

Dataset Source-Kaggle.com Dataset:Iris dataset is a very know dataset used is a good example of classification problem. There are 50 sample on each species. 150 samples. 3 different Species(Iris-setosa, Iris-versicolor, Iris-virginica). Features are SepalLengthCm, SepalWidthCm,PetalLengthCm, Species name. We will not be using the Id, we will drop them.

```
In [31]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
import scipy.stats as stats
import warnings
warnings.filterwarnings("ignore")

In [32]: print(pd.__version__, np.__version__)

2.2.3 1.26.4
```

Loading the Dataset

```
In [33]: df= pd.read_csv("iris.csv")

In [34]: #Display the first few rows
df.head()
```

| | 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 |

```
In [35]: print(df.columns)

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

```
In [36]: #checking dataset shape(rows and columns)
df.shape
```

```
Out[36]: (150, 6)
```

```
In [37]: df.describe(include="all")
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|--------|------------|---------------|--------------|---------------|--------------|-------------|
| count | 150.000000 | 150.000000 | 150.000000 | 150.000000 | 150.000000 | 150 |
| unique | NaN | NaN | NaN | NaN | NaN | 3 |
| top | NaN | NaN | NaN | NaN | NaN | Iris-setosa |
| freq | NaN | NaN | NaN | NaN | NaN | 50 |
| mean | 75.500000 | 5.843333 | 3.054000 | 3.758667 | 1.198667 | NaN |
| std | 43.445368 | 0.828066 | 0.433594 | 1.764420 | 0.763161 | NaN |
| min | 1.000000 | 4.300000 | 2.000000 | 1.000000 | 0.100000 | NaN |
| 25% | 38.250000 | 5.100000 | 2.800000 | 1.600000 | 0.300000 | NaN |
| 50% | 75.500000 | 5.800000 | 3.000000 | 4.350000 | 1.300000 | NaN |
| 75% | 112.750000 | 6.400000 | 3.300000 | 5.100000 | 1.800000 | NaN |
| max | 150.000000 | 7.900000 | 4.400000 | 6.900000 | 2.500000 | NaN |

```
In [38]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  --
0    Id          150 non-null    int64
1    SepalLengthCm  150 non-null    float64
2    SepalWidthCm   150 non-null    float64
3    PetalLengthCm  150 non-null    float64
4    PetalWidthCm   150 non-null    float64
5    Species       150 non-null    object
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB
```

```
In [39]: df.isna().sum()

Out[39]: Id          0
SepalLengthCm      0
SepalWidthCm       0
PetalLengthCm      0
PetalWidthCm       0
Species            0
dtype: int64

No null values in the data.
```

```
In [40]: df.drop('Id', axis=1, inplace=True)
```

