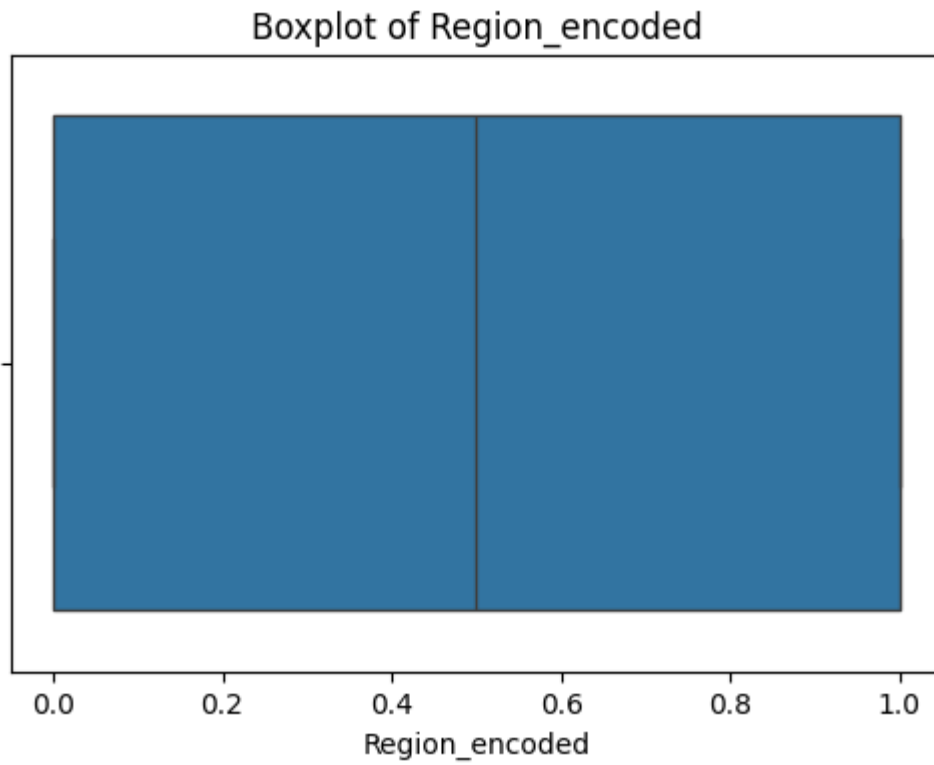


Boxplot of BUI



Boxplot of FWI





Outlier Treatment using IQR

```
Q1 = numeric_df[col].quantile(0.25)
```



```
Q3 = numeric_df[col].quantile(0.75)
```

```
IQR = Q3 - Q1
```

```
df[col] = df[col].clip(lower, upper)
```

Purpose:

Removes extreme values

Prevents model distortion

Makes distributions more stable

Scatter Plots

```
sns.scatterplot(x=df['Temperature'], y=df['FWI'])
```

```
sns.scatterplot(x=df['Ws'], y=df['FWI'])
```

