

Import software libraries

```
import sys                                # Read sys parameters
import os                                 # Interact with the operating system
import numpy as np                       # Work with multi-dimensional arrays and matrices
import pandas as pd                      # Manipulate and analyze data
import matplotlib                         # Create charts
import matplotlib.pyplot as plt          # Perform scientific computing and advanced mathematics
import scipy as sp                       # Perform data mining and analysis
import sklearn                           # Perform data visualization
import seaborn as sns

# Summarize software libraries used
print("Libraries used in this project:")
print("- Numpy {}".format(np.__version__))
print("- Pandas {}".format(pd.__version__))
print("- Matplotlib {}".format(matplotlib.__version__))
print("- SciPy {}".format(sp.__version__))
print("- sklearn {}".format(sklearn.__version__))
print("- Python {}".format(sys.version))

Libraries used in this project:
- Numpy 1.19.2
- Pandas 1.1.3
- Matplotlib 3.3.2
- SciPy 1.5.2
- Sklearn 0.23.2
- Python 3.8.5 (default, Sep 3 2020, 21:29:08) [MSC v.1916 64 bit (AMD64)]
```

Load the dataset

```
In [2]: # Load the dataset as a pandas DataFrame from ./seoul_bike_data/seoul_bike_data.csv
df = pd.read_csv("./seoul_bike_data.csv")
df

Out[2]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
0	254	-5.2	37	2.2	2000	-17.6	0.0	0.0	0.0
1	204	-5.5	38	0.8	2000	-17.6	0.0	0.0	0.0
2	173	-6.0	39	1.0	2000	-17.7	0.0	0.0	0.0
3	107	-6.2	40	0.9	2000	-17.6	0.0	0.0	0.0
4	78	-6.0	36	2.3	2000	-18.6	0.0	0.0	0.0
...
8399	1003	4.2	34	2.6	1894	-10.3	0.0	0.0	0.0
8390	764	3.4	37	2.3	2000	-9.9	0.0	0.0	0.0
8391	694	2.6	39	0.3	1968	-9.9	0.0	0.0	0.0
8392	712	2.1	41	1.0	1859	-9.8	0.0	0.0	0.0
8393	584	1.9	43	1.3	1909	-9.3	0.0	0.0	0.0

8394 rows × 9 columns

Get acquainted with the dataset

```
In [3]: # View data types and see if there are missing entries.
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8394 entries, 0 to 8393
Data columns (total 9 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   bikes_rented  8394 non-null    int64
 1   temp          8394 non-null    float64
 2   humidity      8394 non-null    float64
 3   wind_speed    8394 non-null    float64
 4   visibility     8394 non-null    float64
 5   dew_temp      8394 non-null    float64
 6   solar_rad     8394 non-null    float64
 7   rainfall      8394 non-null    float64
 8   snowfall      8394 non-null    float64
dtypes: float64(6), int64(3)
memory usage: 590.3 KB

In [4]: df.isnull().sum()

Out[4]:
bikes_rented    0
humidity         0
wind_speed       0
visibility        0
dew_temp         0
solar_rad        0
rainfall         0
snowfall         0
dtype: int64
```

Show example records

```
In [5]: # View first 10 records.
df.head(10)

Out[5]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
0	254	-5.2	37	2.2	2000	-17.6	0.00	0.0	0.0
1	204	-5.5	38	0.8	2000	-17.6	0.00	0.0	0.0
2	173	-6.0	39	1.0	2000	-17.7	0.00	0.0	0.0
3	107	-6.2	40	0.9	2000	-17.6	0.00	0.0	0.0
4	78	-6.0	36	2.3	2000	-18.6	0.00	0.0	0.0
5	100	-6.4	37	1.5	2000	-18.7	0.00	0.0	0.0
6	181	-6.6	35	1.3	2000	-19.5	0.00	0.0	0.0
7	460	-7.4	38	0.9	2000	-19.3	0.00	0.0	0.0
8	490	-7.6	37	1.1	2000	-19.8	0.01	0.0	0.0
9	930	-6.5	27	0.5	1928	-22.4	0.23	0.0	0.0