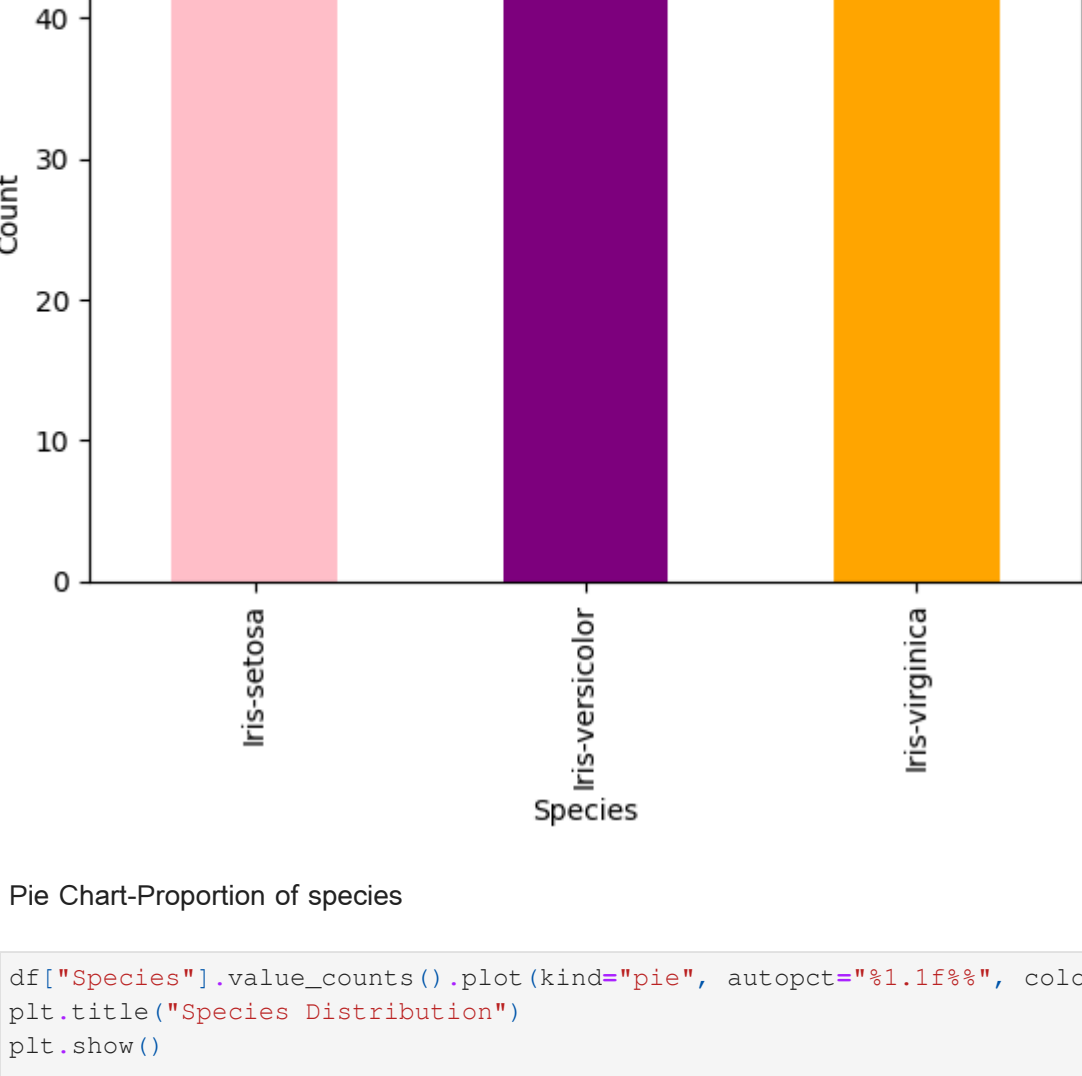
Dropped the Id column

```
In [41]: #checking after dropping the id column
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  --
0    SepalLengthCm  150 non-null    float64
1    SepalWidthCm   150 non-null    float64
