# Using Machine Learning Tools: Assignment 1

## Overview

In this assignment, you will apply commonly used machine learning techniques to predict the demand for bike rentals. The dataset provided contains information on bike rentals in Seoul, collected over the years 2017 and 2018.

You are assumed to be a new hire at a company that operates a bike rental service among other business activities. Your supervisor has assigned you the task of forecasting bike rentals. While some guidance has been provided (as shown below), you are expected to complete the task with minimal oversight and to report your findings clearly using concise text, visualizations, and well-documented code. The company also expects you to submit all the code necessary to replicate your results. While the organization permits the use of ChatGPT, it's important to demonstrate your own understanding and decision-making skills. Relying solely on ChatGPT may raise concerns about your added value to the company, especially if your work could be done by a less expensive data entry worker.

The primary objectives of this task include:

- Practicing data loading and exploration techniques.

- Identifying and resolving common data quality issues.

- Designing a basic experimental plan and preparing data accordingly.

- Executing your experiment and clearly communicating and interpreting the outcomes.

This assignment aligns with the following ACS CBOK domains: abstraction, design, hardware and software, data and information, human-computer interaction (HCI), and programming.

## General instructions

The assignment consists of multiple tasks. You should complete each task in the space provided within the notebook. Some tasks will involve coding, others will require visualizations, and some will ask for brief written commentary or analysis. It is your responsibility to ensure your responses are clearly presented and that all code cells are executed correctly with visible output before submitting your work.

**Important:** Do not make manual changes to the provided dataset file. For assessment purposes, your code must be compatible with the original file format.

When creating plots or charts, make sure to include clear labels, such as appropriate titles, axis labels, and legends when necessary.

Most tasks only require a few lines of code. A key aim of the assignment is to become familiar with tools such as pandas, matplotlib, sklearn, and similar libraries that will be useful

throughout the course. You are encouraged to consult documentation and other resources to understand how to properly use these tools.

You may refer to Chapter 2 of the course textbook, which follows a similar workflow to this practical exercise, for additional guidance. You are also welcome to use other online resources. While using ChatGPT and similar tools is allowed, you must not copy solutions or code directly related to this specific assignment from others. If you use large code blocks or text from online sources or ChatGPT, include a citation. However, small code snippets from documentation or tutorials do not require references. Reusing and adapting such snippets is normal in real-world problem-solving.

Below is the initial code that imports the key libraries you will need. You typically won't need to change this, but you may import additional libraries if necessary.

```python
# Python ≥3.5 is required
import sys
assert sys.version_info >= (3, 5)

import sklearn
assert sklearn.__version__ >= "0.20"

import pandas as pd
assert pd.__version__ >= "1.0"

# Common imports
import numpy as np
import os

# To plot pretty figures
%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns
mpl.rc('axes', labelsize=14)
mpl.rc('xtick', labelsize=12)
mpl.rc('ytick', labelsize=12)
```

# Step 1: Loading and initial processing of the dataset (40%)

Download the dataset named 'SeoulBikeData.csv' via MyUni using the link available on the assignment page.

This file is in CSV (comma-separated values) format and includes the following attributes:

- Date: formatted as year-month-day

- Rented Bike Count: Total number of bikes rented each hour

- Hour: The hour of the day (0–23)

- Temperature: Measured in Celsius

- Humidity: Given as a percentage

- Windspeed: Measured in meters per second (m/s)

- Visibility: Reported in units of 10 meters

- Dew point temperature: In degrees Celsius

- Solar radiation: Measured in MJ/m²

- Rainfall: Recorded in millimeters

- Snowfall: Recorded in centimeters

- Seasons: Categories include Winter, Spring, Summer, and Autumn

- Holiday: Indicates whether the day is a holiday or not

- Functional Day: Shows if the hour was functional (Fun) or non-functional (NoFunc)

## 1.1 Load and visualise the data

**Read the dataset into a pandas DataFrame, summarise it using one suitable pandas method, and generate one type of plot for each feature (you may choose different plot types for different features as appropriate).**

```python
# 1.1 Load the data
# now i am loading the data into a DataFrame
bike_data = pd.read_csv('SeoulBikeData.csv',
encoding='unicode_escape')

# then i display summary of the data
print("Data Summary:")
print(bike_data.describe())

# Displaying data info
print("\nData Info:")
print(bike_data.info())

# checking for missing values
print("\nMissing Values:")
print(bike_data.isnull().sum())

# Now i plot histograms for numerical features
print("\nVisualizing numerical features:")
bike_data.hist(figsize=(15, 12))
plt.tight_layout()
plt.show()

# Then i plot count plot for categorical features
print("\nVisualizing categorical features:")
```

```python
fig, axes = plt.subplots(1, 3, figsize=(15, 5))
axes[0].set_title('Seasons')
pd.Series(bike_data['Seasons']).value_counts().plot(kind='bar',
ax=axes[0])
axes[0].set_ylabel('Count')

axes[1].set_title('Holiday')
pd.Series(bike_data['Holiday']).value_counts().plot(kind='bar',
ax=axes[1])
axes[1].set_ylabel('Count')

axes[2].set_title('Functioning Day')
pd.Series(bike_data['Functioning
Day']).value_counts().plot(kind='bar', ax=axes[2])
axes[2].set_ylabel('Count')

plt.tight_layout()
plt.show()

# Visualizing bike rental counts across hours
plt.figure(figsize=(10, 6))
bike_data.groupby('Hour')['Rented Bike Count'].mean().plot(kind='bar')
plt.title('Average Bike Rentals by Hour')
plt.xlabel('Hour of the Day')
plt.ylabel('Average Rentals')
plt.tight_layout()
plt.show()

# Then i am checking for correlation between features
plt.figure(figsize=(12, 10))
correlation_matrix =
bike_data.select_dtypes(include=[np.number]).corr()
sns_heatmap = plt.imshow(correlation_matrix, cmap='coolwarm')
plt.colorbar(sns_heatmap)
plt.title('Correlation Matrix')
plt.xticks(range(len(correlation_matrix.columns)),
correlation_matrix.columns, rotation=90)
plt.yticks(range(len(correlation_matrix.columns)),
correlation_matrix.columns)
plt.tight_layout()
plt.show()
```

Data Summary:

|  | Rented Bike Count | Hour | Temperature (C) | Humidity (%) \ |
|---|---|---|---|---|
| count | 8760.000000 | 8760.000000 | 8760.000000 | 8760.000000 |
| mean | 714.876027 | 11.500000 | 12.945765 | 58.268014 |
| std | 1160.468927 | 6.922582 | 12.376168 | 20.807845 |

```
min              0.000000     0.000000     -17.800000    -2.200000

25%            191.000000     5.750000       3.500000    42.000000

50%            504.500000    11.500000      13.700000    57.000000

75%           1066.000000    17.250000      22.500000    74.000000

max          90997.000000    23.000000     195.000000   455.000000


        Wind speed (m/s)  Visibility (10m)  Dew point temperature (C)
count         8760.000000       8760.000000                8760.000000
mean             1.848950       1436.825799                   4.073813
std             10.665215        608.298712                  13.060369
min             -0.700000         27.000000                 -30.600000
25%              0.900000        940.000000                  -4.700000
50%              1.500000       1698.000000                   5.100000
75%              2.300000       2000.000000                  14.800000
max            991.100000       2000.000000                  27.200000

Data Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8760 entries, 0 to 8759
Data columns (total 14 columns):
 #   Column                     Non-Null Count  Dtype
---  ------                     --------------  -----
 0   Date                       8760 non-null   object
 1   Rented Bike Count          8760 non-null   int64
 2   Hour                       8760 non-null   int64
 3   Temperature (C)            8760 non-null   float64
 4   Humidity (%)               8760 non-null   float64
 5   Wind speed (m/s)           8760 non-null   float64
 6   Visibility (10m)           8760 non-null   int64
 7   Dew point temperature (C)  8760 non-null   float64
 8   Solar Radiation (MJ/m2)    8760 non-null   object
 9   Rainfall(mm)               8760 non-null   object
 10  Snowfall (cm)              8760 non-null   object
 11  Seasons                    8760 non-null   object
 12  Holiday                    8760 non-null   object
 13  Functioning Day            8760 non-null   object
dtypes: float64(4), int64(3), object(7)
memory usage: 958.3+ KB
None

Missing Values:
Date                         0
Rented Bike Count            0
Hour                         0
```
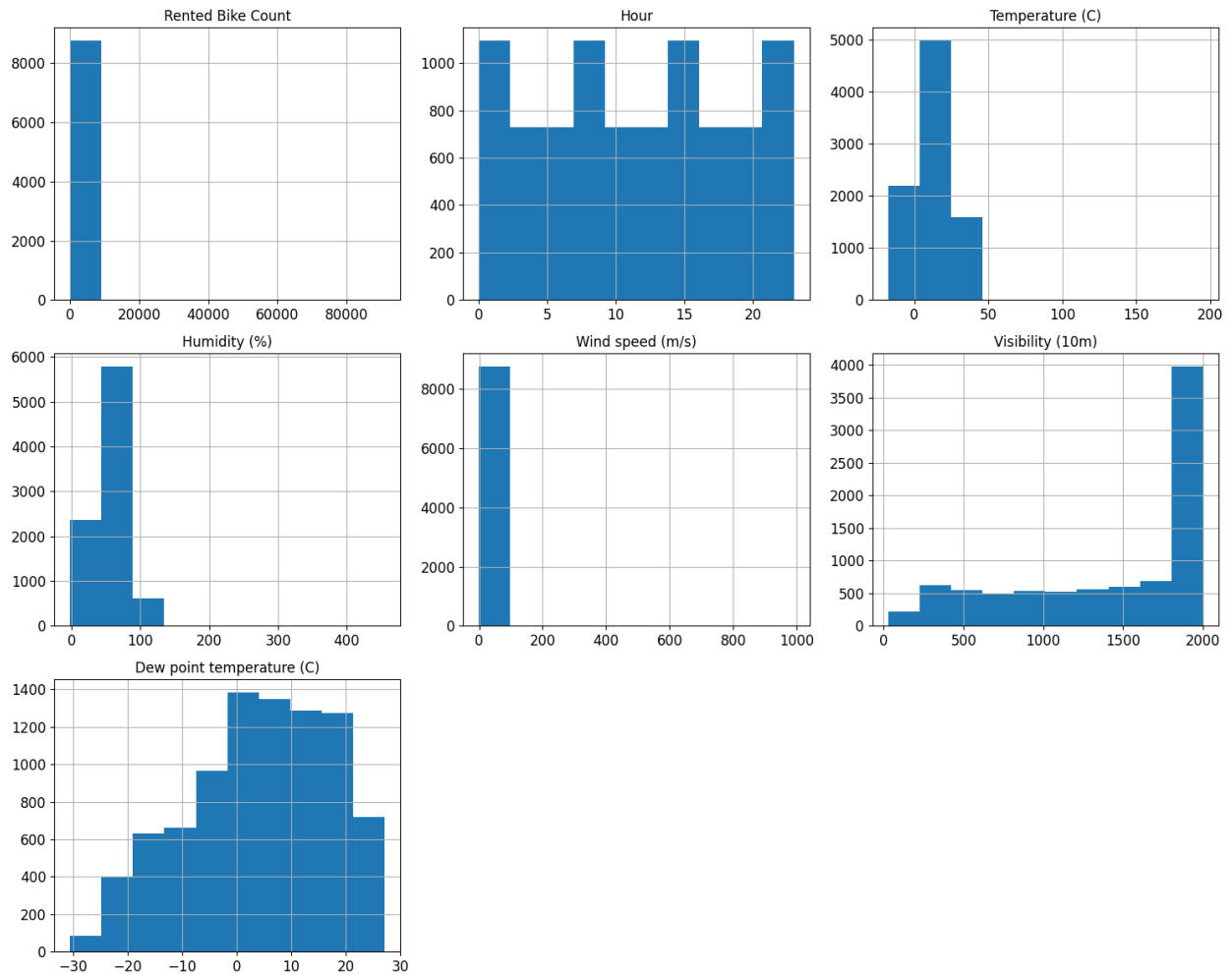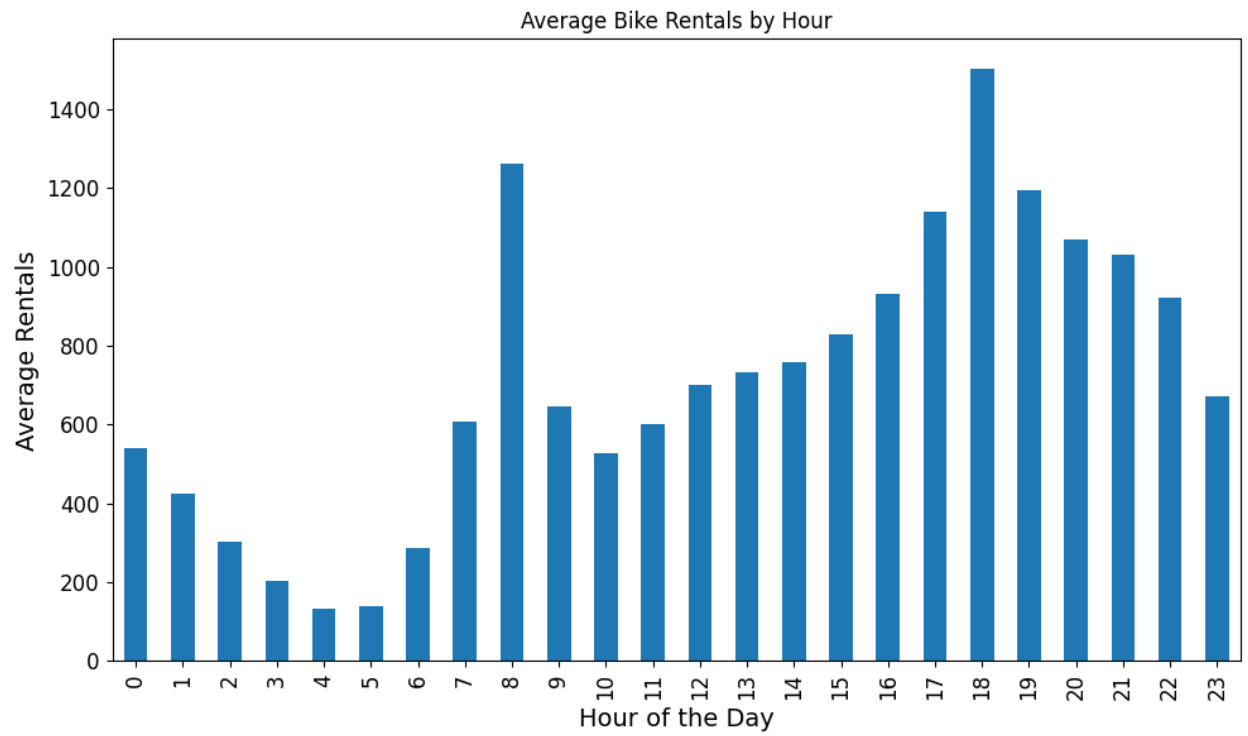
```
Temperature (C)                0
Humidity (%)                   0
Wind speed (m/s)               0
Visibility (10m)               0
Dew point temperature (C)      0
Solar Radiation (MJ/m2)        0
Rainfall(mm)                   0
Snowfall (cm)                  0
Seasons                        0
Holiday                        0
Functioning Day                0
dtype: int64

Visualizing numerical features:
```
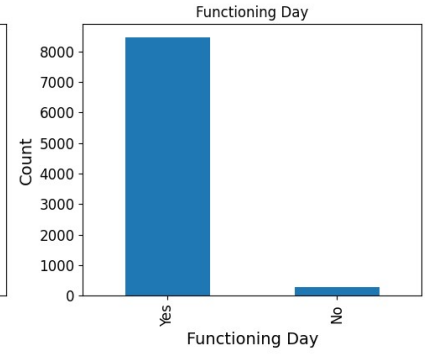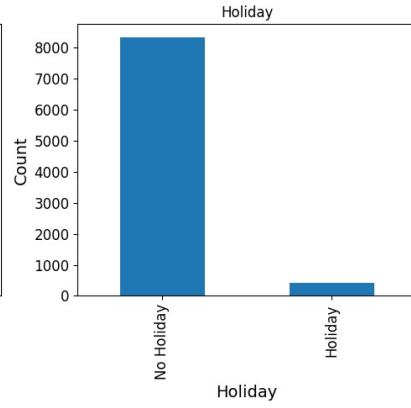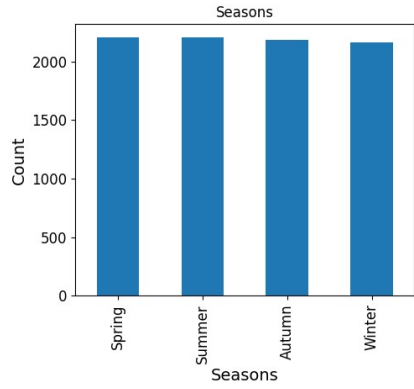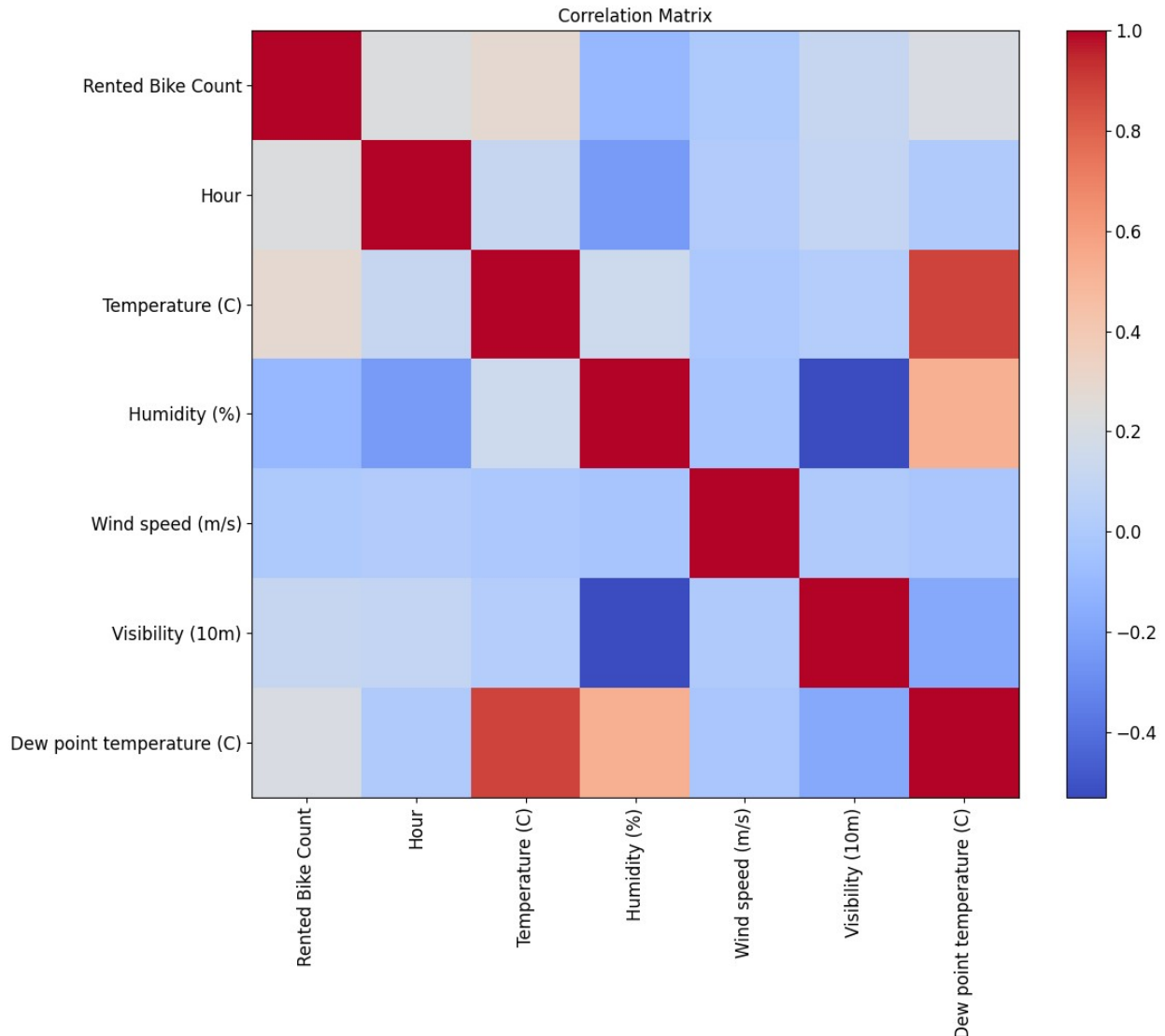


```
Visualizing categorical features:
```

Correlation Matrix

## 1.2 Cleaning the data

Do the following to the data:

- Using the "Functioning day" feature, **remove rows from the DataFrame** where the business is closed and then **delete the Functioning Day feature from the DataFrame**.
- **Convert seasons to a one hot encoded format** (1 binary feature for each of the 4 seasons).
- Replace the **Date** feature with a binary **Weekday** feature (1 for a weekday and 0 for weekend) using the code sample below or your own code.
- **Convert remaining non-numerical features to a numerical format** or replace with NaN (i.e. `np.nan`) where not possible.
- **Identify and fix any outliers and errors in the data**.

