

```
In [130]: # Importing Libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [131]: # Importing the Iris Dataset
iris = pd.read_csv("C:/Users/jasika/OneDrive/Documents/Datasets/Iris/Iris.csv")

In [132]: # Looking the dataset, its shape, datatype and descriptive statistics
iris.head()

Out[132]:
   sepal.length  petal.length  petal.width  variety
0          5.1          3.5          1.4      0.2  Setosa
1          4.9          3.0          1.4      0.2  Setosa
2          4.7          3.2          1.3      0.2  Setosa
3          4.6          3.1          1.5      0.2  Setosa
4          5.0          3.6          1.4      0.2  Setosa

In [133]: iris.shape
Out[133]: (150, 5)

In [134]: iris.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 150 entries, 0 to 149
Data columns (total 5 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   sepal.length  150 non-null     float64
 1   sepal.width   150 non-null     float64
 2   petal.length  150 non-null     float64
 3   petal.width   150 non-null     float64
 4   variety       150 non-null     object
dtypes: float64(4), object(1)
memory usage: 8.7+ KB

In [135]: iris.describe()

Out[135]:
   sepal.length  sepal.width  petal.length  petal.width
count  150.000000    150.000000    150.000000    150.000000
mean     5.843333     3.057333     3.758000     1.190333
std      0.828066     0.433066     0.762266     0.762236
min      4.300000     2.000000     1.000000     0.100000
25%      5.100000     2.800000     1.600000     0.300000
50%      5.800000     3.000000     3.300000     1.300000
75%      6.400000     3.300000     4.500000     1.800000
max      7.900000     4.600000     6.900000     2.000000

In [136]: # Renaming the Columns
df = iris.rename(columns = {'sepal.length': 'SepalLength',
                             'sepal.width': 'SepalWidth',
                             'petal.length': 'PetalLength',
                             'petal.width': 'PetalWidth'})

In [137]: # Data Cleaning - Looking at the null values and the duplicates
df.isnull().sum()

Out[137]:
SepalLength    0
SepalWidth     0
PetalLength    0
PetalWidth     0
variety        0
dtype: int64

In [138]: df.dropna(inplace=True)

Out[138]: 1

In [139]: df.drop_duplicates(inplace=True)

Out[139]:
   SepalLength  SepalWidth  PetalLength  PetalWidth  variety
0            5.1          3.5          1.4          0.2    Setosa
50           7.6          3.2          4.7          1.4  Versicolour
100          6.3          3.3          6.0          2.5  Virginica

In [140]: # Counting the distinct variety
df.value_counts('variety')

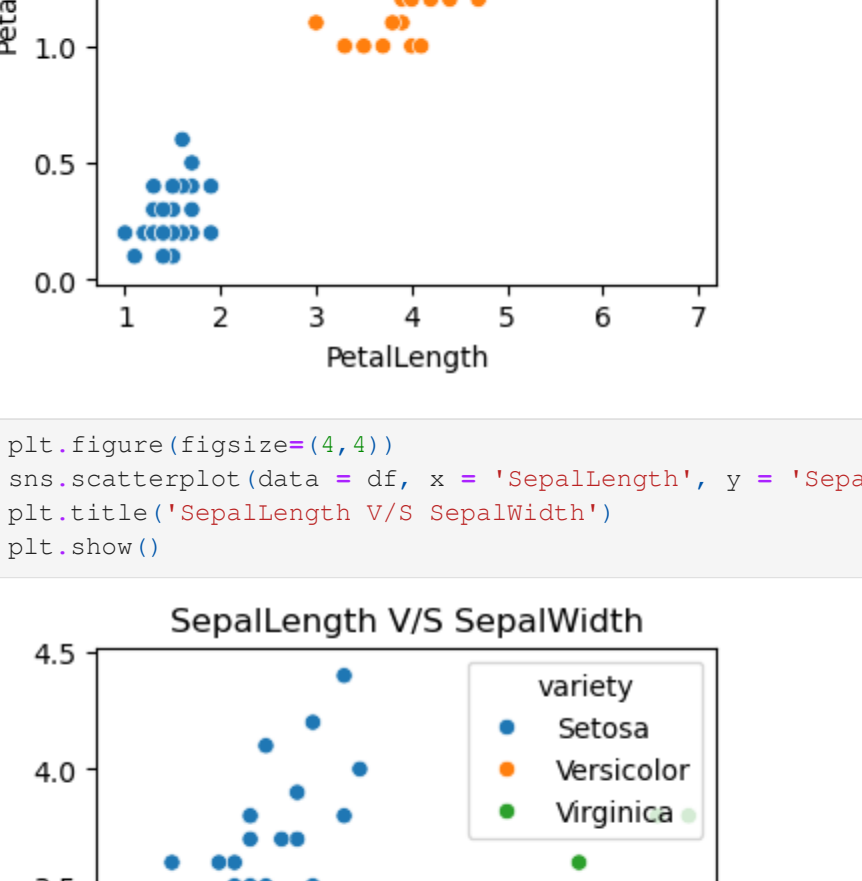
Out[140]:
variety
Setosa      50
Versicolour 50
Virginica   50
Name: count, dtype: int64

In [141]: # Graphics Representation
plt.figure(figsize=(5,3))
sns.countplot(data = df, x = 'variety', hue = 'variety')
plt.title('Bar chart Distribution')
plt.xlabel('variety')
plt.ylabel('No. of Iris')

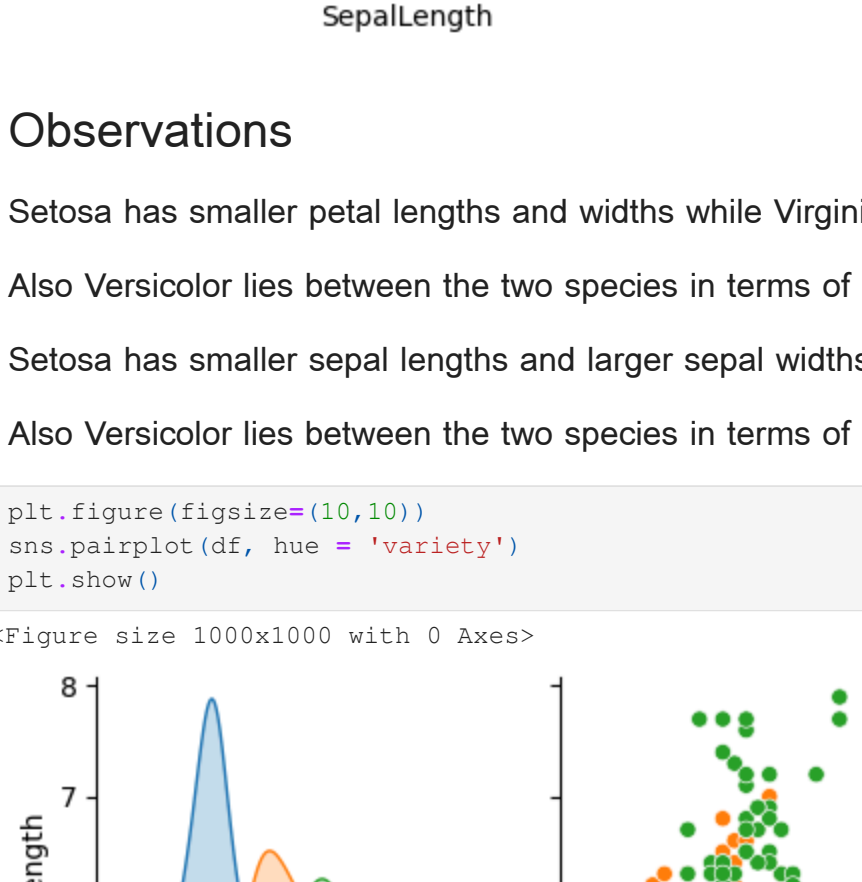
Out[141]:
Text(0, 0.5, 'No. of Iris')
```