Examine a general summary of statistics

```
In [6]: # View summary statistics (mean, standard deviation, min, max, etc.) for each feature.
df.describe()

Out[6]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
count	8394.000000	8394.000000	8394.000000	8394.000000	8394.000000	8394.000000	8394.000000	8394.000000	8394.000000
mean	731.374738	12.812099	58.074696	1.740481	1433.226590	3.964260	0.572427	0.149261	0.077949
std	643.616638	12.810977	20.483539	1.026341	609.802729	13.242399	0.870429	1.126075	0.445880
min	2.000000	-17.800000	0.000000	0.100000	27.000000	-30.600000	0.000000	0.000000	0.000000
25%	214.000000	4.200000	42.000000	1.000000	932.250000	-5.100000	0.000000	0.000000	0.000000
50%	546.000000	13.600000	57.000000	1.500000	1690.000000	4.800000	0.010000	0.000000	0.000000
75%	1088.000000	22.700000	74.000000	2.300000	2000.000000	15.200000	0.940000	0.000000	0.000000
max	3556.000000	39.400000	98.000000	7.400000	2000.000000	27.200000	3.520000	35.000000	8.800000

Look for columns that correlate with bikes_rented

```
In [7]: # View the correlation values for each feature compared to the label.
df.corr()

Out[7]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
bikes_rented	1.000000	0.563440	-0.201466	0.120961	0.213989	0.401160	0.272748	-0.128794	-0.151881
temp	0.563440	1.000000	0.165484	-0.044827	0.031410	0.914372	0.354692	0.052120	-0.218070
humidity	-0.201466	0.165484	1.000000	-0.336857	-0.549300	0.338730	-0.457904	0.237436	0.110487
wind_speed	0.120961	-0.044827	-0.336857	1.000000	0.184935	-0.182518	0.321812	-0.025538	-0.004840
visibility	0.213989	0.031410	-0.549300	0.184935	1.000000	-0.180199	0.154676	-0.169727	-0.123300
dew_temp	0.401160	0.914372	0.538730	-0.182518	-0.180199	1.000000	0.098152	0.127034	-0.149969
solar_rad	0.272748	0.354692	-0.457904	0.321812	0.154676	0.098152	1.000000	-0.074607	-0.073923
rainfall	-0.128794	0.052120	0.237436	-0.025538	-0.169727	0.127034	-0.074607	1.000000	0.008712
snowfall	-0.151881	-0.218070	0.110487	-0.004840	-0.123300	-0.149969	-0.073923	0.008712	1.000000

Visually analyze cross correlations

```
In [8]: # Use Seaborn to plot the correlation matrix as a heatmap.
plt.figure(figsize=(16,9))
sns.heatmap(df.corr(), cmap="viridis", annot=True, fmt=".2g", linewidths=2)
plt.show()

Out[8]:
```

Use histograms to visualize the distribution of all features

```
In [9]: # Use Matplotlib to plot distribution histograms for all features.

plt.figure(figsize=(20,20))
plt.subplot(4,2,1)
df["bikes_rented"].plot(kind="hist")
plt.title("Bikes Rented")
plt.subplot(4,2,2)
df["temp"].plot(kind="hist")
plt.title("temp")
plt.subplot(4,2,3)
df["humidity"].plot(kind="hist")
plt.title("humidity")
plt.subplot(4,2,4)
df["wind_speed"].plot(kind="hist")
plt.title("wind_speed")
plt.subplot(4,2,5)
df["visibility"].plot(kind="hist")
plt.title("visibility")
plt.subplot(4,2,6)
df["dew_temp"].plot(kind="hist")
plt.title("dew_temp")
plt.subplot(4,2,7)
df["solar_rad"].plot(kind="hist")
plt.title("solar_rad")
plt.subplot(4,2,8)
df["rainfall"].plot(kind="hist")
plt.title("rainfall")
plt.show()

Out[9]:
```