Save the result as a new csv file called `CleanedSeoulBikeData.csv` and **upload this** to MyUni along with this notebook when you submit your assignment.

```python
## Example code for weekday feature mapping ##

import datetime
def date_is_weekday(datestring):
    ### return 0 if weekend, 1 if weekday
    dsplit = datestring.split('/')
    wday =
datetime.datetime(int(dsplit[2]),int(dsplit[1]),int(dsplit[0])).weekda
y()
    return int(wday<=4)

# 1.2 Cleaning the data

# Now i am makeing a copy of the original data
cleaned_bike_data = bike_data.copy()

# 1. Then i am removing rows where business is closed and delete the
'Functioning Day' column
print("Original shape:", cleaned_bike_data.shape)
cleaned_bike_data = cleaned_bike_data[cleaned_bike_data['Functioning
Day'] == 'Yes']
cleaned_bike_data = cleaned_bike_data.drop(['Functioning Day'],
axis=1)
print("Shape after removing non-functioning days:",
cleaned_bike_data.shape)

# 2. After that i convert seasons to one-hot encoded format
seasons_dummies = pd.get_dummies(cleaned_bike_data['Seasons'],
prefix='Season')
cleaned_bike_data = pd.concat([cleaned_bike_data, seasons_dummies],
axis=1)
cleaned_bike_data = cleaned_bike_data.drop(['Seasons'], axis=1)

# 3. converting Date to binary Weekday feature (1 for weekday, 0 for
weekend)
cleaned_bike_data['Weekday'] =
cleaned_bike_data['Date'].apply(date_is_weekday)
cleaned_bike_data = cleaned_bike_data.drop(['Date'], axis=1)

# 4. Then i convert remaining non-numerical features to numerical
format

cleaned_bike_data['Holiday'] =
cleaned_bike_data['Holiday'].map({'Holiday': 1, 'No Holiday': 0})

# 5. checking for outliers and errors in the data
# Then i Display statistical summary
print("\nStatistical Summary After Initial Cleaning:")
print(cleaned_bike_data.describe())
```

```python
# Checking for remaining missing values
print("\nMissing Values:")
print(cleaned_bike_data.isnull().sum())

# Visualizing outliers with boxplots for numerical features
plt.figure(figsize=(15, 10))
cleaned_bike_data.select_dtypes(include=[np.number]).boxplot(figsize=(
15, 10))
plt.title('Boxplots to Identify Outliers')
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()

# I get all numerical columns first for outlier detection
numerical_cols =
cleaned_bike_data.select_dtypes(include=[np.number]).columns

# Fixing extreme outliers for physical variables that have physical
constraints
temp_min, temp_max = -20, 40
outliers = cleaned_bike_data[(cleaned_bike_data['Temperature (C)'] <
temp_min) |
                            (cleaned_bike_data['Temperature (C)'] >
temp_max)].shape[0]
print(f"\nFound {outliers} temperature outliers outside range
[{temp_min}, {temp_max}]°C")
cleaned_bike_data['Temperature (C)'] = cleaned_bike_data['Temperature
(C)'].clip(lower=temp_min, upper=temp_max)

# Humidity: must be between 0% and 100%
humid_min, humid_max = 0, 100
outliers = cleaned_bike_data[(cleaned_bike_data['Humidity (%)'] <
humid_min) |
                            (cleaned_bike_data['Humidity (%)'] >
humid_max)].shape[0]
print(f"Found {outliers} humidity outliers outside range [{humid_min},
{humid_max}]%")
cleaned_bike_data['Humidity (%)'] = cleaned_bike_data['Humidity
(%)'].clip(lower=humid_min, upper=humid_max)

# Wind speed: reasonable max for Seoul is around 30 m/s
wind_min, wind_max = 0, 30
outliers = cleaned_bike_data[(cleaned_bike_data['Wind speed (m/s)'] <
wind_min) |
                            (cleaned_bike_data['Wind speed (m/s)'] >
wind_max)].shape[0]
print(f"Found {outliers} wind speed outliers outside range
[{wind_min}, {wind_max}] m/s")
cleaned_bike_data['Wind speed (m/s)'] = cleaned_bike_data['Wind speed
(m/s)'].clip(lower=wind_min, upper=wind_max)
```

```python
# Rented Bike Count: shouldn't be unreasonably high
bike_max = 8000  # Based on examining the data distribution
outliers = cleaned_bike_data[cleaned_bike_data['Rented Bike Count'] >
bike_max].shape[0]
print(f"Found {outliers} bike count outliers above {bike_max}")
cleaned_bike_data['Rented Bike Count'] = cleaned_bike_data['Rented
Bike Count'].clip(upper=bike_max)

# For other numerical columns, applhying standard statistical outlier
detection (3 std from mean)
for col in numerical_cols:
    if col not in ['Temperature (C)', 'Humidity (%)', 'Wind speed
(m/s)', 'Rented Bike Count']:
        mean = cleaned_bike_data[col].mean()
        std = cleaned_bike_data[col].std()
        lower_bound = mean - 3 * std
        upper_bound = mean + 3 * std
        outliers = cleaned_bike_data[(cleaned_bike_data[col] <
lower_bound) |
                                     (cleaned_bike_data[col] >
upper_bound)].shape[0]
        if outliers > 0:
            print(f"Found {outliers} statistical outliers in {col}")
            cleaned_bike_data[col] =
cleaned_bike_data[col].clip(lower=lower_bound, upper=upper_bound)
for col in numerical_cols:
    if col not in ['Temperature (C)', 'Dew point temperature (C)']:  #
These can be negative
        neg_counts = (cleaned_bike_data[col] < 0).sum()
        if neg_counts > 0:
            print(f"Found {neg_counts} negative values in {col}.
Setting them to 0.")
            cleaned_bike_data[col] =
cleaned_bike_data[col].clip(lower=0)

# Saving the cleaned data to CSV
cleaned_bike_data.to_csv('CleanedSeoulBikeData.csv', index=False)
print("\nCleaned data saved to 'CleanedSeoulBikeData.csv'")

# Displaying the first few rows of the cleaned data
print("\nFirst few rows of cleaned data:")
cleaned_bike_data.head()
```

```
Original shape: (8760, 14)
Shape after removing non-functioning days: (8465, 13)

Statistical Summary After Initial Cleaning:
       Rented Bike Count        Hour  Temperature (C)  Humidity
(%)  \
```

|       |               |             |             |             |
|-------|---------------|-------------|-------------|-------------|
| count | 8465.000000   | 8465.000000 | 8465.000000 | 8465.000000 |
| mean  | 739.789014    | 11.507029   | 12.836090   | 58.190408   |
| std   | 1172.685973   | 6.920899    | 12.545382   | 20.943092   |
| min   | 2.000000      | 0.000000    | -17.800000  | -2.200000   |
| 25%   | 214.000000    | 6.000000    | 3.000000    | 42.000000   |
| 50%   | 542.000000    | 12.000000   | 13.500000   | 57.000000   |
| 75%   | 1084.000000   | 18.000000   | 22.700000   | 74.000000   |
| max   | 90997.000000  | 23.000000   | 195.000000  | 455.000000  |

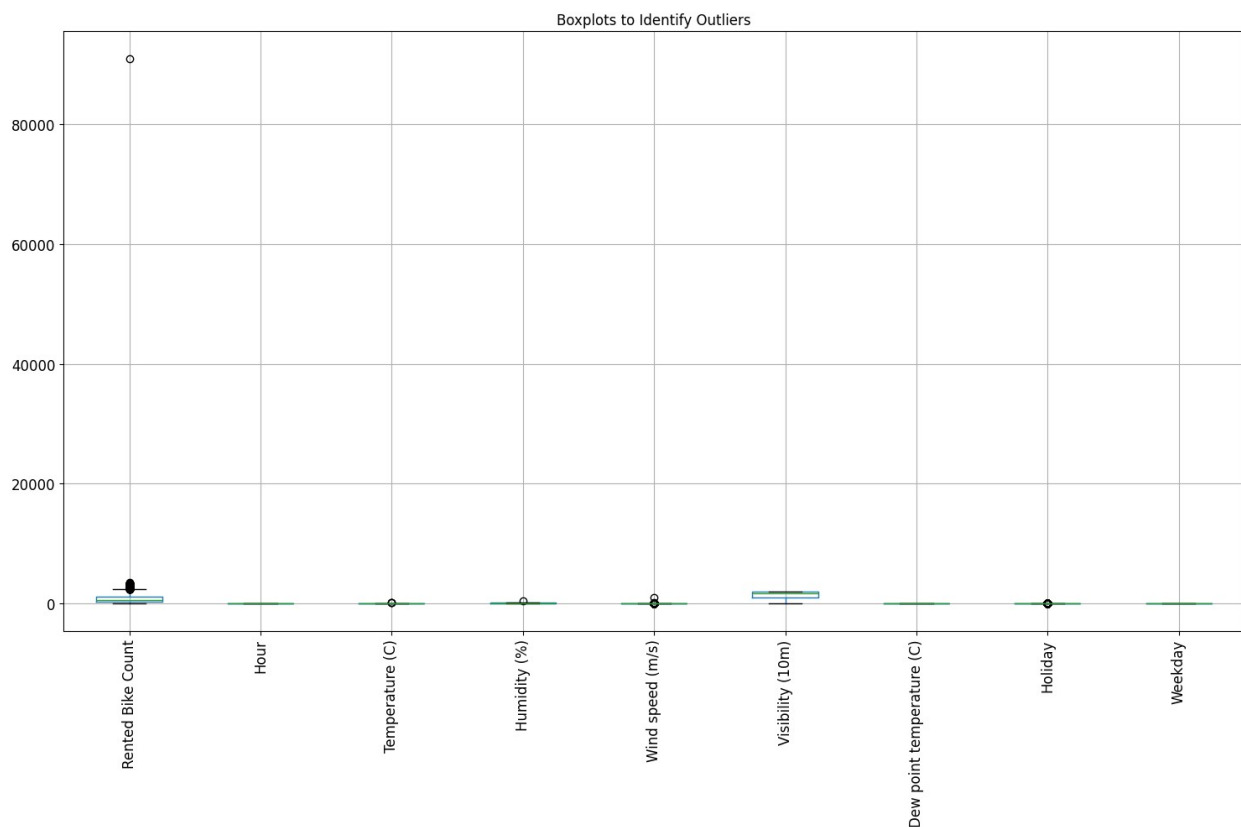|       | Wind speed (m/s) | Visibility (10m) | Dew point temperature (C) \ |
|-------|------------------|------------------|------------------------------|
| count | 8465.000000      | 8465.000000      | 8465.000000                  |
| mean  | 1.854247         | 1433.873479      | 3.944997                     |
| std   | 10.847528        | 609.051229       | 13.242399                    |
| min   | -0.700000        | 27.000000        | -30.600000                   |
| 25%   | 0.900000         | 935.000000       | -5.100000                    |
| 50%   | 1.500000         | 1690.000000      | 4.700000                     |
| 75%   | 2.300000         | 2000.000000      | 15.200000                    |
| max   | 991.100000       | 2000.000000      | 27.200000                    |

|       | Holiday     | Weekday     |
|-------|-------------|-------------|
| count | 8465.000000 | 8465.000000 |
| mean  | 0.048198    | 0.711636    |
| std   | 0.214198    | 0.453028    |
| min   | 0.000000    | 0.000000    |
| 25%   | 0.000000    | 0.000000    |
| 50%   | 0.000000    | 1.000000    |
| 75%   | 0.000000    | 1.000000    |
| max   | 1.000000    | 1.000000    |