We can see that there is equal numbers of each variety in our dataset.

```
In [142]: # Relationship between the PetalLength, PetalWidth, SepalLength and the SepalLength
plt.figure(figsize=(4,4))
sns.scatterplot(data = df, x = 'PetalLength', y = 'PetalWidth', hue = 'variety')
plt.title('PetalLength V/S PetalWidth')
plt.show()
```



```
In [143]: plt.figure(figsize=(4,4))
sns.scatterplot(data = df, x = 'SepalLength', y = 'SepalWidth', hue = 'variety')
plt.title('SepalLength V/S SepalWidth')
plt.show()
```



### Observations

- Setosa has smaller petal lengths and widths while Virginica has the largest of petal lengths and widths.
- Also Versicolour lies between the two species in terms of petal length and width.
- Setosa has smaller sepal lengths and larger sepal widths while Virginica has larger sepal lengths but smaller sepal widths.
- Also Versicolour lies between the two species in terms of sepal length and width.s.s.

```
In [144]: plt.figure(figsize=(10,10))
sns.pairplot(df, hue = 'variety')
plt.show()

#Figure size 100x100 with 0 axis
```



```
In [145]: fig, axes = plt.subplots(2, 2, figsize=(8,8))

axes[0,0].set_title('Sepal Length')
axes[0,0].hist(df['SepalLength'], bins=5)

axes[0,1].set_title('Sepal Width')
axes[0,1].hist(df['SepalWidth'], bins=5)

axes[1,0].set_title('Petal Length')
axes[1,0].hist(df['PetalLength'], bins=5)

axes[1,1].set_title('Petal Width')
axes[1,1].hist(df['PetalWidth'], bins=5)

Out[145]:
(array([16., 2., 15., 37., 15., 23.]),
 array([0.1, 0.5, 0.5, 1.5, 1.5, 1.7, 2.1, 2.5]),
 OutContourSet-0x2c4572750bb67751)
```