2    PetalLengthCm  150 non-null    float64
3    PetalWidthCm   150 non-null    float64
4    Species       150 non-null    object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```

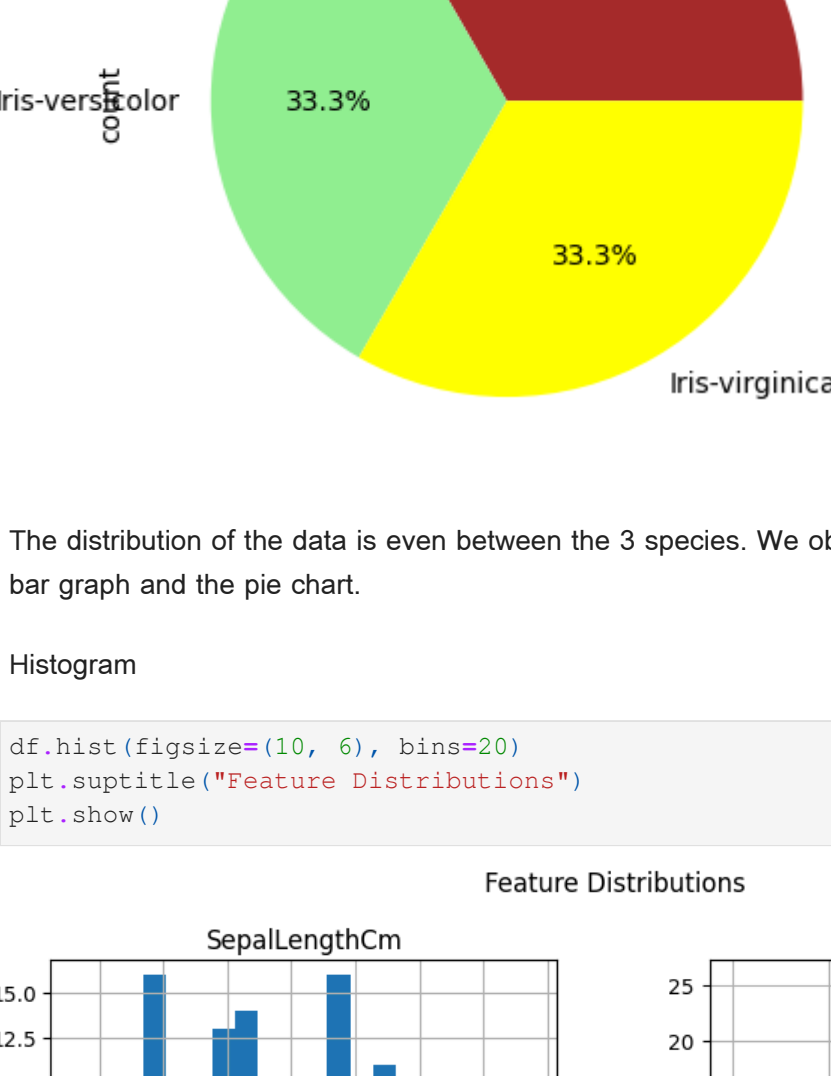
Bar Graph-count of each species

