

```
In [99]: # DSBDA Practical A-10:
# Download the Iris flower dataset or any other dataset into a DataFrame.
# Scan the dataset and give the inference as:
# 1. List down the features and their types (e.g., numeric, nominal) available in the dataset.
# 2. Create a histogram for each feature in the dataset to illustrate the feature distributions.
# 3. Create a boxplot for each feature in the dataset.
# 4. Compare distributions and identify outliers.
```

```
In [100]: # https://www.kaggle.com/datasets/uciml/iris?resource=download
```

```
In [101]: df=pd.read_csv('Iris.csv')
```

```
In [102]: df
```

Out[102]:

	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

150 rows × 6 columns

```
In [103]: 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 [104]: df.dtypes
```

```
Out[104]: Id              int64
SepalLengthCm          float64
SepalWidthCm           float64
PetalLengthCm          float64
PetalWidthCm           float64
Species                object
dtype: object
```

```
In [105]: df.describe()
```

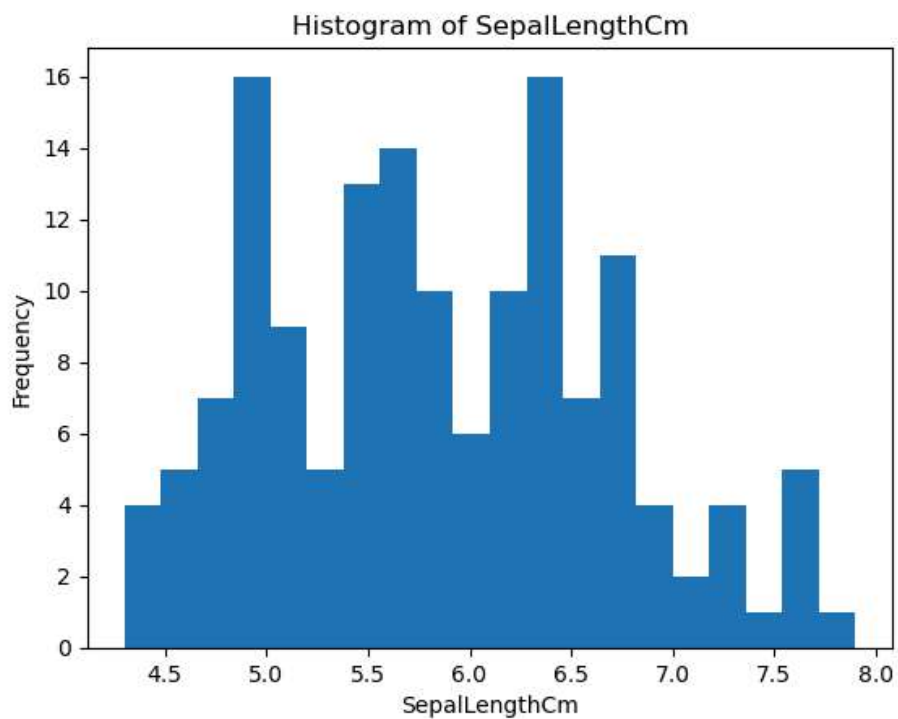
```
Out[105]:
```

	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

```
In [106]: # Draw Histograms for each feature in dataset to illustrate the feature distributions
import matplotlib.pyplot as plt
```