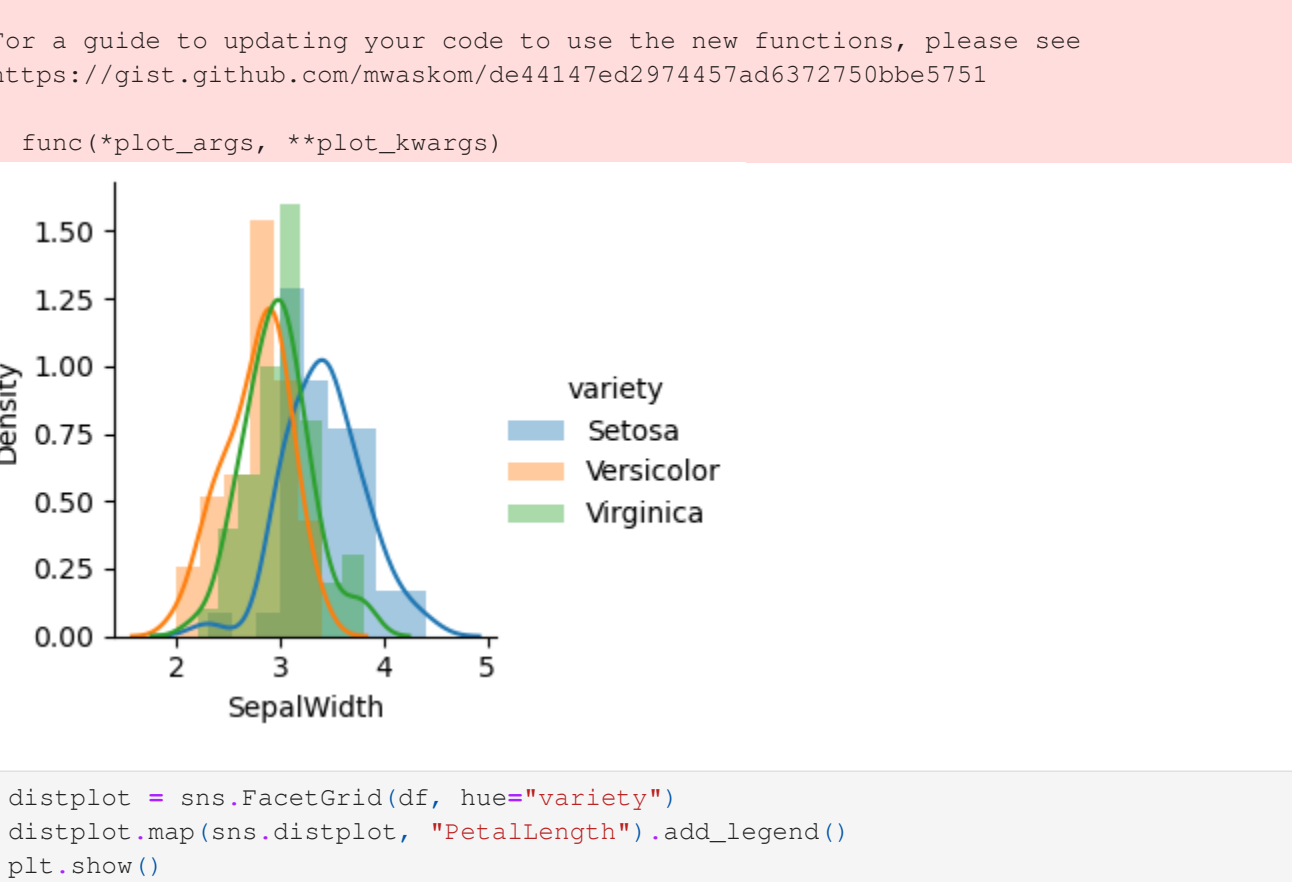


- The highest frequency of the sepal length is between 30 and 35 which is between 5.5 and 6
- The highest frequency of the sepal width is around 70 which is between 3.0 and 3.5
- The highest frequency of the petal length is around 50 which is between 1 and 2
- The highest frequency of the petal width is between 40 and 50 which is between 0.0 and 0.5