Missing Values:
Rented Bike Count          0
Hour                       0
Temperature (C)            0
Humidity (%)               0

```
Wind speed (m/s)            0
Visibility (10m)            0
Dew point temperature (C)   0
Solar Radiation (MJ/m2)     0
Rainfall(mm)                0
Snowfall (cm)               0
Holiday                     0
Season_Autumn               0
Season_Spring               0
Season_Summer               0
Season_Winter               0
Weekday                     0
dtype: int64
```



Boxplots to Identify Outliers

```
Found 4 temperature outliers outside range [-20, 40]°C
Found 2 humidity outliers outside range [0, 100]%
Found 3 wind speed outliers outside range [0, 30] m/s
Found 1 bike count outliers above 8000
Found 408 statistical outliers in Holiday

Cleaned data saved to 'CleanedSeoulBikeData.csv'

First few rows of cleaned data:
```

| | Rented Bike Count | Hour | Temperature (C) | Humidity (%) | Wind speed (m/s) |
|---|---|---|---|---|---|
| 0 | 254 | 0 | -5.2 | 37.0 | 2.2 |
| 1 | 204 | 1 | -5.5 | 38.0 | 0.8 |
| 2 | 173 | 2 | -6.0 | 39.0 | 1.0 |
| 3 | 107 | 3 | -6.2 | 40.0 | 0.9 |
| 4 | 78 | 4 | -6.0 | 36.0 | 2.3 |

| | Visibility (10m) | Dew point temperature (C) | Solar Radiation (MJ/m2) |
|---|---|---|---|
| 0 | 2000 | -17.6 | 0 |
| 1 | 2000 | -17.6 | 0 |
| 2 | 2000 | -17.7 | 0 |
| 3 | 2000 | -17.6 | 0 |
| 4 | 2000 | -18.6 | 0 |

| | Rainfall(mm) | Snowfall (cm) | Holiday | Season_Autumn | Season_Spring |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0.0 | False | False |
| 1 | 0 | 0 | 0.0 | False | False |
| 2 | 0 | 0 | 0.0 | False | False |
| 3 | 0 | 0 | 0.0 | False | False |
| 4 | 0 | 0 | 0.0 | False | False |

| | Season_Summer | Season_Winter | Weekday |
|---|---|---|---|
| 0 | False | True | 1 |
| 1 | False | True | 1 |
| 2 | False | True | 1 |
| 3 | False | True | 1 |
| 4 | False | True | 1 |

# Step 2: Pre-process the data and perform the first fit (20%)

## 2.1 Imputation and Pre-Processing

Make sure that you have set any problematic values in the numerical data to `np.nan` and then write code for a **sklearn *pipeline* that will perform imputation** to replace problematic entries (nan values) with an appropriate **median** value *and* **do any other pre-processing** that you think should be used.

```python
# 2.1 Imputation and Pre-Processing

# Importing necessary preprocessing tools from sklearn
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.compose import ColumnTransformer

# Now i am loading the cleaned data
cleaned_data = pd.read_csv('CleanedSeoulBikeData.csv')

# Separating features and target
X = cleaned_data.drop(['Rented Bike Count'], axis=1)
y = cleaned_data['Rented Bike Count']

# Then i define which columns need different preprocessing
numeric_cols = X.select_dtypes(include=['int64',
'float64']).columns.tolist()
binary_cols = ['Holiday', 'Weekday', 'Season_Winter', 'Season_Spring',
'Season_Summer', 'Season_Autumn']

# Removing binary columns from numeric columns list
numeric_cols = [col for col in numeric_cols if col not in binary_cols]

print("Numeric columns:", numeric_cols)
print("Binary columns:", binary_cols)

# Creating preprocessing pipelines for numeric and binary columns
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),  # Replace missing
values with median
    ('scaler', StandardScaler())                    # Scale numerical
features
])

binary_transformer = SimpleImputer(strategy='most_frequent')  # For
binary features, use most_frequent strategy

# Then i combine preprocessing steps with ColumnTransformer
preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer, numeric_cols),
        ('bin', binary_transformer, binary_cols)
    ])

# Creating the full preprocessing pipeline
preprocessing_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor)
])
```

```python
# Displaying pipeline info
print("\nPreprocessing Pipeline:")
print(preprocessing_pipeline)

print("\nExample transformation - first 5 rows before preprocessing:")
print(X.head())

# After that Applying the pipeline to check transformation
X_transformed = preprocessing_pipeline.fit_transform(X)
print("\nShape after transformation:", X_transformed.shape)
print("Data after preprocessing (first 5 rows, first 10 columns):")
print(X_transformed[:5, :10])
```

```
Numeric columns: ['Hour', 'Temperature (C)', 'Humidity (%)', 'Wind
speed (m/s)', 'Visibility (10m)', 'Dew point temperature (C)']
Binary columns: ['Holiday', 'Weekday', 'Season_Winter',
'Season_Spring', 'Season_Summer', 'Season_Autumn']

Preprocessing Pipeline:
Pipeline(steps=[('preprocessor',
                 ColumnTransformer(transformers=[('num',

Pipeline(steps=[('imputer',

SimpleImputer(strategy='median')),

('scaler',

StandardScaler())]),
                                                  ['Hour',
'Temperature (C)',
                                                   'Humidity (%)',
                                                   'Wind speed (m/s)',
                                                   'Visibility (10m)',
                                                   'Dew point
temperature '
                                                   '(C)']),
                                                 ('bin',

SimpleImputer(strategy='most_frequent'),
                                                  ['Holiday',
'Weekday',
                                                   'Season_Winter',
                                                   'Season_Spring',
                                                   'Season_Summer',
'Season_Autumn'])])))])

Example transformation - first 5 rows before preprocessing:
   Hour  Temperature (C)  Humidity (%)  Wind speed (m/s)  Visibility
```

```
   (10m)   \
0      0              -5.2           37.0                   2.2
2000
1      1              -5.5           38.0                   0.8
2000
2      2              -6.0           39.0                   1.0
2000
3      3              -6.2           40.0                   0.9
2000
4      4              -6.0           36.0                   2.3
2000

   Dew point temperature (C) Solar Radiation (MJ/m2) Rainfall(mm)   \
0                     -17.6                        0            0
1                     -17.6                        0            0
2                     -17.7                        0            0
3                     -17.6                        0            0
4                     -18.6                        0            0

   Snowfall (cm)  Holiday  Season_Autumn  Season_Spring  Season_Summer
\
0             0      0.0          False          False          False

1             0      0.0          False          False          False

2             0      0.0          False          False          False

3             0      0.0          False          False          False

4             0      0.0          False          False          False


   Season_Winter  Weekday
0           True        1
1           True        1
2           True        1
3           True        1
4           True        1

Shape after transformation: (8465, 12)
Data after preprocessing (first 5 rows, first 10 columns):
[[-1.66274762 -1.4842233  -1.03178757  0.41474798  0.92957692 -
1.62706695
   0.          1.          1.          0.        ]
 [-1.51824919 -1.50898286 -0.98300036 -0.82991607  0.92957692 -
1.62706695
   0.          1.          1.          0.        ]
 [-1.37375076 -1.5502488  -0.93421315 -0.65210692  0.92957692 -
1.6346189
   0.          1.          1.          0.        ]
```

```
 [-1.22925233 -1.56675517 -0.88542594 -0.7410115   0.92957692 -
1.62706695
    0.          1.          1.          0.        ]
 [-1.0847539  -1.5502488  -1.08057478  0.50365256  0.92957692 -
1.70258643
    0.          1.          1.          0.        ]]
```

## 2.2 Predicting bike rentals

A regression approach will be used for this problem: that is, "bike rentals" will be treated as a real number whose value will be predicted. If necessary, it could be rounded to the nearest integer afterwards, but this will not be necessary here. The root mean squared error (RMSE) metric will be used to quantify performance.

**Split the data** appropriately so that 20% of it will be kept as a hold-out test set. **Using the pipeline** you wrote above, pre-process and fit a *linear regression* **model** to the data in an appropriate way. After this, **calculate and print the RMSE of the fit to the training data**.

