Data Wrangling, I Perform the following operations using Python on any open source dataset (e.g., data.csv)

1. Import all the required Python Libraries.
2. Locate an open source data from the web (e.g., https://www.kaggle.com). Provide a clear description of the data and its source (i.e., URL of the web site).
3. Load the Dataset into pandas dataframe.
4. Data Preprocessing: check for missing values in the data using pandas isnull(), describe() function to get some initial statistics. Provide variable descriptions. Types of variables etc. Check the dimensions of the data frame.
5. Data Formatting and Data Normalization: Summarize the types of variables by checking the data types (i.e., character, numeric, integer, factor, and logical) of the variables in the data set. If variables are not in the correct data type, apply proper type conversions.
6. Turn categorical variables into quantitative variables in Python.

```python
import pandas as pd
import numpy as np
```

```python
df = pd.read_csv("/content/Iris.csv")
```

```python
df
```

|     | Id  | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species        |
| --- | --- | ------------- | ------------ | ------------- | ------------ | -------------- |
| 0   | 1   | 5.1           | 3.5          | 1.4           | 0.2          | Iris-setosa    |
| 1   | 2   | 4.9           | 3.0          | 1.4           | 0.2          | Iris-setosa    |
| 2   | 3   | 4.7           | 3.2          | 1.3           | 0.2          | Iris-setosa    |
| 3   | 4   | 4.6           | 3.1          | 1.5           | 0.2          | Iris-setosa    |
| 4   | 5   | 5.0           | 3.6          | 1.4           | 0.2          | Iris-setosa    |
| ... | ... | ...           | ...          | ...           | ...          | ...            |
| 145 | 146 | 6.7           | 3.0          | 5.2           | 2.3          | Iris-virginica |
| 146 | 147 | 6.3           | 2.5          | 5.0           | 1.9          | Iris-virginica |
| 147 | 148 | 6.5           | 3.0          | 5.2           | 2.0          | Iris-virginica |
| 148 | 149 | 6.2           | 3.4          | 5.4           | 2.3          | Iris-virginica |
| 149 | 150 | 5.9           | 3.0          | 5.1           | 1.8          | Iris-virginica |

Saved successfully!

150 rows × 6 columns

```python
df['Species']
```

```
0        Iris-setosa
1        Iris-setosa
2        Iris-setosa
3        Iris-setosa
4        Iris-setosa
            ...
145    Iris-virginica
146    Iris-virginica
147    Iris-virginica
148    Iris-virginica
149    Iris-virginica
Name: Species, Length: 150, dtype: object
```

```python
df.iloc[1]
```

```
Id                       2
SepalLengthCm          4.9
SepalWidthCm           3.0
PetalLengthCm          1.4
PetalWidthCm           0.2
Species        Iris-setosa
Name: 1, dtype: object
```

```python
df["PetalLengthCm"].iloc[2]
```

```
1.3
```

```
idx = [1, 2, 3]
sample = df.iloc[idx]
sample
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| **1** | 2 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| **2** | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| **3** | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |

```
df.describe()
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|---|---|---|---|---|---|
| **count** | 150.000000 | 150.000000 | 150.000000 | 150.000000 | 150.000000 |
| **mean** | 75.500000 | 5.843333 | 3.054000 | 3.758667 | 1.198667 |
| **std** | 43.445368 | 0.828066 | 0.433594 | 1.764420 | 0.763161 |
| **min** | 1.000000 | 4.300000 | 2.000000 | 1.000000 | 0.100000 |
| **25%** | 38.250000 | 5.100000 | 2.800000 | 1.600000 | 0.300000 |
| **50%** | 75.500000 | 5.800000 | 3.000000 | 4.350000 | 1.300000 |
| **75%** | 112.750000 | 6.400000 | 3.300000 | 5.100000 | 1.800000 |
| **max** | 150.000000 | 7.900000 | 4.400000 | 6.900000 | 2.500000 |

```
df['PetalLengthCm'].mean()
```

```
3.758666666666666
```

Saved successfully!                          ✕

```
array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)
```

```
df.isnull()
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| **0** | False | False | False | False | False | False |
| **1** | False | False | False | False | False | False |
| **2** | False | False | False | False | False | False |
| **3** | False | False | False | False | False | False |
| **4** | False | False | False | False | False | False |
| **...** | ... | ... | ... | ... | ... | ... |
| **145** | False | False | False | False | False | False |
| **146** | False | False | False | False | False | False |
| **147** | False | False | False | False | False | False |
| **148** | False | False | False | False | False | False |
| **149** | False | False | False | False | False | False |

150 rows × 6 columns

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

```
Id               0
SepalLengthCm    0
SepalWidthCm     0
PetalLengthCm    0
PetalWidthCm     0
Species          0
dtype: int64
```

```
df.notnull()
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| 0 | True | True | True | True | True | True |
| 1 | True | True | True | True | True | True |
| 2 | True | True | True | True | True | True |
| 3 | True | True | True | True | True | True |
| 4 | True | True | True | True | True | True |
| ... | ... | ... | ... | ... | ... | ... |
| 145 | True | True | True | True | True | True |
| 146 | True | True | True | True | True | True |
| 147 | True | True | True | True | True | True |
| 148 | True | True | True | True | True | True |
| 149 | True | True | True | True | True | True |

150 rows × 6 columns

```
df.notnull().sum()
```

```
Id              150
SepalLengthCm   150
SepalWidthCm    150
PetalLengthCm   150
PetalWidthCm    150
Species         150
dtype: int64
```

```
df['Species'].replace({'Iris-setosa':1, 'Iris-versicolor':2, 'Iris-virginica':3}, inplace = True)
df
```

Saved successfully!                    ✕

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| 0 | 1 | 5.1 | 3.5 | 1.4 | 0.2 | 1 |
| 1 | 2 | 4.9 | 3.0 | 1.4 | 0.2 | 1 |
| 2 | 3 | 4.7 | 3.2 | 1.3 | 0.2 | 1 |
| 3 | 4 | 4.6 | 3.1 | 1.5 | 0.2 | 1 |
| 4 | 5 | 5.0 | 3.6 | 1.4 | 0.2 | 1 |
| ... | ... | ... | ... | ... | ... | ... |
| 145 | 146 | 6.7 | 3.0 | 5.2 | 2.3 | 3 |
| 146 | 147 | 6.3 | 2.5 | 5.0 | 1.9 | 3 |
| 147 | 148 | 6.5 | 3.0 | 5.2 | 2.0 | 3 |
| 148 | 149 | 6.2 | 3.4 | 5.4 | 2.3 | 3 |
| 149 | 150 | 5.9 | 3.0 | 5.1 | 1.8 | 3 |

150 rows × 6 columns

```
df.columns
```

```
Index(['Id', 'SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm',
       'Species'],
      dtype='object')
```

```
df.dtypes
```

```
Id              int64
SepalLengthCm   float64
SepalWidthCm    float64
PetalLengthCm   float64
PetalWidthCm    float64
Species         int64
dtype: object
```

Saved successfully!