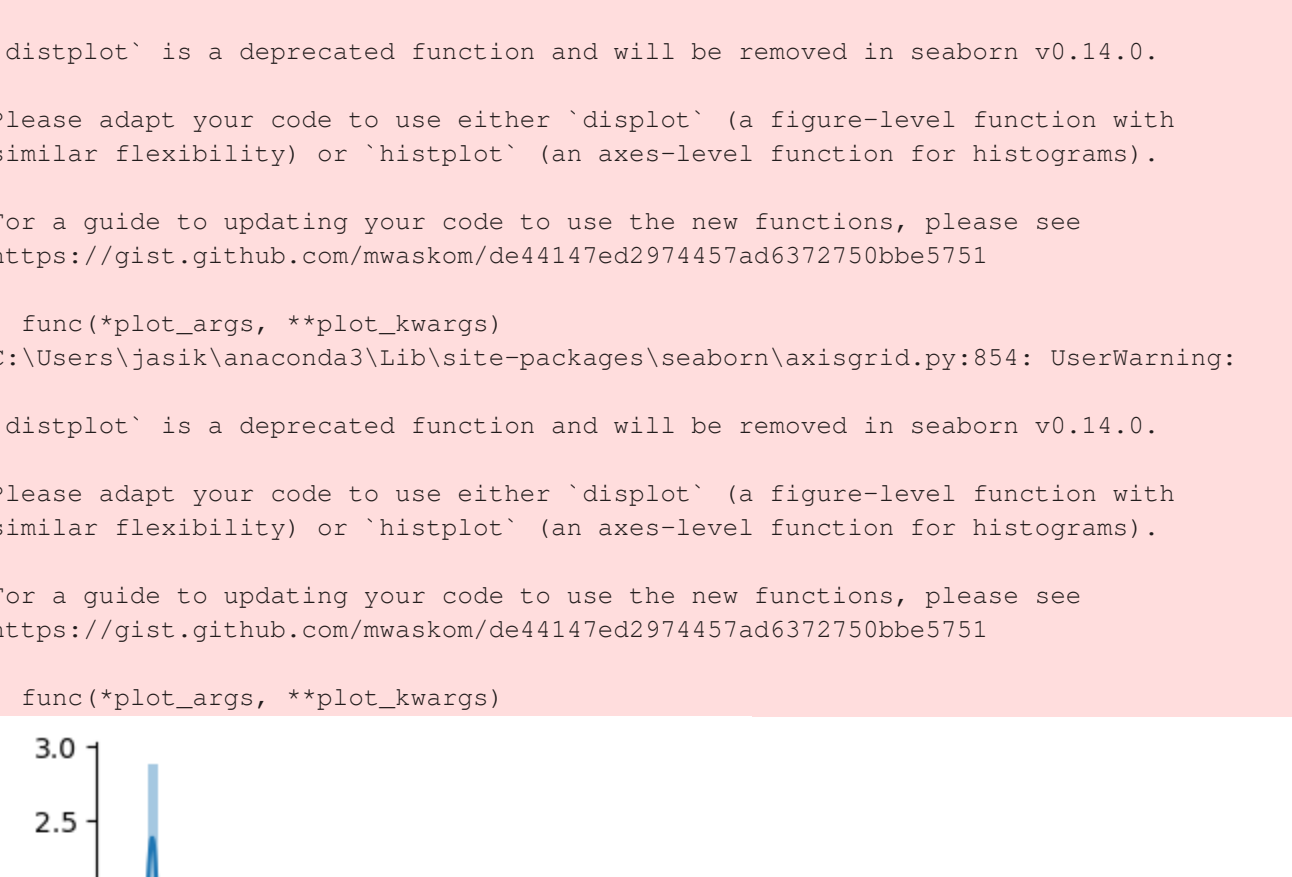
```
In [146]: distplot = sns.FacetGrid(df, hue='variety')
distplot.map(sns.distplot, "PetalLength", add_legend=True)
plt.show()
```