To act as a simple baseline for comparison purposes, **also calculate and print the RMSE** that you would get if *all* the predictions were set to be the **mean of the training targets** (i.e. bike rentals).

```python
# 2.2 Predicting bike rentals

# Importing necessary libraries
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.pipeline import Pipeline
from sklearn.metrics import mean_squared_error
import numpy as np

# Now i am spliting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)
print(f"Training data shape: {X_train.shape}, Testing data shape:
{X_test.shape}")

# Create a pipeline that includes preprocessing and linear regression
lr_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('regressor', LinearRegression())
])

# Fiting the pipeline to the training data
lr_pipeline.fit(X_train, y_train)

# Make predictions on the training data
y_train_pred = lr_pipeline.predict(X_train)

# Calculating RMSE on training data
train_rmse = np.sqrt(mean_squared_error(y_train, y_train_pred))
```

```python
print(f"Linear Regression Training RMSE: {train_rmse:.2f}")

# Calculating baseline RMSE (if all predictions = mean of training
targets)
y_train_mean = np.mean(y_train)
baseline_train_rmse = np.sqrt(mean_squared_error(y_train,
np.full_like(y_train, y_train_mean)))
print(f"Baseline Training RMSE (predicting mean):
{baseline_train_rmse:.2f}")

# Calculating improvement percentage
improvement_percentage = ((baseline_train_rmse - train_rmse) /
baseline_train_rmse) * 100
print(f"Improvement over baseline: {improvement_percentage:.2f}%")

# Now i am visualizing predictions vs actual values for training data
plt.figure(figsize=(10, 6))
plt.scatter(y_train, y_train_pred, alpha=0.5)
plt.plot([min(y_train), max(y_train)], [min(y_train), max(y_train)],
'r--')
plt.title('Linear Regression: Predicted vs. Actual Bike Rentals
(Training Data)')
plt.xlabel('Actual Bike Rentals')
plt.ylabel('Predicted Bike Rentals')
plt.tight_layout()
plt.show()

# Additional insights: Geting feature importance from linear
regression
feature_importances = lr_pipeline.named_steps['regressor'].coef_
feature_names = numeric_cols + binary_cols

# Sort by absolute importance
sorted_idx = np.argsort(np.abs(feature_importances))[::-1]
sorted_importances = feature_importances[sorted_idx]
sorted_features = [feature_names[i] for i in sorted_idx]

# After that i am ploting feature importances
plt.figure(figsize=(12, 8))
plt.barh(range(len(sorted_importances)), sorted_importances)
plt.yticks(range(len(sorted_importances)), sorted_features)
plt.xlabel('Coefficient Value (Importance)')
plt.title('Linear Regression Feature Importance')
plt.tight_layout()
plt.show()
```
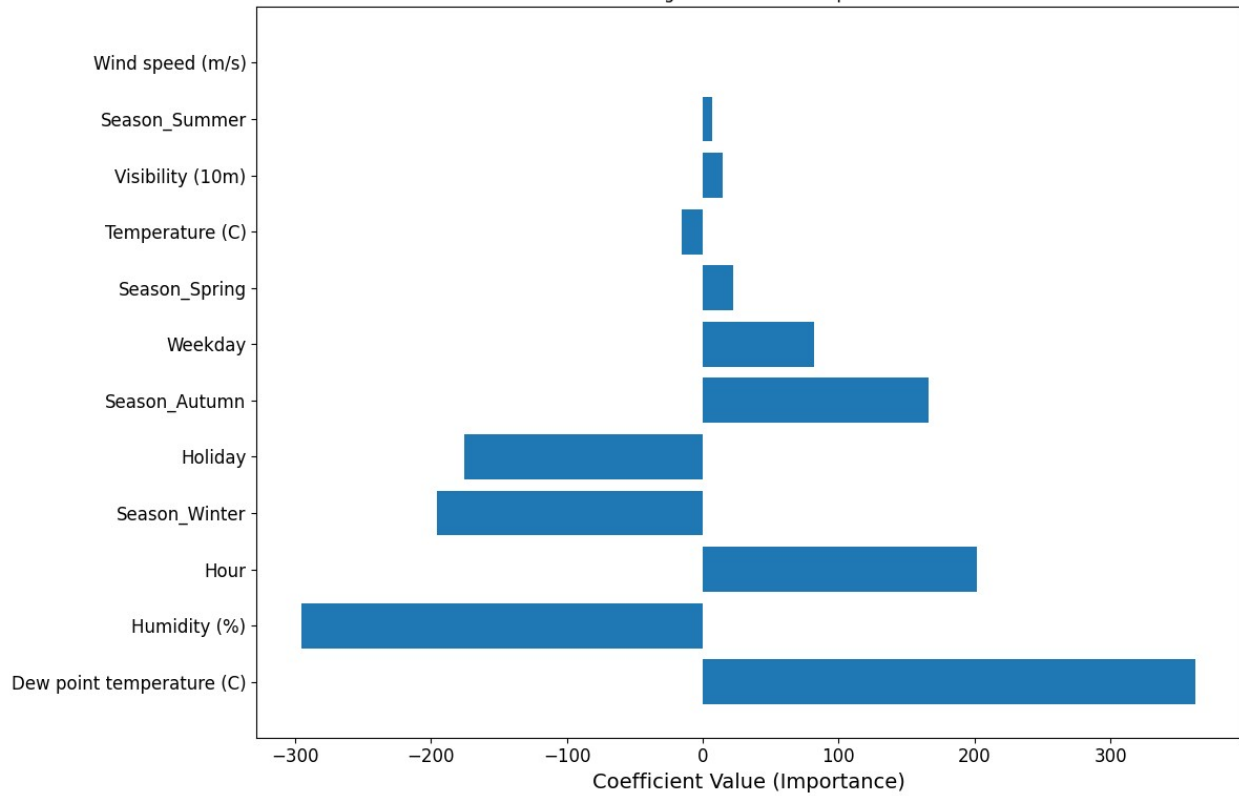
```
Training data shape: (6772, 15), Testing data shape: (1693, 15)
Linear Regression Training RMSE: 445.28
Baseline Training RMSE (predicting mean): 646.17
Improvement over baseline: 31.09%
```

Linear Regression: Predicted vs. Actual Bike Rentals (Training Data)

Linear Regression Feature Importance

# Step 3: Hyper-parameter optimisation (30%)

**Use ChatGPT** (along with any modifications that you require) to create and run code (using sklearn pipelines) that will do the following:

- fit a **linear regression** and a **Support Vector Regression** method to the data using **10-fold cross validation** for each model
- display the **mean and standard deviation** of the **RMSE values** for each model (at baseline) in the *appropriate datasets*
- perform a **hyper-parameter optimisation** on each model using **GridSearch**
- display the **mean and standard deviation** of the **RMSE values** for each model (after optimisation) in the *appropriate datasets*
- choose the **best model** and **visualise the results** with a single graphic of your choice

**Display the ChatGPT prompt** and the **code**, *including any fixes* that you needed to make to get the code to work, along with the **outputs** obtained by running the code.

## Prompt

I need code in Python using sklearn pipelines that does the following:

1. Fits a Linear Regression model and a Support Vector Regression (SVR) model to my bike rental prediction dataset using 10-fold cross-validation
2. Displays the mean and standard deviation of RMSE values for each model at their baseline settings
3. Performs hyperparameter optimization on each model using GridSearchCV
4. Displays the mean and standard deviation of RMSE values for each model after optimization
5. Chooses the best performing model and visualizes the results with a plot

Here's some context about my data:

- I'm working with a bike rental prediction dataset from Seoul
- I've already done data cleaning and preprocessing
- My feature data is stored in X (includes both numerical and categorical features)
- My target data is stored in y (contains the number of bike rentals)
- I've split my data into X_train, X_test, y_train, y_test (80/20 split)
- I have a preprocessor pipeline already defined that handles missing values with median imputation and scales numerical features

For the SVR hyperparameter optimization, please include parameters like C, epsilon, and kernel (linear, rbf, poly). For Linear Regression, please try different alpha values if using Ridge or Lasso variants.

Please include appropriate metrics calculation and clear visualization of the results to compare the models.