```
In [42]: df["Species"].value_counts().plot(kind="bar", color=["pink","purple","orange"])
plt.title("Count of Each Iris Species")
plt.xlabel("Species")
plt.ylabel("Count")
plt.show()
```



Pie Chart-Proportion of species

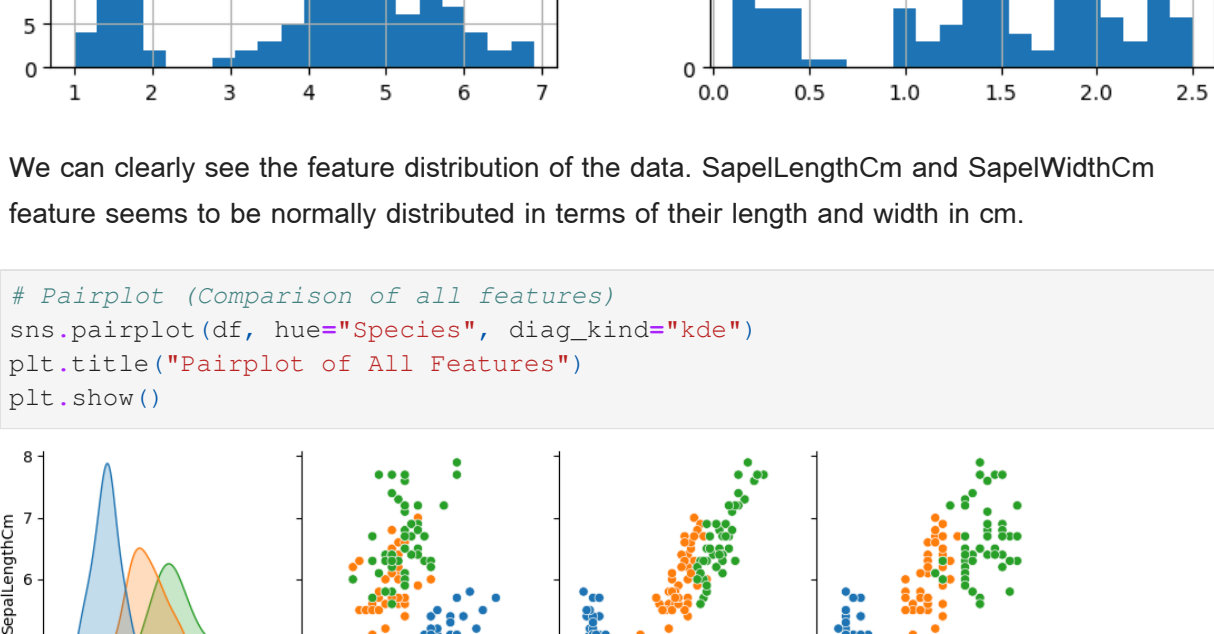
```
In [43]: df["Species"].value_counts().plot(kind="pie", autopct="%1.1f%%", colors=["brown","green","yellow"])
plt.title("Species Distribution")
plt.show()
```



The distribution of the data is even between the 3 species. We observe that from the above bar graph and the pie chart.

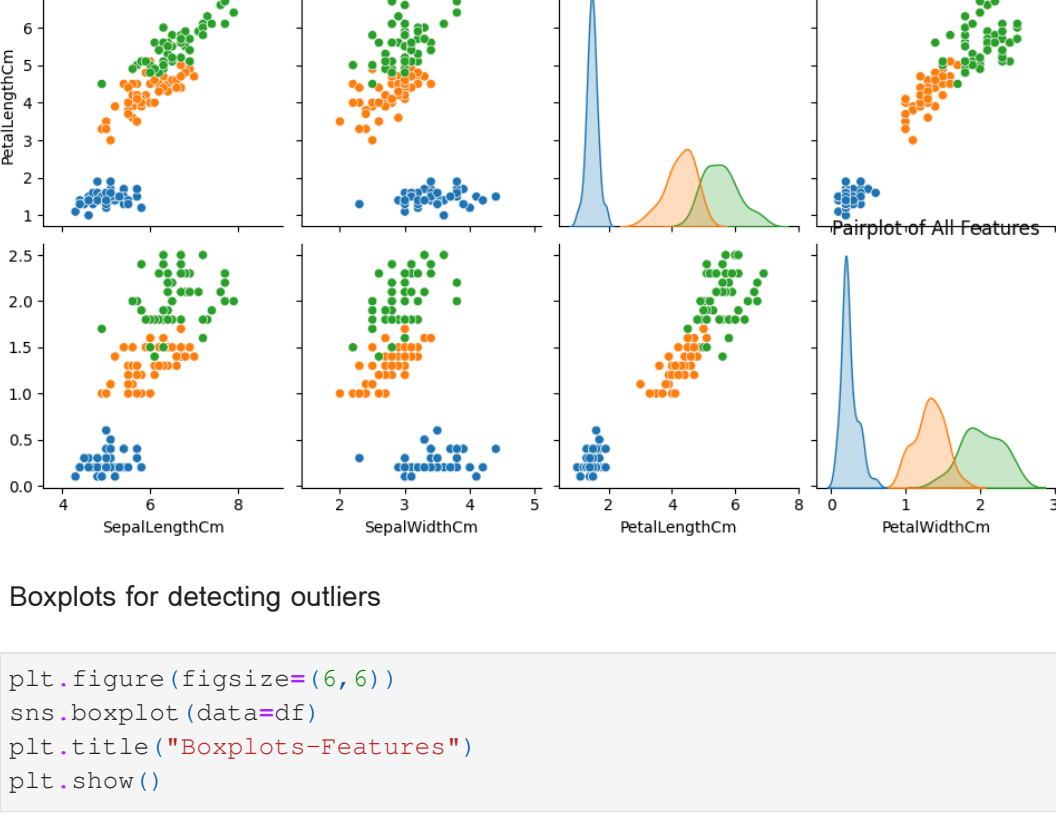
Histogram