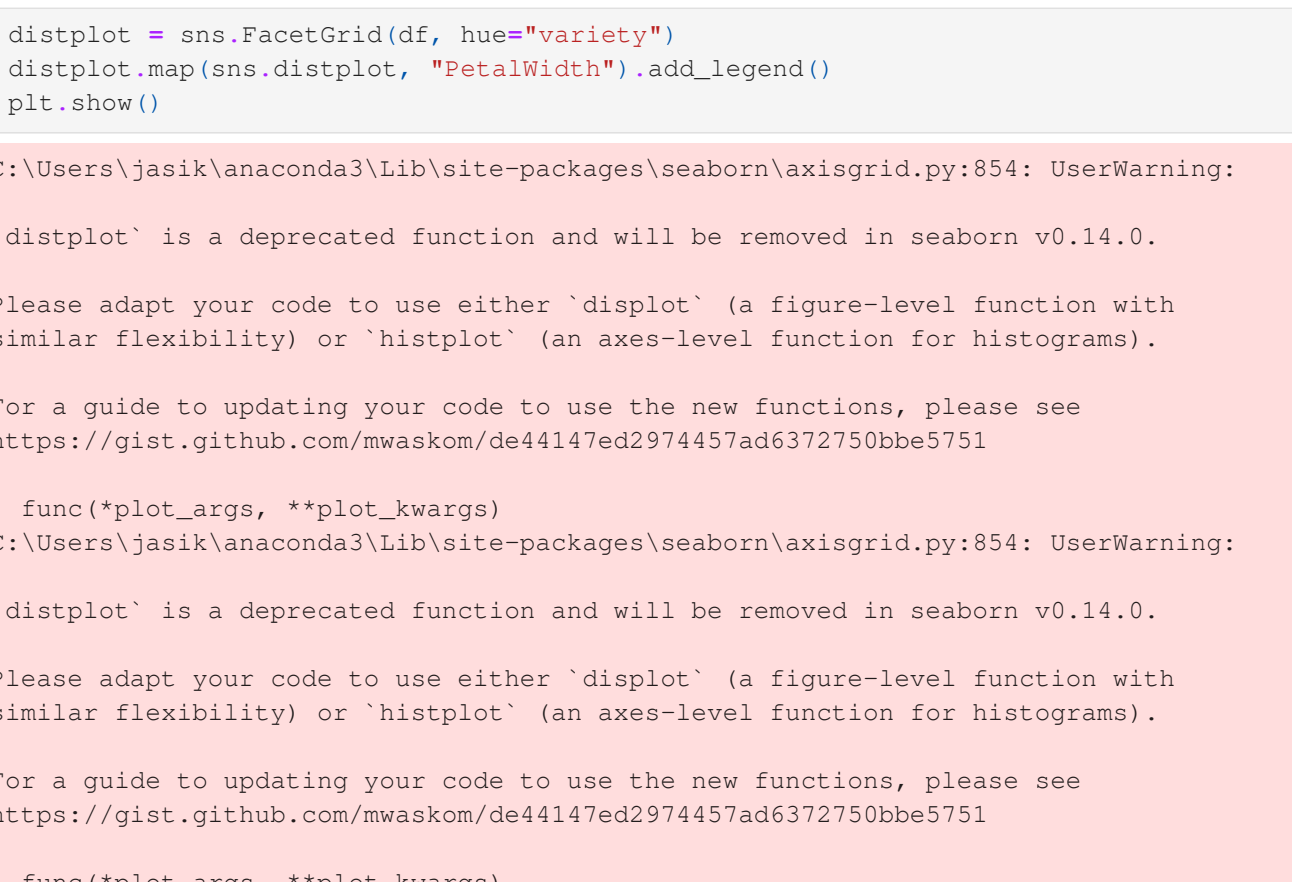
```
In [147]: distplot = sns.FacetGrid(df, hue='variety')
distplot.map(sns.distplot, "PetalWidth", add_legend=True)
plt.show()
```



```
In [148]: distplot = sns.FacetGrid(df, hue='variety')
distplot.map(sns.distplot, "SepalLength", add_legend=True)
plt.show()
```



```
In [149]: distplot = sns.FacetGrid(df, hue='variety')
distplot.map(sns.distplot, "SepalWidth", add_legend=True)
plt.show()
```