Split the data into training and testing sets and labels

```
In [10]: # Split the training and test datasets and their labels.

# Compare the number of rows and columns in the original data to the training and test sets.
from sklearn.model_selection import train_test_split

Out[10]:
df.shape

In [11]:
df.shape

Out[11]:
(8394, 9)

In [12]:
X = df.iloc[:,1:]
y = df.iloc[:,0]

Out[12]:
X.values, y.values

In [13]:
(array([[ -5.2, 37., 2.2, ..., 0., 0., 0., 0.],
       [ -5.5, 38., 0.8, ..., 0., 0., 0., 0.],
       [ -6., 39., 1., ..., 0., 0., 0., 0.],
       ...,
       [ 2.6, 39., 0.3, ..., 0., 0., 0., 0.],
       [ 3., 41., 1., ..., 0., 0., 0., 0.],
       [ 1.9, 43., 1.3, ..., 0., 0., 0., 0.]],
      array([254, 204, 173, ..., 694, 712, 584], dtype=int64))

In [14]:
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

Out[14]:
X_train.shape, y_train.shape

Out[15]:
(6715, 8), (6715, 1)
```

Build and test an initial linear regression model

```
In [16]: # Create a linear regression model.

# Fit the model using training data and labels.
from sklearn.linear_model import LinearRegression

In [17]:
lr = LinearRegression()

In [18]:
lr.fit(X_train, y_train)

Out[18]:
LinearRegression()
```

Use the holdout dataset to test the model

```
In [19]: # Print the regressor model's score using the test data and labels.

lr.score(X_test, y_test)

Out[19]:
0.4478627291734715
```

Compare the first ten predictions to actual values

```
In [20]: # Make predictions on the test set.

# View examples comparing actual bike rentals to predicted bike rentals.

In [21]:
lr_pred = lr.predict(X_test)

In [22]:
lr_pred[:10]

Out[22]:
array([ 857.323119, 751.8098849, 655.0806967, 1226.26361471,
        1590.38483993, 285.89222901, 1454.77870504, 711.79165875,
        325.36163583, 598.10222966])

In [23]:
y_test[:10]

Out[23]:
2417    946
2933    103
7268    1255
6342    2022
5697    1897
1259    184
5622    920
2431    616
1916    388
2650    688
Name: bikes_rented, dtype: int64
```

Identify outliers

```
In [24]: # Use Matplotlib to create box plot distributions for bikes_rented and wind_speed.

plt.boxplot(x=df.bikes_rented)
plt.title("Bikes Rented Boxplot")
plt.show()

Out[24]:
```

```
In [25]:
plt.boxplot(x=df.wind_speed)
plt.title("Wind Speed Boxplot")
plt.show()

Out[25]:
```

Examine data values in the outliers

```
In [26]: # Show rows that exceed 3,500 bikes rented.
df[df["bikes_rented"] > 3500]

Out[26]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
4743	3556	24.1	57	2.9	1301	15.0	0.56	0.0	0.0

```
In [27]: # Show rows with wind speed greater than 6 meters per second.
df[df["wind_speed"] > 6]

Out[27]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
909	146	0.7	77	6.7	692	-2.8	0.0	0.9	1.0
3108	1805	21.2	35	7.4	1992	5.1	1.8	0.0	0.0
3112	913	19.7	52	7.2	2000	9.5	0.2	0.0	0.0
3114	336	19.1	58	6.1	2000	10.6	0.0	0.0	0.0
151	133	17.5	70	7.3	1634	11.9	0.0	0.4	0.0
6230	49	25.3	70	6.9	925	19.4	0.0	0.5	0.0