```
In [44]: df.hist(figsize=(10, 6), bins=20)
plt.suptitle("Feature Distributions")
plt.show()
```



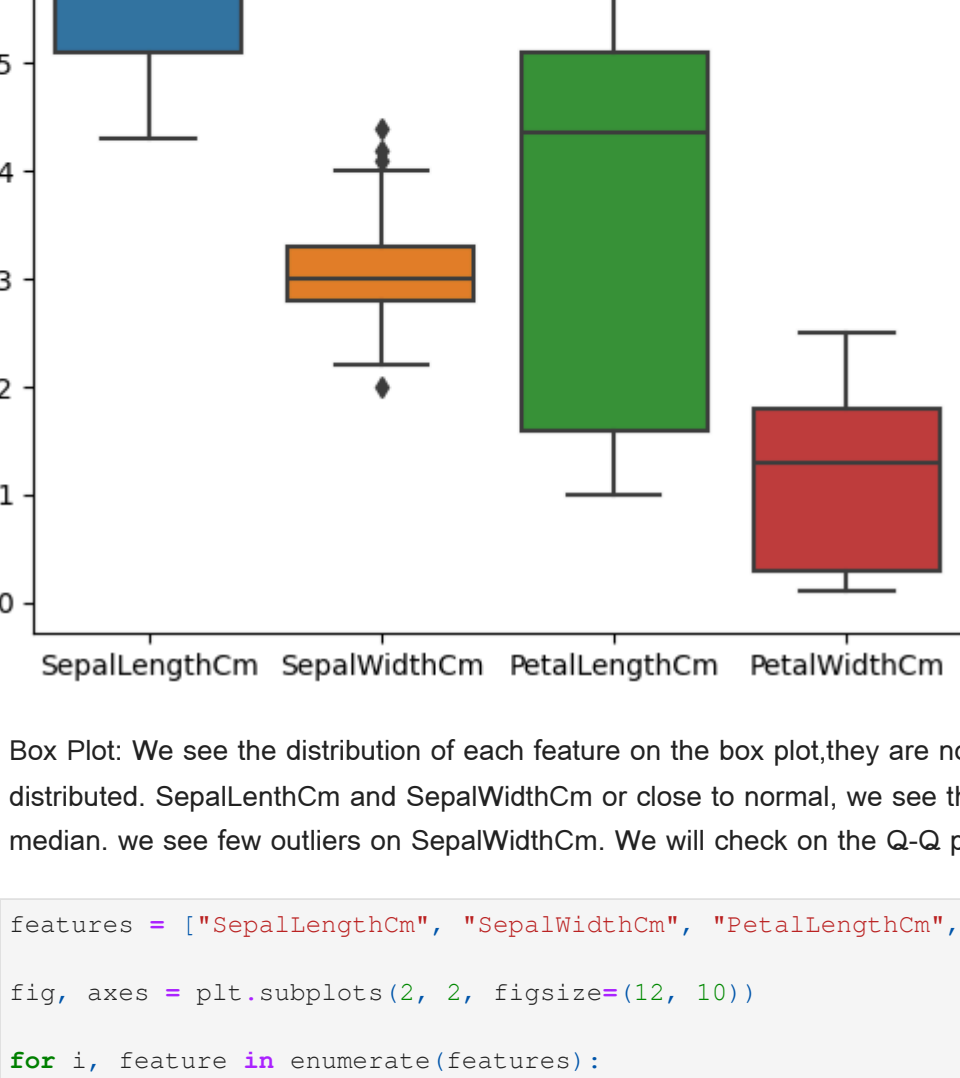
We can clearly see the feature distribution of the data. SepalLengthCm and SepalWidthCm feature seems to be normally distributed in terms of their length and width in cm.