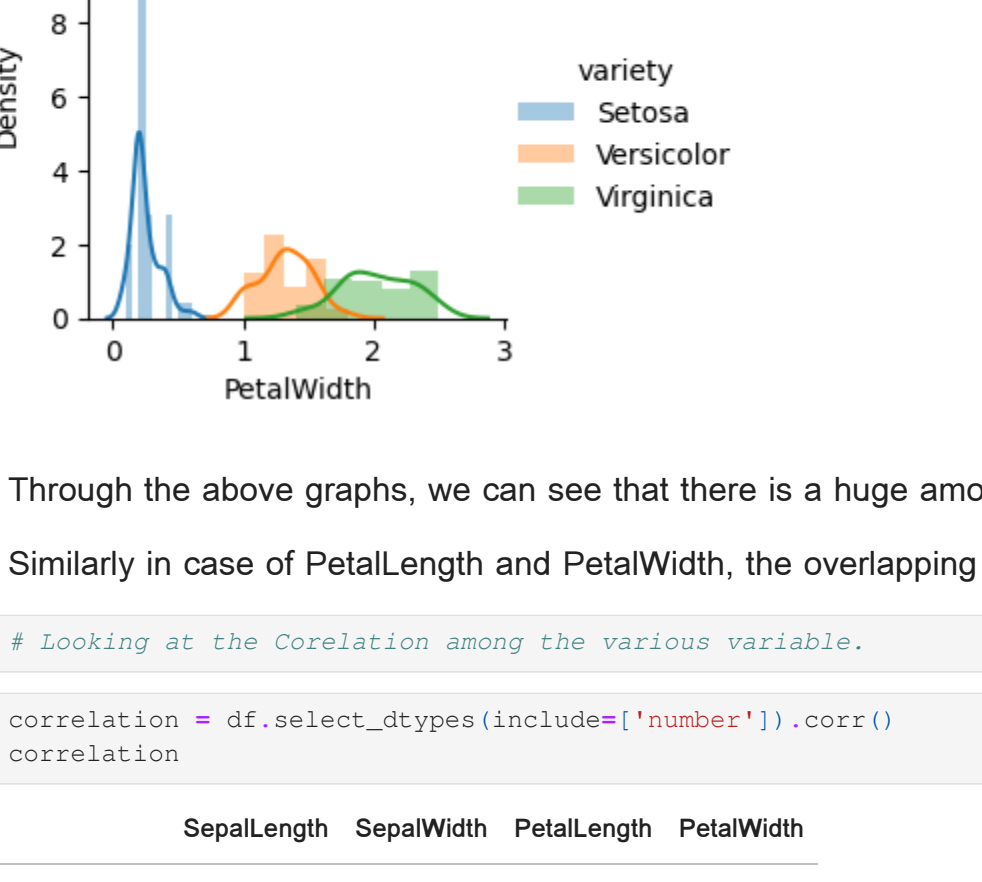
Through the above graphs, we can see that there is a huge amount of overlapping in case of SepalLength and SepalWidth. Similarly in case of PetalLength and PetalWidth, the overlapping reduces.

```
In [150]: # Looking at the Correlation among the various variable.
correlation = df.select_dtypes(include='number').corr()
correlation

Out[150]:
   SepalLength  SepalWidth  PetalLength  PetalWidth
SepalLength    1.000000    -0.117570    0.871754    0.817241
SepalWidth     -0.117570    1.000000   -0.428440   -0.368126
PetalLength     0.871754   -0.428440    1.000000    0.902865
PetalWidth      0.817241   -0.368126    0.902865    1.000000

In [151]: sns.heatmap(correlation, annot = True)

Out[151]:
<axes> >
```



Observation

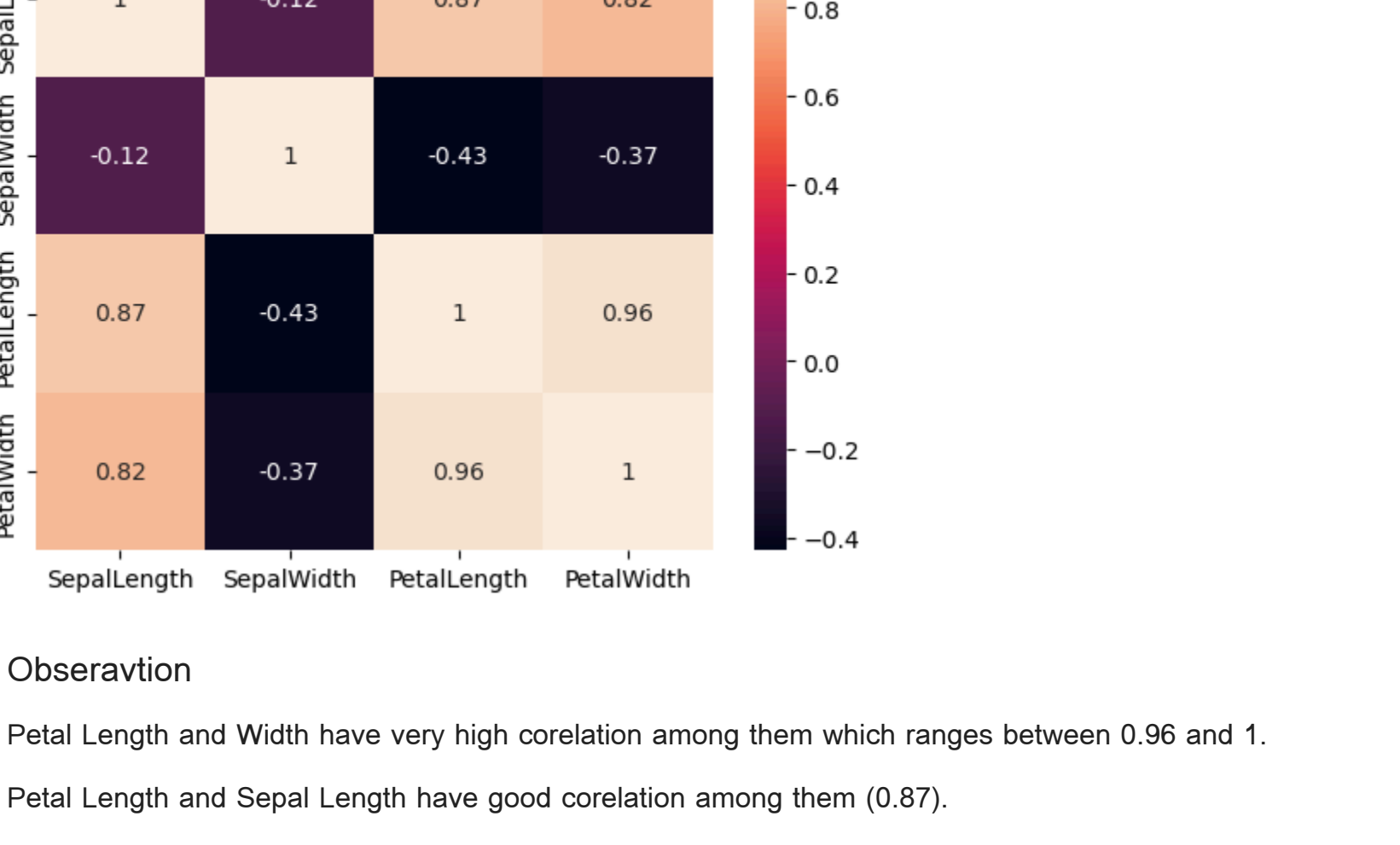
Petal Length and Width have very high correlation among them which ranges between 0.96 and 1.