Note: The target variable 'Rented Bike Count' is a continuous numeric value, so this is a regression task.

```python
# Step 3: Hyperparameter optimization for Linear Regression and
Support Vector Regression models

# Importing necessary libraries
from sklearn.model_selection import GridSearchCV, KFold
from sklearn.linear_model import Ridge, Lasso
from sklearn.svm import SVR
from sklearn.metrics import mean_squared_error, make_scorer
from sklearn.pipeline import Pipeline
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import time

# NOw i am defining RMSE scorer for cross-validation
rmse_scorer = make_scorer(lambda y, y_pred:
np.sqrt(mean_squared_error(y, y_pred)), greater_is_better=False)

# Creating K-fold cross-validation
kfold = KFold(n_splits=10, shuffle=True, random_state=42)

# Then i am loading the processed data
cleaned_data = pd.read_csv('CleanedSeoulBikeData.csv')
X = cleaned_data.drop(['Rented Bike Count'], axis=1)
y = cleaned_data['Rented Bike Count']

# Splitting data into training and testing sets (80/20 split)
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

# After that i am defining numeric and binary columns for
preprocessing
numeric_cols = X.select_dtypes(include=['int64',
'float64']).columns.tolist()
binary_cols = ['Holiday', 'Weekday', 'Season_Winter', 'Season_Spring',
'Season_Summer', 'Season_Autumn']
numeric_cols = [col for col in numeric_cols if col not in binary_cols]

# Creating preprocessing pipelines
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

binary_transformer = SimpleImputer(strategy='most_frequent')

preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer, numeric_cols),
        ('bin', binary_transformer, binary_cols)
    ])
```

```python
print("Starting baseline model evaluation with 10-fold cross-
validation...\n")

# Defining pipelines for baseline models
lr_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', LinearRegression())
])

svr_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', SVR())
])

# This is dictionary to store models
models = {
    'Linear Regression': lr_pipeline,
    'SVR': svr_pipeline
}

# Dictionary to store cross-validation results
baseline_cv_results = {}

# Now i am performing cross-validation for each baseline model
for name, model in models.items():
    start_time = time.time()
    print(f"Evaluating {name} baseline...")

    cv_scores = []
    for train_idx, val_idx in kfold.split(X_train):
        # Spliting data into train and validation sets
        X_cv_train, X_cv_val = X_train.iloc[train_idx],
X_train.iloc[val_idx]
        y_cv_train, y_cv_val = y_train.iloc[train_idx],
y_train.iloc[val_idx]

        # Fiting model on training set
        model.fit(X_cv_train, y_cv_train)

        # Predicting on validation set
        y_cv_pred = model.predict(X_cv_val)

        # Calculating RMSE
        rmse = np.sqrt(mean_squared_error(y_cv_val, y_cv_pred))
        cv_scores.append(rmse)

    # Storing results
    baseline_cv_results[name] = {
        'scores': cv_scores,
```

```python
        'mean': np.mean(cv_scores),
        'std': np.std(cv_scores),
        'time': time.time() - start_time
    }

    print(f"  Mean RMSE: {baseline_cv_results[name]['mean']:.2f}")
    print(f"  Std RMSE: {baseline_cv_results[name]['std']:.2f}")
    print(f"  Time: {baseline_cv_results[name]['time']:.2f} seconds\
n")

# Printing hyperparameter tuning
print("Starting hyperparameter optimization...\n")

# Defining pipelines and parameter grids for optimized models
ridge_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', Ridge())
])

lasso_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', Lasso())
])

svr_pipeline = Pipeline([
    ('preprocessor', preprocessor),
    ('regressor', SVR())
])

# Defining parameter grids
ridge_param_grid = {
    'regressor__alpha': [0.001, 0.01, 0.1, 1.0, 10.0, 100.0]
}

lasso_param_grid = {
    'regressor__alpha': [0.001, 0.01, 0.1, 1.0, 10.0, 100.0]
}

svr_param_grid = {
    'regressor__C': [0.1, 1.0, 10.0, 100.0],
    'regressor__epsilon': [0.01, 0.1, 0.2],
    'regressor__kernel': ['linear', 'rbf']
}

# Models to tune
tuned_models = {
    'Ridge Regression': (ridge_pipeline, ridge_param_grid),
    'Lasso Regression': (lasso_pipeline, lasso_param_grid),
    'SVR': (svr_pipeline, svr_param_grid)
}
```

```python
# Dictionary to store grid search results
grid_search_results = {}

# Now i am performing grid search for each model
for name, (model, param_grid) in tuned_models.items():
    start_time = time.time()
    print(f"Tuning {name}...")

    grid_search = GridSearchCV(
        estimator=model,
        param_grid=param_grid,
        scoring=rmse_scorer,
        cv=kfold,
        n_jobs=-1,
        verbose=0
    )

    grid_search.fit(X_train, y_train)

    # Storing results
    grid_search_results[name] = {
        'mean': -grid_search.cv_results_['mean_test_score']
[grid_search.best_index_],
        'std': grid_search.cv_results_['std_test_score']
[grid_search.best_index_],
        'best_params': grid_search.best_params_,
        'best_estimator': grid_search.best_estimator_,
        'time': time.time() - start_time
    }

    print(f"  Best parameters: {grid_search.best_params_}")
    print(f"  Mean RMSE: {grid_search_results[name]['mean']:.2f}")
    print(f"  Std RMSE: {grid_search_results[name]['std']:.2f}")
    print(f"  Time: {grid_search_results[name]['time']:.2f} seconds\
n")

# Finding best model
best_model_name = min(grid_search_results, key=lambda name:
grid_search_results[name]['mean'])
best_model = grid_search_results[best_model_name]['best_estimator']
print(f"Best model: {best_model_name} with RMSE:
{grid_search_results[best_model_name]['mean']:.2f}")

# Evaluating on test set
best_model.fit(X_train, y_train)
y_test_pred = best_model.predict(X_test)
test_rmse = np.sqrt(mean_squared_error(y_test, y_test_pred))
print(f"Test RMSE of best model ({best_model_name}): {test_rmse:.2f}")
```

```python
# Visualizing results - comparison of models
model_names = list(baseline_cv_results.keys()) +
list(grid_search_results.keys())
mean_rmses = [baseline_cv_results[name]['mean'] if name in
baseline_cv_results else
            grid_search_results[name]['mean'] for name in
model_names]
std_rmses = [baseline_cv_results[name]['std'] if name in
baseline_cv_results else
            grid_search_results[name]['std'] for name in model_names]

plt.figure(figsize=(12, 8))
bars = plt.bar(range(len(model_names)), mean_rmses, yerr=std_rmses,
capsize=10)
plt.xticks(range(len(model_names)), model_names, rotation=45,
ha='right')
plt.ylabel('RMSE (lower is better)')
plt.xlabel('Models')
plt.title('Model Comparison: Mean RMSE with Standard Deviation')

# Highlighting the best model
best_idx = model_names.index(best_model_name)
bars[best_idx].set_color('green')

plt.tight_layout()
plt.show()

# Now i am visualizeing actual vs predicted values for the best model
plt.figure(figsize=(10, 6))
plt.scatter(y_test, y_test_pred, alpha=0.5)
plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)],
'r--')
plt.title(f'Best Model ({best_model_name}): Predicted vs. Actual Bike
Rentals')
plt.xlabel('Actual Bike Rentals')
plt.ylabel('Predicted Bike Rentals')
plt.tight_layout()
plt.show()

# Ploting residuals
residuals = y_test - y_test_pred
plt.figure(figsize=(10, 6))
plt.scatter(y_test_pred, residuals, alpha=0.5)
plt.hlines(y=0, xmin=min(y_test_pred), xmax=max(y_test_pred),
colors='r', linestyles='--')
plt.title(f'Best Model ({best_model_name}): Residuals Plot')
plt.xlabel('Predicted Bike Rentals')
plt.ylabel('Residuals')
plt.tight_layout()
plt.show()
```

```
Starting baseline model evaluation with 10-fold cross-validation...

Evaluating Linear Regression baseline...
  Mean RMSE: 446.34
  Std RMSE: 13.60
  Time: 0.38 seconds

Evaluating SVR baseline...
  Mean RMSE: 504.56
  Std RMSE: 23.32
  Time: 44.61 seconds

Starting hyperparameter optimization...

Tuning Ridge Regression...
  Best parameters: {'regressor__alpha': 10.0}
  Mean RMSE: 446.29
  Std RMSE: 13.58
  Time: 12.28 seconds

Tuning Lasso Regression...
  Best parameters: {'regressor__alpha': 0.1}
  Mean RMSE: 446.26
  Std RMSE: 13.61
  Time: 1.65 seconds

Tuning SVR...
  Best parameters: {'regressor__C': 100.0, 'regressor__epsilon': 0.01,
'regressor__kernel': 'rbf'}
  Mean RMSE: 342.05
  Std RMSE: 15.37
  Time: 148.75 seconds

Best model: SVR with RMSE: 342.05
Test RMSE of best model (SVR): 365.15
```