```
In [46]: # Pairplot (Comparison of all features)
sns.pairplot(df, hue="Species", diag_kind="kde")
plt.title("Pair plot of All-Features")
plt.show()
```



Boxplots for detecting outliers

```
In [18]: plt.figure(figsize=(6,6))
sns.boxplot(data=df)
plt.title("Boxplots=Features")
plt.show()
```

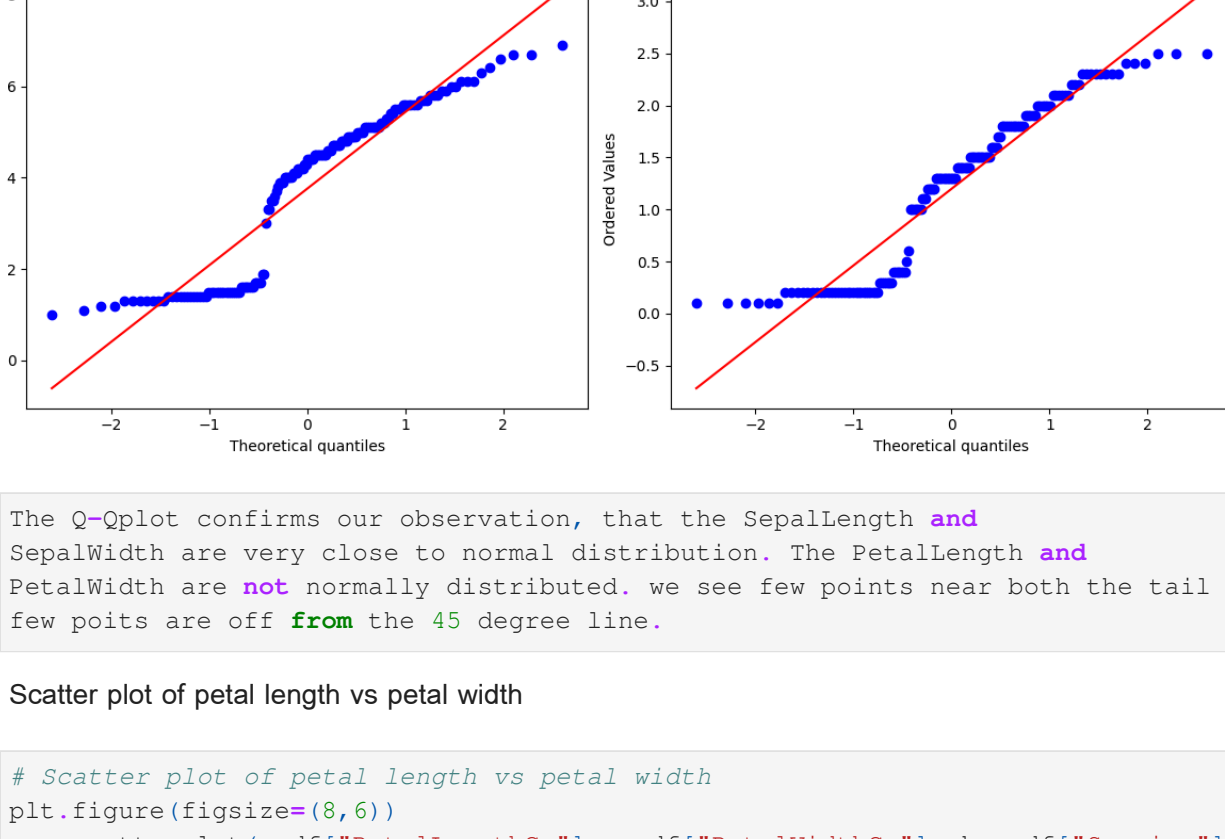


Box Plot: We see the distribution of each feature on the box plot, they are not normally distributed. SepalLengthCm and SepalWidthCm are close to normal, we see that from the median. We see few outliers on SepalWidthCm. We will check on the Q-Q plot to confirm this.

```
In [22]: features = ["SepalLengthCm", "SepalWidthCm", "PetalLengthCm", "PetalWidthCm"]
fig, axes = plt.subplots(2, 2, figsize=(12, 10))

for i, feature in enumerate(features):
    stats.probplot(df[feature], dist='norm', plot=axes[i//2, i%2])
    axes[i//2, i%2].set_title(f"Q-Q Plot for {feature}")