Petal Length and Sepal Length have good correlation among them (0.87).

Petal Width and Sepal Length have good correlation among them (0.82).

```
In [152]: # Skplot to look at the distribution
df.groupby('variety').boxplot()

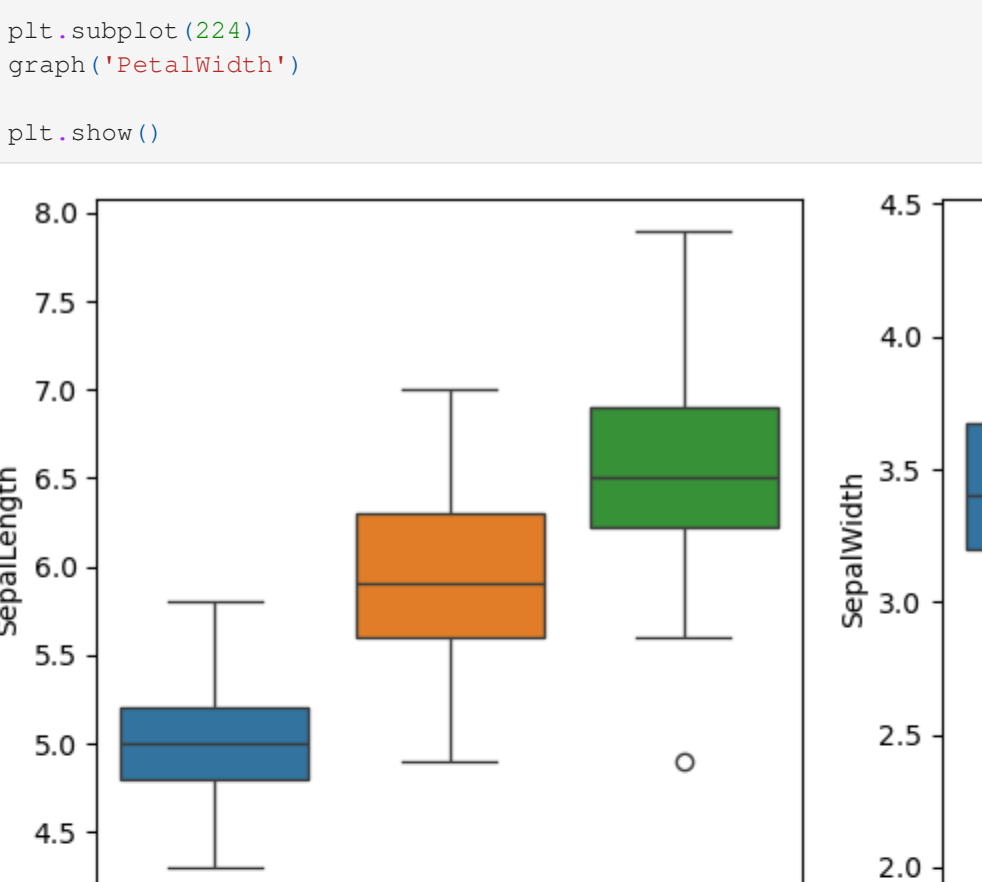
In [153]: plt.figure(figsize = (10, 10))
plt.subplot(221)
plt.title('SepalLength')
plt.subplot(222)
plt.title('SepalWidth')
plt.subplot(223)
plt.title('PetalLength')
plt.subplot(224)
plt.title('PetalWidth')
plt.show()
```