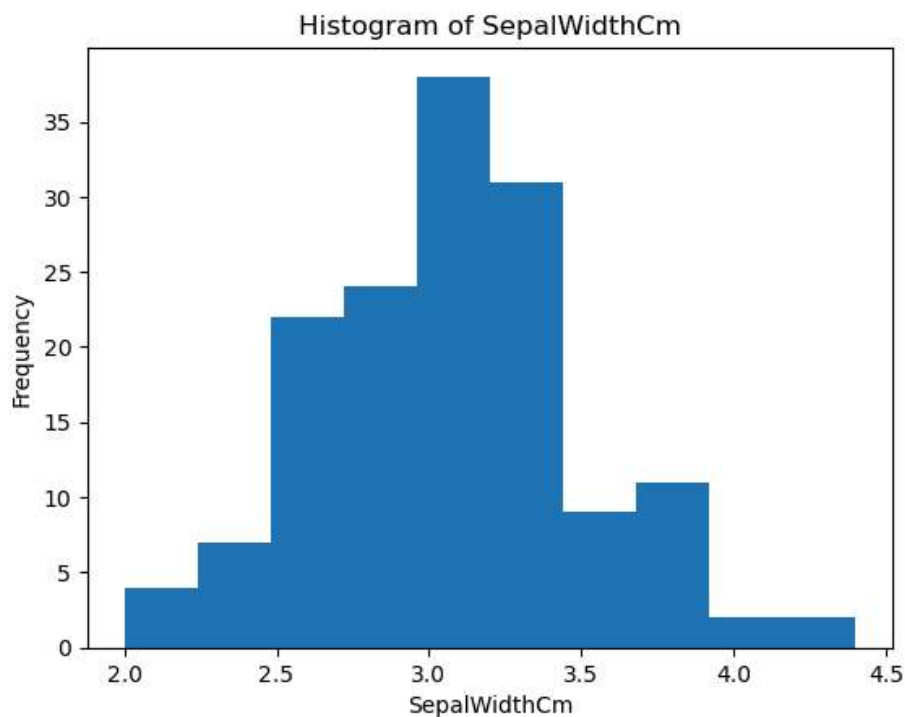
```
In [107]: plt.hist(df['SepalLengthCm'], bins=20) # Adjust the number of bins as needed
plt.title('Histogram of SepalLengthCm')
plt.xlabel('SepalLengthCm')
plt.ylabel('Frequency')
```

```
Out[107]: Text(0, 0.5, 'Frequency')
```



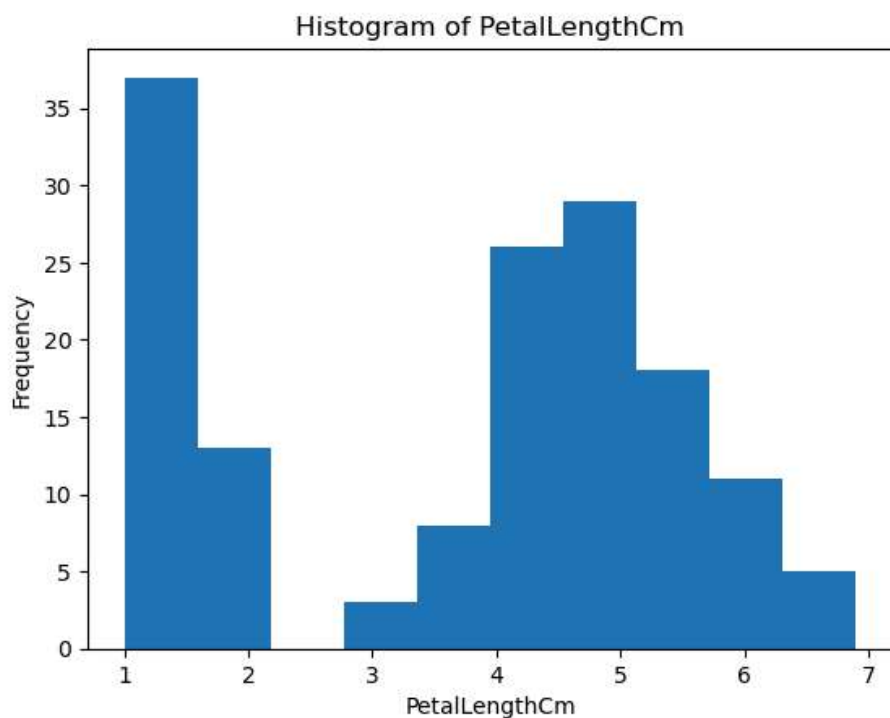
```
In [108]: plt.hist(df['SepalWidthCm'], bins=10) # Adjust the number of bins as needed
plt.title(f'Histogram of SepalWidthCm')
plt.xlabel('SepalWidthCm')
plt.ylabel('Frequency')
```