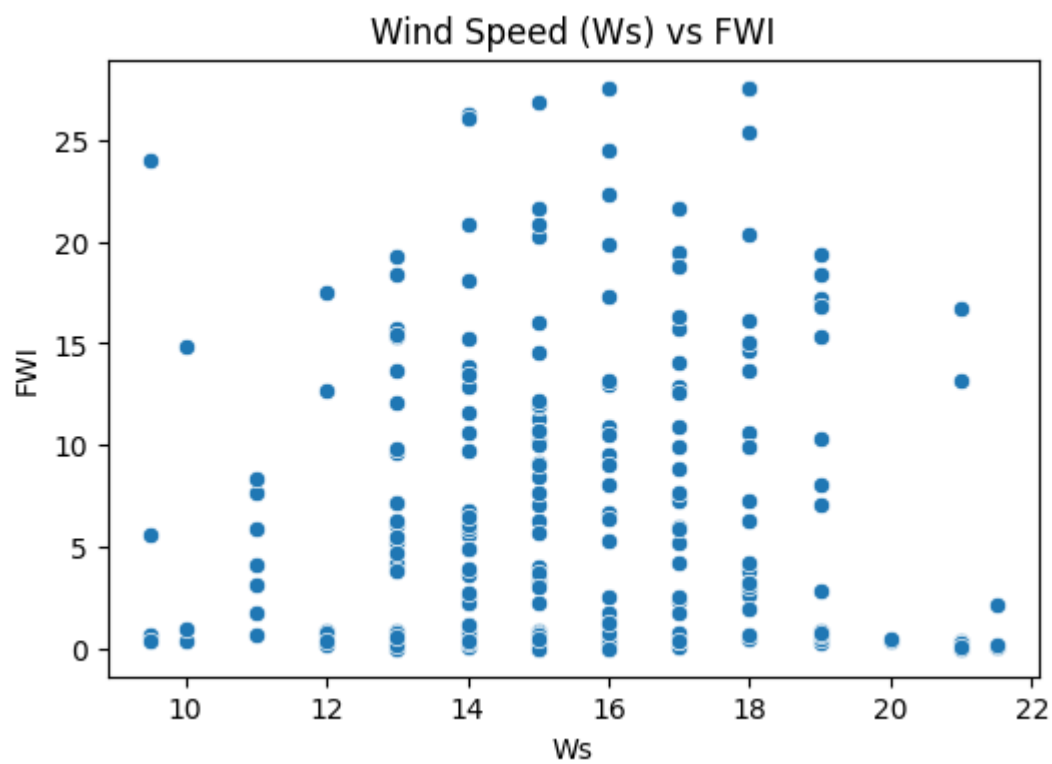
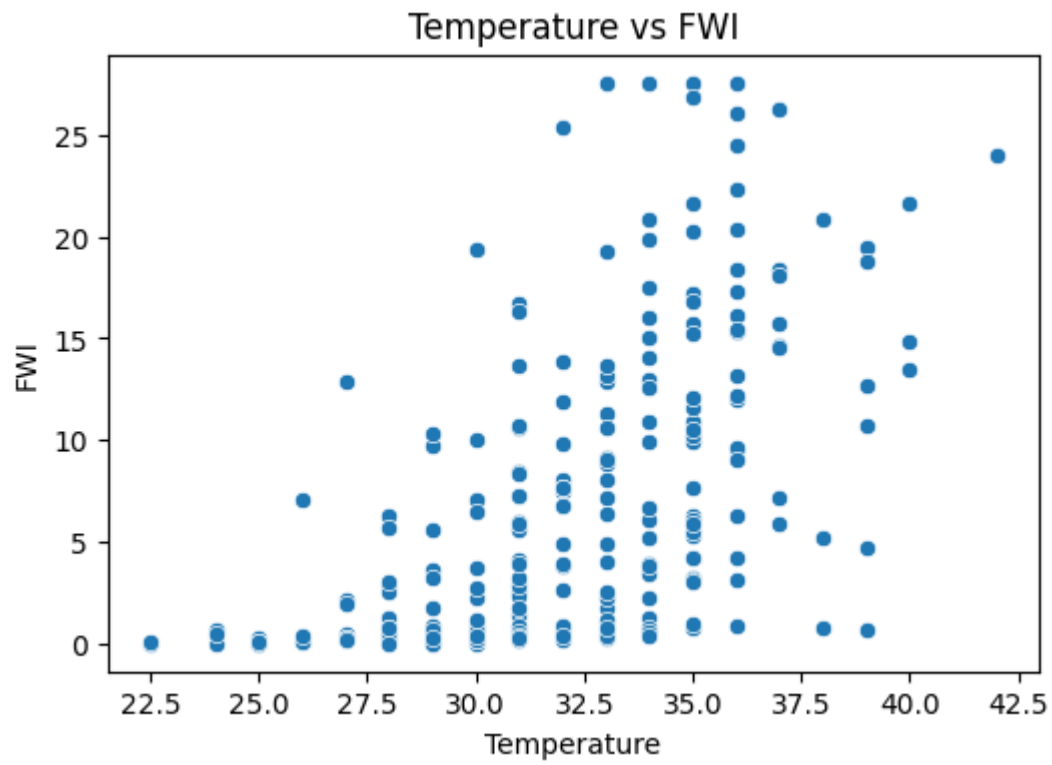
```
sns.scatterplot(x=df['RH'], y=df['FWI'])
```

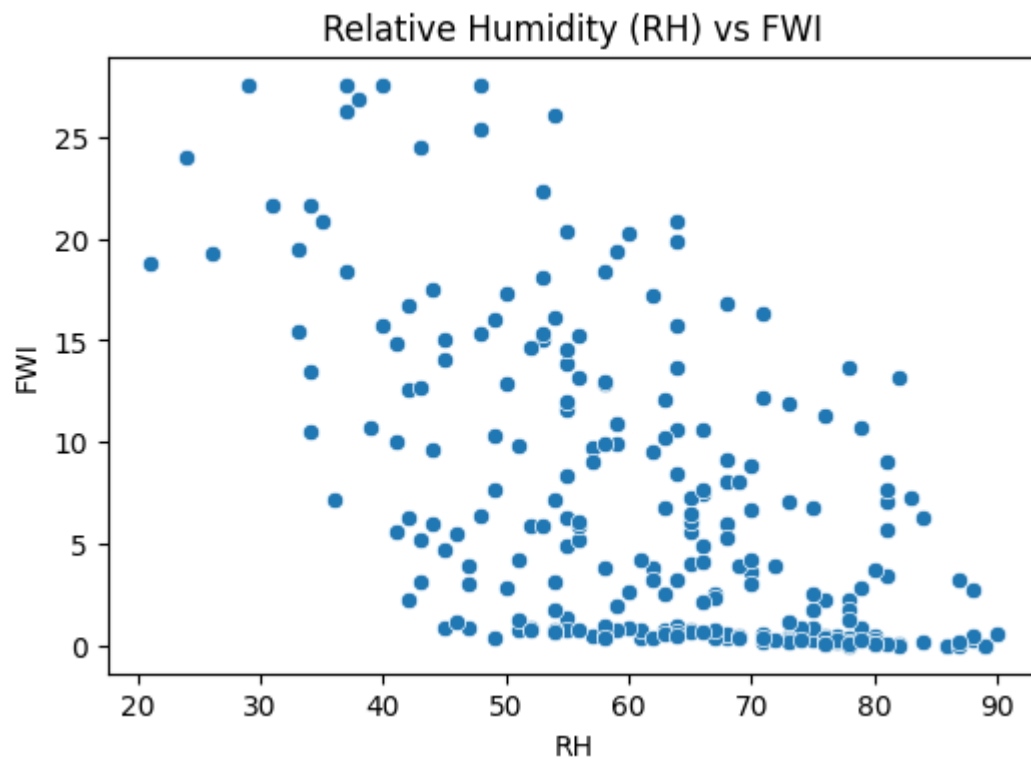
Shows:

How individual features impact FWI

Linear / non-linear trends

Data clusters





To save the cleaned dataset.

```
df.to_csv("FWI Cleaned.csv", index=False)
```