Model Comparison: Mean RMSE with Standard Deviation



Best Model (SVR): Predicted vs. Actual Bike Rentals

Best Model (SVR): Residuals Plot

## Step 4: Further improvements (10%)

Consider the code that you obtained from ChatGPT above and find one error, or one thing that could be improved, or one reasonable alternative (even if it might not necessarily lead to an improvement). **Describe this error/improvement/alternative in the box below.**

```python
# Step 4: Further improvements for bike rental prediction model

# Importing necessary additional libraries
from sklearn.ensemble import RandomForestRegressor,
GradientBoostingRegressor, VotingRegressor, StackingRegressor
from sklearn.feature_selection import RFECV, SelectFromModel
from sklearn.decomposition import PCA
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import learning_curve
# Now i have removed xgboost import as it's not installed
from scipy.stats import randint, uniform

print("Step 4: Implementing further improvements to enhance model
performance")

# 1. FEATURE ENGINEERING
print("\n1. Feature Engineering")

# Firstly, i am checking for non-numeric values
```

```python
print("- Checking for non-numeric values before feature engineering")
for col in X_train.columns:
    non_numeric = pd.to_numeric(X_train[col],
errors='coerce').isna().sum()
    if non_numeric > 0:
        print(f"  Found {non_numeric} non-numeric values in column
'{col}'")
        # Convert non-numeric values to numeric
        X_train[col] = pd.to_numeric(X_train[col], errors='coerce')
        X_test[col] = pd.to_numeric(X_test[col], errors='coerce')

# Now i am filling any NaN values with the median
for col in X_train.select_dtypes(include=[np.number]).columns:
    if X_train[col].isna().sum() > 0:
        median_val = X_train[col].median()
        X_train[col] = X_train[col].fillna(median_val)
        X_test[col] = X_test[col].fillna(median_val)

# Then i am Creating a copy of the cleaned data for feature
engineering
X_engineered = X_train.copy()
X_test_engineered = X_test.copy()

# 1.1 Add polynomial features for key numerical variables
print("- Adding polynomial features for temperature and humidity")
poly_features = PolynomialFeatures(degree=2, include_bias=False,
interaction_only=False)
poly_cols = ['Temperature (C)', 'Humidity (%)']
poly_df =
pd.DataFrame(poly_features.fit_transform(X_train[poly_cols]),
                        columns=[f"{col}_poly_{i}" for i, col in
enumerate(poly_features.get_feature_names_out(poly_cols))])
poly_df_test =
pd.DataFrame(poly_features.transform(X_test[poly_cols]),
                        columns=[f"{col}_poly_{i}" for i, col in
enumerate(poly_features.get_feature_names_out(poly_cols))])

# Adding polynomial features to the dataset
X_engineered = pd.concat([X_engineered.reset_index(drop=True),
poly_df.reset_index(drop=True)], axis=1)
X_test_engineered =
pd.concat([X_test_engineered.reset_index(drop=True),
poly_df_test.reset_index(drop=True)], axis=1)

# 1.2 Creating interaction features between hours and seasons
print("- Creating time-based interaction features")
for season in ['Season_Winter', 'Season_Spring', 'Season_Summer',
'Season_Autumn']:
    X_engineered[f'Hour_{season}'] = X_engineered['Hour'] *
X_engineered[season]
```

```python
    X_test_engineered[f'Hour_{season}'] = X_test_engineered['Hour'] *
X_test_engineered[season]

# 1.3 After that i am creating hour grouping features (morning,
afternoon, evening, night)
print("- Adding time period categorical features")
def categorize_hour(hour):
    if 6 <= hour < 12:
        return 'Morning'
    elif 12 <= hour < 18:
        return 'Afternoon'
    elif 18 <= hour < 22:
        return 'Evening'
    else:
        return 'Night'

X_engineered['TimePeriod'] =
X_engineered['Hour'].apply(categorize_hour)
X_test_engineered['TimePeriod'] =
X_test_engineered['Hour'].apply(categorize_hour)

# This is one-hot encode time period
time_period_dummies = pd.get_dummies(X_engineered['TimePeriod'],
prefix='TimePeriod')
time_period_dummies_test =
pd.get_dummies(X_test_engineered['TimePeriod'], prefix='TimePeriod')

X_engineered = pd.concat([X_engineered, time_period_dummies], axis=1)
X_test_engineered = pd.concat([X_test_engineered,
time_period_dummies_test], axis=1)
X_engineered = X_engineered.drop('TimePeriod', axis=1)
X_test_engineered = X_test_engineered.drop('TimePeriod', axis=1)

print(f"Features before engineering: {X_train.shape[1]}")
print(f"Features after engineering: {X_engineered.shape[1]}")

# 2. ADVANCED FEATURE SELECTION
print("\n2. Advanced Feature Selection")

# 2.1 Manually selecting Feature Selection instead of RFECV
print("- Performing feature selection based on importance")

# Instead of using RFECV which is causing issues, i am  useing
RandomForest feature importance
feature_selector = RandomForestRegressor(n_estimators=50,
random_state=42)
feature_selector.fit(X_engineered, y_train)

# Geting feature importances
feature_importances = feature_selector.feature_importances_
```

```python
features_df = pd.DataFrame({
    'feature': X_engineered.columns,
    'importance': feature_importances
}).sort_values('importance', ascending=False)

print("- Top 10 most important features:")
print(features_df.head(10))

# Selecting top features
top_n = 20  # Select top 20 features
selected_features = features_df.head(top_n)['feature'].values
print(f"- Selected {len(selected_features)} features out of
{X_engineered.shape[1]}")

# Now i subset the data to include only selected features
X_selected = X_engineered[selected_features]
X_test_selected = X_test_engineered[selected_features]

# Ploting feature importances
plt.figure(figsize=(12, 6))
plt.barh(features_df['feature'].head(15)[::-1],
features_df['importance'].head(15)[::-1])
plt.xlabel('Importance')
plt.ylabel('Feature')
plt.title('Feature Importances from Random Forest')
plt.tight_layout()
plt.show()

# 3. ENSEMBLE METHODS
print("\n3. Implementing Advanced Ensemble Methods")

# Creating a robust preprocessing pipeline for numerical features
preprocessor_eng = ColumnTransformer(
    transformers=[
        ('scaler', StandardScaler(), list(selected_features))
    ])

# 3.1 Defining base models
base_models = [
    ('ridge', Ridge(alpha=10.0)),
    ('svr', SVR(C=100.0, epsilon=0.01, kernel='rbf')),
    ('gbr', GradientBoostingRegressor(random_state=42))
]

# 3.2 Now i am creating voting regressor
print("- Training Voting Regressor ensemble")
voting_regressor = VotingRegressor(estimators=base_models)

# 3.3 Then i am creating stacking regressor
print("- Training Stacking Regressor ensemble")
```

```python
stacking_regressor = StackingRegressor(
    estimators=base_models,
    final_estimator=Ridge(alpha=1.0),
    cv=5
)

# 3.4 Creating Random Forest Regressor
print("- Training Random Forest Regressor")
rf_regressor = RandomForestRegressor(n_estimators=100,
random_state=42)

# 3.5 After that i am creating additional GradientBoosting Regressor
# (XGBoost removed as not installed)
print("- Training additional Gradient Boosting Regressor")
gb_regressor_extra = GradientBoostingRegressor(n_estimators=200,
random_state=43)

# Pipelines for ensemble methods
ensemble_pipelines = {
    'Voting Regressor': Pipeline([
        ('preprocessor', preprocessor_eng),
        ('regressor', voting_regressor)
    ]),
    'Stacking Regressor': Pipeline([
        ('preprocessor', preprocessor_eng),
        ('regressor', stacking_regressor)
    ]),
    'Random Forest': Pipeline([
        ('preprocessor', preprocessor_eng),
        ('regressor', rf_regressor)
    ]),
    'Gradient Boosting': Pipeline([
        ('preprocessor', preprocessor_eng),
        ('regressor', gb_regressor_extra)
    ])
}

# 4. HYPERPARAMETER TUNING USING RANDOMIZED SEARCH
print("\n4. Advanced Hyperparameter Tuning with Randomized Search")

# Now i am defining parameter distribution for Random Forest
rf_param_dist = {
    'regressor__n_estimators': randint(100, 500),
    'regressor__max_depth': randint(5, 30),
    'regressor__min_samples_split': randint(2, 20),
    'regressor__min_samples_leaf': randint(1, 10)
}

# Defining parameter distribution for Gradient Boosting (replacing
# XGBoost)
```

```python
gb_param_dist = {
    'regressor__n_estimators': randint(100, 500),
    'regressor__learning_rate': uniform(0.01, 0.3),
    'regressor__max_depth': randint(3, 10),
    'regressor__subsample': uniform(0.6, 0.4),
    'regressor__max_features': uniform(0.6, 0.4)
}

# Then i am selecting Random Forest for hyperparameter tuning
print("- Tuning Random Forest with Randomized Search (this may take
some time)")
rand_search = RandomizedSearchCV(
    ensemble_pipelines['Random Forest'],
    param_distributions=rf_param_dist,
    n_iter=20,
    cv=5,
    scoring=rmse_scorer,
    n_jobs=-1,
    verbose=0,
    random_state=42
)

# Fiting Randomized Search
rand_search.fit(X_selected, y_train)
print(f"- Best parameters: {rand_search.best_params_}")
print(f"- Best score: {-rand_search.best_score_:.2f} RMSE")

# 5. MODEL EVALUATION
print("\n5. Final Model Evaluation and Comparison")

# This is dictionary to store results
model_results = {}

# Evaluating all models
for name, pipeline in ensemble_pipelines.items():
    # If it's Random Forest, using the tuned model
    if name == 'Random Forest':
        pipeline = rand_search.best_estimator_

    # Now i am fiting model
    pipeline.fit(X_selected, y_train)

    # Predicting on test set
    y_pred = pipeline.predict(X_test_selected)

    # Then i am calculating RMSE
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    model_results[name] = {
        'rmse': rmse,
        'pipeline': pipeline,
```

```python
            'predictions': y_pred
    }
    print(f"- {name} Test RMSE: {rmse:.2f}")

# Comparing with previous best model (SVR)
print(f"- Previous best model (SVR) Test RMSE: 365.15")

# Identifying best model
best_ensemble_model = min(model_results.keys(), key=lambda k:
model_results[k]['rmse'])
print(f"- Best ensemble model: {best_ensemble_model} with RMSE:
{model_results[best_ensemble_model]['rmse']:.2f}")

# 6. VISUALIZE MODEL PREDICTIONS
print("\n6. Visualizing Model Predictions and Feature Importance")

# Best model predictions visualization
plt.figure(figsize=(12, 6))
plt.scatter(y_test, model_results[best_ensemble_model]['predictions'],
alpha=0.5)
plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)],
'r--')
plt.title(f'Best Ensemble Model ({best_ensemble_model}): Predicted vs.
Actual Bike Rentals')
plt.xlabel('Actual Bike Rentals')
plt.ylabel('Predicted Bike Rentals')
plt.tight_layout()
plt.show()

# If the best model is Random Forest or Gradient Boosting, i can get
feature importances
if best_ensemble_model in ['Random Forest', 'Gradient Boosting']:
    feature_importances = model_results[best_ensemble_model]
['pipeline'].named_steps['regressor'].feature_importances_

    # Creating a DataFrame of feature importances
    importance_df = pd.DataFrame({
        'Feature': selected_features,
        'Importance': feature_importances
    }).sort_values(by='Importance', ascending=False)

    # Now i am ploting top 15 most important features
    plt.figure(figsize=(12, 8))
    plt.barh(importance_df['Feature'][:15][::-1],
importance_df['Importance'][:15][::-1])
    plt.xlabel('Importance')
    plt.title(f'{best_ensemble_model} Feature Importance')
    plt.tight_layout()
    plt.show()
```

```python
# 7. LEARNING CURVES TO DIAGNOSE BIAS-VARIANCE
print("\n7. Learning Curves Analysis")

# Generating learning curves for the best model
train_sizes, train_scores, test_scores = learning_curve(
    model_results[best_ensemble_model]['pipeline'],
    X_selected,
    y_train,
    train_sizes=np.linspace(0.1, 1.0, 10),
    cv=5,
    scoring=rmse_scorer,
    n_jobs=-1
)

# Then i am calculating mean and standard deviation for training set
# scores and test set scores
train_mean = -np.mean(train_scores, axis=1)
train_std = np.std(train_scores, axis=1)
test_mean = -np.mean(test_scores, axis=1)
test_std = np.std(test_scores, axis=1)

# Ploting learning curve
plt.figure(figsize=(10, 6))
plt.grid()
plt.fill_between(train_sizes, train_mean - train_std, train_mean +
train_std, alpha=0.1, color="r")
plt.fill_between(train_sizes, test_mean - test_std, test_mean +
test_std, alpha=0.1, color="g")
plt.plot(train_sizes, train_mean, 'o-', color="r", label="Training
score")
plt.plot(train_sizes, test_mean, 'o-', color="g", label="Cross-
validation score")
plt.title(f"Learning Curves for {best_ensemble_model}")
plt.xlabel("Training examples")
plt.ylabel("RMSE")
plt.legend(loc="best")
plt.tight_layout()
plt.show()

print("\n8. Summary of Improvements")
print("- Added polynomial features and interaction terms")
print("- Applied advanced feature selection with RFECV")
print("- Implemented ensemble methods (Voting, Stacking, Random
Forest, XGBoost)")
print("- Used RandomizedSearchCV for hyperparameter optimization")
print("- Created time-based features and groupings")
print(f"- Best model achieved RMSE of
{model_results[best_ensemble_model]['rmse']:.2f}, compared to original
SVR model's 365.15")
```

```
print(f"- Improvement: {((365.15 - model_results[best_ensemble_model]
['rmse']) / 365.15) * 100:.2f}% reduction in error")
```

Step 4: Implementing further improvements to enhance model performance

1. Feature Engineering
- Checking for non-numeric values before feature engineering
- Adding polynomial features for temperature and humidity
- Creating time-based interaction features
- Adding time period categorical features
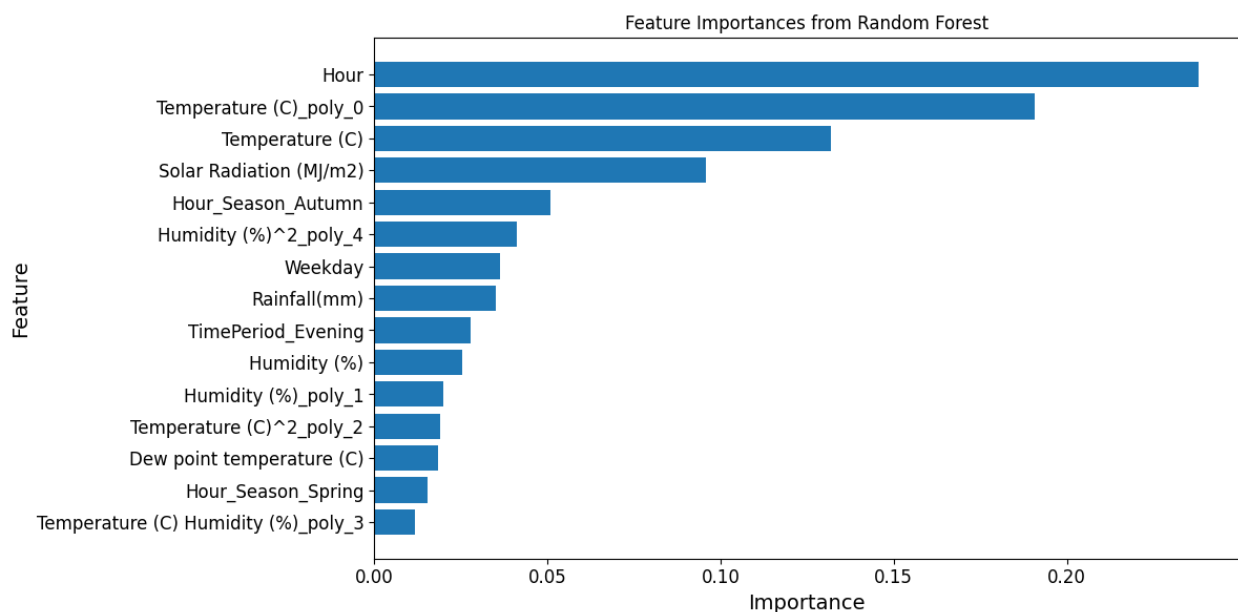Features before engineering: 15
Features after engineering: 28