Out[108]: Text(0, 0.5, 'Frequency')



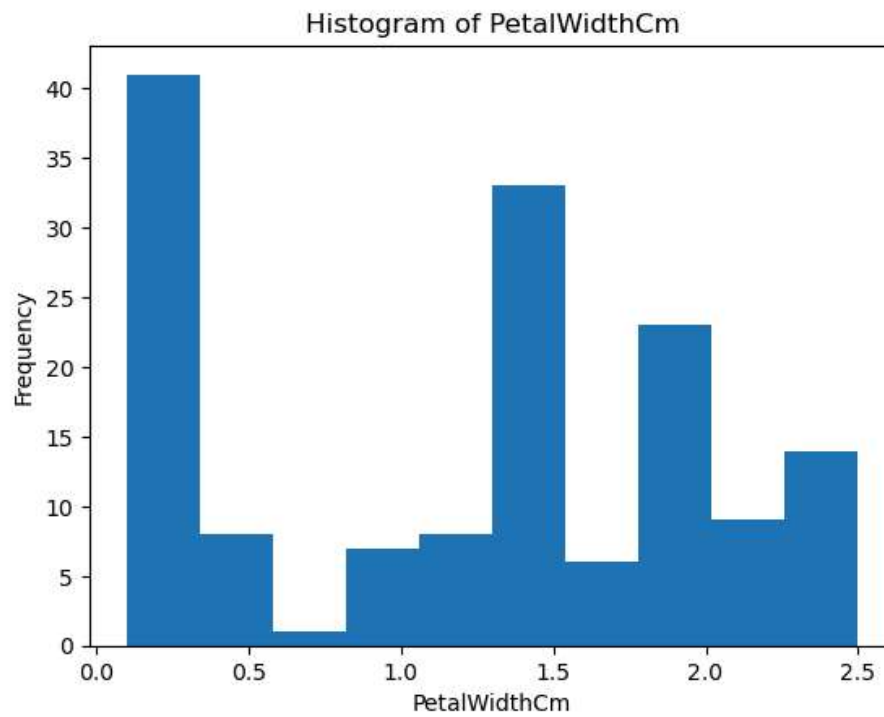
```
In [109]: plt.hist(df['PetalLengthCm'], bins=10) # Adjust the number of bins as needed
plt.title(f'Histogram of PetalLengthCm')
plt.xlabel('PetalLengthCm')
plt.ylabel('Frequency')
```

Out[109]: Text(0, 0.5, 'Frequency')