8393 rows × 9 columns

```
In [28]: # Keep only the rows where number of bikes rented is less than 3,500.
df1 = df[df["bikes_rented"] < 3500]
df1

Out[28]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
0	254	-5.2	37	2.2	2000	-17.6	0.0	0.0	0.0
1	204	-5.5	38	0.8	2000	-17.6	0.0	0.0	0.0
2	173	-6.0	39	1.0	2000	-17.7	0.0	0.0	0.0
3	107	-6.2	40	0.9	2000	-17.6	0.0	0.0	0.0
4	78	-6.0	36	2.3	2000	-18.6	0.0	0.0	0.0
...
8399	1003	4.2	34	2.6	1894	-10.3	0.0	0.0	0.0
8390	764	3.4	37	2.3	2000	-9.9	0.0	0.0	0.0
8391	694	2.6	39	0.3	1968	-9.9	0.0	0.0	0.0
8392	712	2.1	41	1.0	1859	-9.8	0.0	0.0	0.0
8393	584	1.9	43	1.3	1909	-9.3	0.0	0.0	0.0

8393 rows × 9 columns

```
In [29]: df2 = df1[df1["wind_speed"] < 6]

In [30]: df2

Out[30]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
0	254	-5.2	37	2.2	2000	-17.6	0.0	0.0	0.0
1	204	-5.5	38	0.8	2000	-17.6	0.0	0.0	0.0
2	173	-6.0	39	1.0	2000	-17.7	0.0	0.0	0.0
3	107	-6.2	40	0.9	2000	-17.6	0.0	0.0	0.0
4	78	-6.0	36	2.3	2000	-18.6	0.0	0.0	0.0
...
8399	1003	4.2	34	2.6	1894	-10.3	0.0	0.0	0.0
8390	764	3.4	37	2.3	2000	-9.9	0.0	0.0	0.0
8391	694	2.6	39	0.3	1968	-9.9	0.0	0.0	0.0
8392	712	2.1	41	1.0	1859	-9.8	0.0	0.0	0.0
8393	584	1.9	43	1.3	1909	-9.3	0.0	0.0	0.0

8385 rows × 9 columns

```
In [31]: df2.reset_index(drop=True)

Out[31]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar_rad	rainfall	snowfall
0	254	-5.2	37	2.2	2000	-17.6	0.0	0.0	0.0
1	204	-5.5	38	0.8	2000	-17.6	0.0	0.0	0.0
2	173	-6.0	39	1.0	2000	-17.7	0.0	0.0	0.0
3	107	-6.2	40	0.9	2000	-17.6	0.0	0.0	0.0
4	78	-6.0	36	2.3	2000	-18.6	0.0	0.0	0.0
...
8380	1003	4.2	34	2.6	1894	-10.3	0.0	0.0	0.0
8381	764	3.4	37	2.3	2000	-9.9	0.0	0.0	0.0
8382	694	2.6	39	0.3	1968	-9.9	0.0	0.0	0.0
8383	712	2.1	41	1.0	1859	-9.8	0.0	0.0	0.0
8384	584	1.9	43	1.3	1909	-9.3	0.0	0.0	0.0

8385 rows × 9 columns

```
In [32]: #df2.to_csv("part2.csv", index=False)

In [33]: # Define a function that uses Matplotlib to visually compare the scale and distribution of bikes_rented and
# Call the function.

plt.figure(figsize=(20,5))
plt.subplot(1,2,1)
df2["bikes_rented"].plot(kind="hist")
plt.title("Bikes Rented")
plt.subplot(1,2,2)
df2["wind_speed"].plot(kind="hist")
plt.title("Wind Speed")
plt.show()

Out[33]:
```

Transform bikes_rented and wind_speed, and compare results

```
In [34]: df2 = pd.read_csv("part2.csv")

In [35]: # Apply a log transformation (np.log) to scale bikes_rented and wind speed.

# Compare scale and distribution of bikes_rented and wind speed by calling the function you defined earlier.
df2.head()

Out[35]:
```

	bikes_rented	temp	humidity	wind_speed	visibility	dew_temp	solar
--	--------------	------	----------	------------	------------	----------	-------