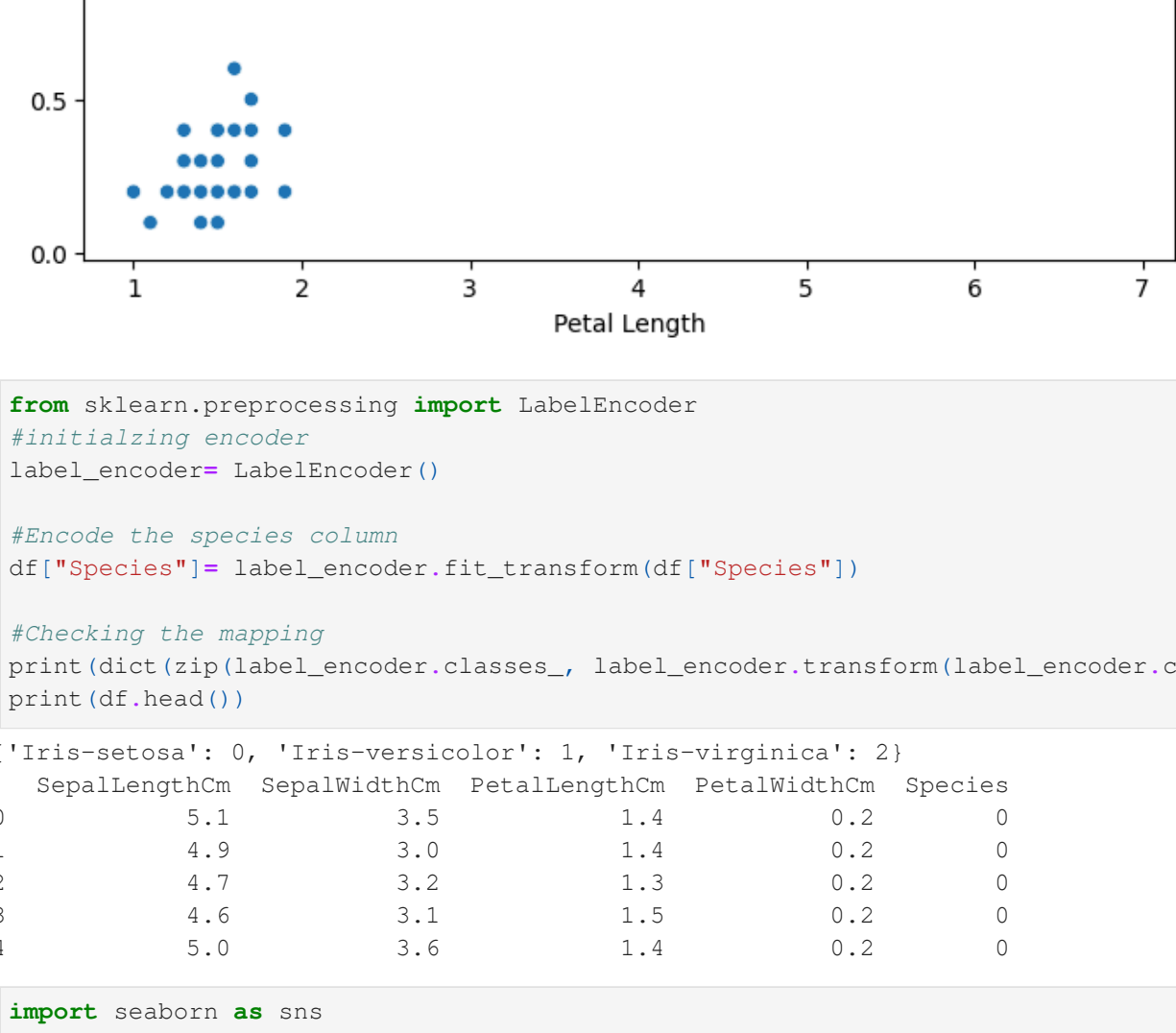
plt.tight_layout()
plt.show()
```



The Q-Q plot confirms our observation, that the SepalLength and SepalWidth are very close to normal distribution. The PetalLength and PetalWidth are not normally distributed. We see few points near both the tail and the head of the distribution. We will check on the Q-Q plot to confirm this.

Scatter plot of petal length vs petal width