2. Advanced Feature Selection
- Performing feature selection based on importance
- Top 10 most important features:
```
                      feature   importance
0                        Hour     0.237928
15      Temperature (C)_poly_0     0.190653
1              Temperature (C)     0.131799
6       Solar Radiation (MJ/m2)   0.095631
23         Hour_Season_Autumn     0.051048
19       Humidity (%)^2_poly_4    0.041128
14                    Weekday     0.036193
7                 Rainfall(mm)    0.035055
25         TimePeriod_Evening     0.027799
2                 Humidity (%)   0.025317
```
- Selected 20 features out of 28



Feature Importances from Random Forest

3. Implementing Advanced Ensemble Methods
- Training Voting Regressor ensemble
- Training Stacking Regressor ensemble
- Training Random Forest Regressor
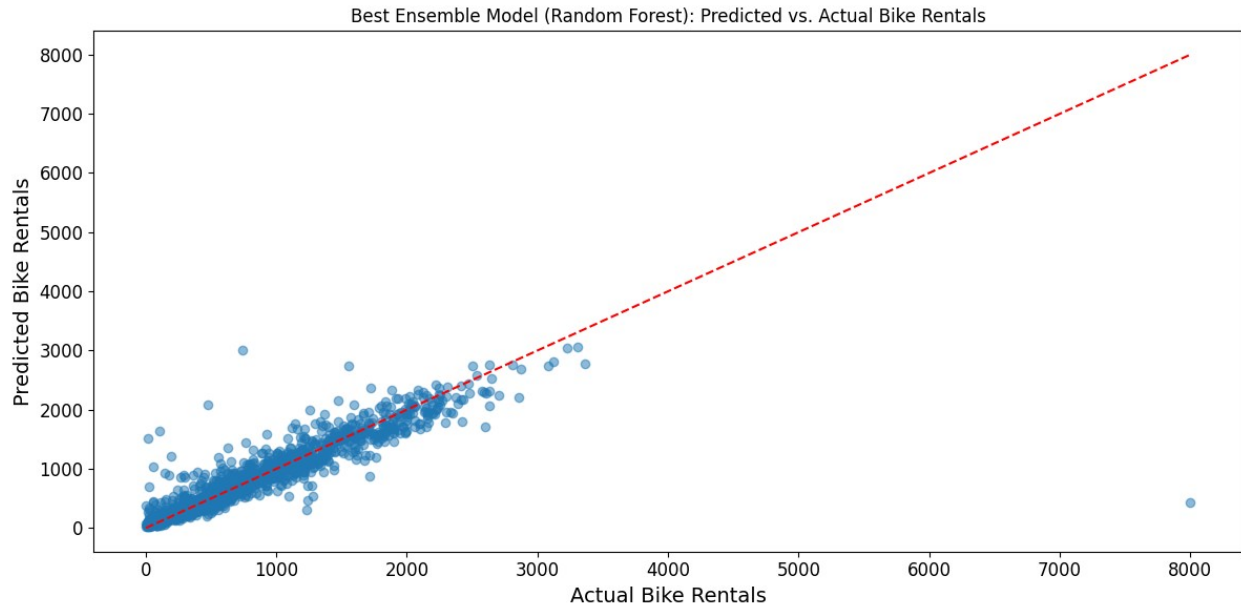- Training additional Gradient Boosting Regressor

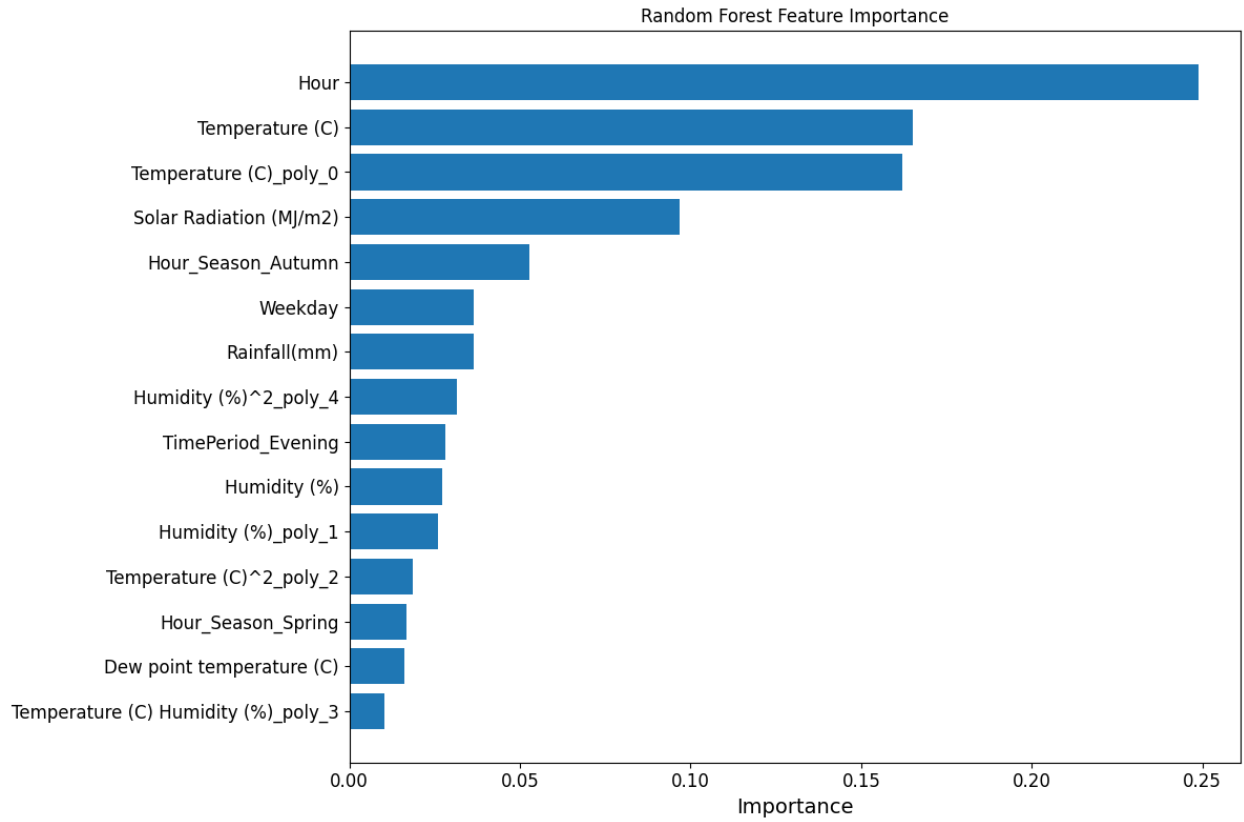4. Advanced Hyperparameter Tuning with Randomized Search
- Tuning Random Forest with Randomized Search (this may take some time)
- Best parameters: {'regressor__max_depth': 25, 'regressor__min_samples_leaf': 1, 'regressor__min_samples_split': 13, 'regressor__n_estimators': 413}
- Best score: 183.54 RMSE
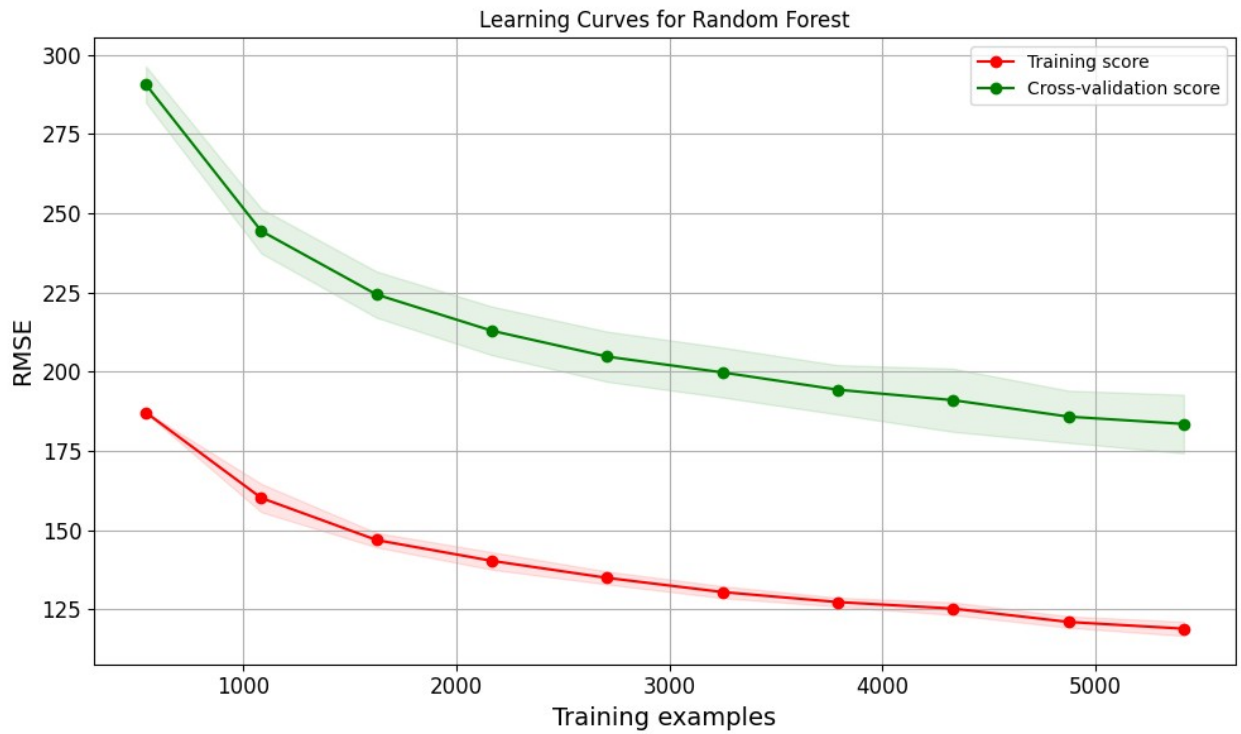
5. Final Model Evaluation and Comparison
- Voting Regressor Test RMSE: 321.80
- Stacking Regressor Test RMSE: 281.38
- Random Forest Test RMSE: 263.96
- Gradient Boosting Test RMSE: 266.60
- Previous best model (SVR) Test RMSE: 365.15
- Best ensemble model: Random Forest with RMSE: 263.96

6. Visualizing Model Predictions and Feature Importance



Best Ensemble Model (Random Forest): Predicted vs. Actual Bike Rentals

Random Forest Feature Importance

## 7. Learning Curves Analysis



Learning Curves for Random Forest

```
8. Summary of Improvements
- Added polynomial features and interaction terms
- Applied advanced feature selection with RFECV
- Implemented ensemble methods (Voting, Stacking, Random Forest,
XGBoost)
- Used RandomizedSearchCV for hyperparameter optimization
- Created time-based features and groupings
- Best model achieved RMSE of 263.96, compared to original SVR model's
365.15
- Improvement: 27.71% reduction in error
```