Species Setosa has the smallest features and less distributed with some outliers whereas Virginica has the highest features.

```
In [154]: # Looking for outliers and handling them
sns.boxplot(x='SepalWidth', data=df)

Out[154]:
<axes: xlabel='SepalWidth'>
```



```
In [155]: Q1 = df['SepalWidth'].quantile(0.25)
Q3 = df['SepalWidth'].quantile(0.75)
IQR = Q3 - Q1
IQR

Out[155]: 0.5

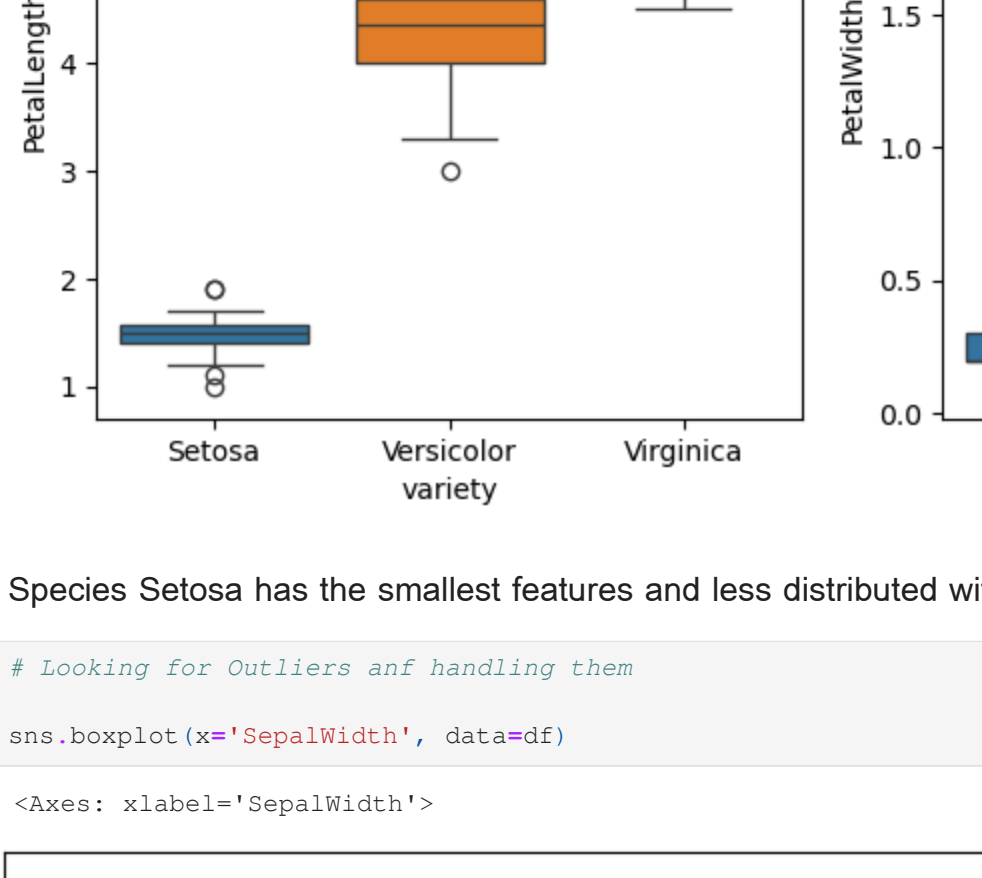
In [156]: print('Old Shape: ', df.shape)
upper = np.where(df['SepalWidth'] >= (Q3+1.5*IQR))
lower = np.where(df['SepalWidth'] <= (Q1-1.5*IQR))

# Removing the Outliers
df.drop_upper(0).inplace = True
df.drop_lower(0).inplace = True

print('New Shape: ', df.shape)

sns.boxplot(x='SepalWidth', data=df)

Out[156]:
Old Shape: (150, 5)
New Shape: (144, 5)
<axes: xlabel='SepalWidth'>
```



In the previous graph we could see, values less than 2 and values more than 4 are acting as Outliers.

Removing them, brings mean closer to the central Values.