```
In [56]: # Scatter plot of petal length vs petal width
plt.figure(figsize=(8,6))
sns.scatterplot(x=df["PetalLengthCm"], y=df["PetalWidthCm"], hue=df["Species"])
plt.title("Petal Length vs Petal Width")
plt.xlabel("Petal Length")
plt.ylabel("Petal Width")
plt.show()
```



```
In [24]: from sklearn.preprocessing import LabelEncoder
#initializing encoder
label_encoder=LabelEncoder()

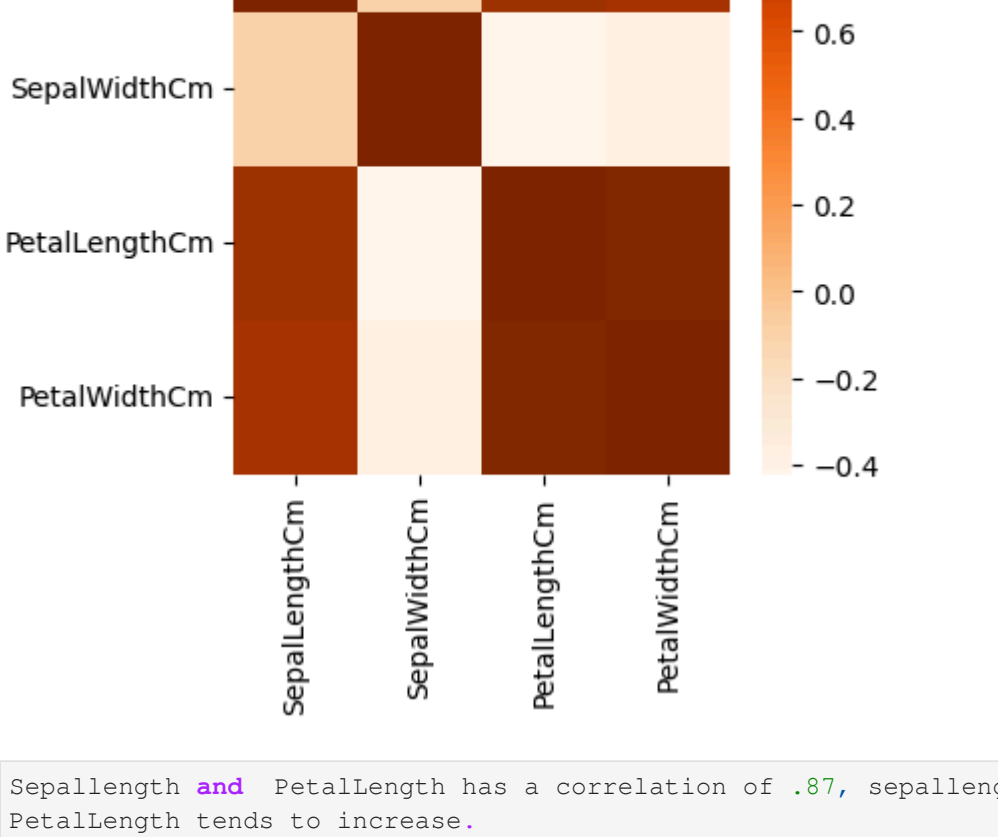
#Encode the species column
df["Species"]=label_encoder.fit_transform(df["Species"])

#Checking the mapping
print(dict(zip(label_encoder.classes_, label_encoder.transform(label_encoder.classes_))))
print(df.head())
```

```
In [50]: import seaborn as sns
import matplotlib.pyplot as plt

# Ensure only numeric data is used
numeric_df = df.select_dtypes(include=["float64", "int64"])

# Compute correlation matrix
corr_matrix = numeric_df.corr()
```



SepalLength and PetalLength has a correlation of .87, sepallength increases PetalLength tends to increase.

SepalLength and petalWidth has next high correlation of 0.82. Longer sepals are associated with petals.

SepalWidth and SepalWidth has a correlation of -0.11, which is a slight negative correlation. Wider sepals might correspond to shorter lengths in some cases.