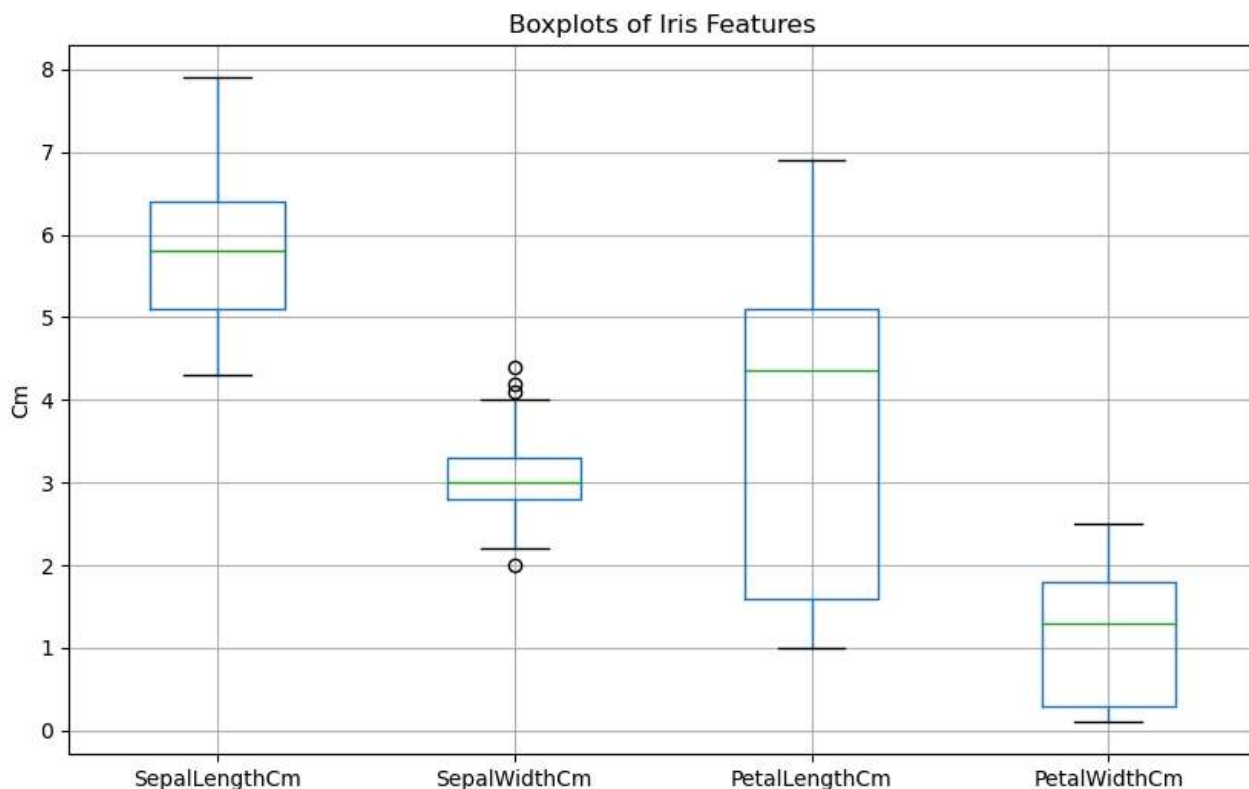


```
In [110]: plt.hist(df['PetalWidthCm'], bins=10) # Adjust the number of bins as needed
plt.title('Histogram of PetalWidthCm')
plt.xlabel('PetalWidthCm')
plt.ylabel('Frequency')
```

Out[110]: Text(0, 0.5, 'Frequency')



```
In [111]: # Create boxplots for SepalLengthCm, SepalWidthCm, PetalLengthCm, PetalWidthCm
plt.figure(figsize=(10, 6))
df.boxplot(column=['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm'])
plt.title('Boxplots of Iris Features')
plt.ylabel('Cm')
plt.show()
```



```
In [112]: # Detect outlier in Boxplot of SepalWidthCm Column
# The boxplot itself visually represents outliers.
# We can use IQR to programmatically find them.
```

```
In [113]: # Calculate Q1, Q3, and IQR
Q1 = df['SepalWidthCm'].quantile(0.25)
Q3 = df['SepalWidthCm'].quantile(0.75)
IQR = Q3 - Q1
```

```
In [114]: # Define bounds for outliers
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
```

```
In [115]: lower_bound
```

```
Out[115]: 2.05
```

```
In [116]: upper_bound
```

```
Out[116]: 4.05
```

```
In [117]: # Identify outliers
outliers = df[(df['SepalWidthCm'] < lower_bound) | (df['SepalWidthCm'] > upper_bound)]
```

```
In [118]: # Print or further process the outliers
print("Outliers in SepalWidthCm:")
outliers
```

Outliers in SepalWidthCm:

Out[118]:

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
15	16	5.7	4.4	1.5	0.4	Iris-setosa
32	33	5.2	4.1	1.5	0.1	Iris-setosa
33	34	5.5	4.2	1.4	0.2	Iris-setosa
60	61	5.0	2.0	3.5	1.0	Iris-versicolor

```
In [119]: # Printing the Row Indexes of Outliers
outlier_indices = outliers.index
print("Outliers RowIndex:")
outlier_indices
```

Outliers RowIndex:

Out[119]: Int64Index([15, 32, 33, 60], dtype='int64')

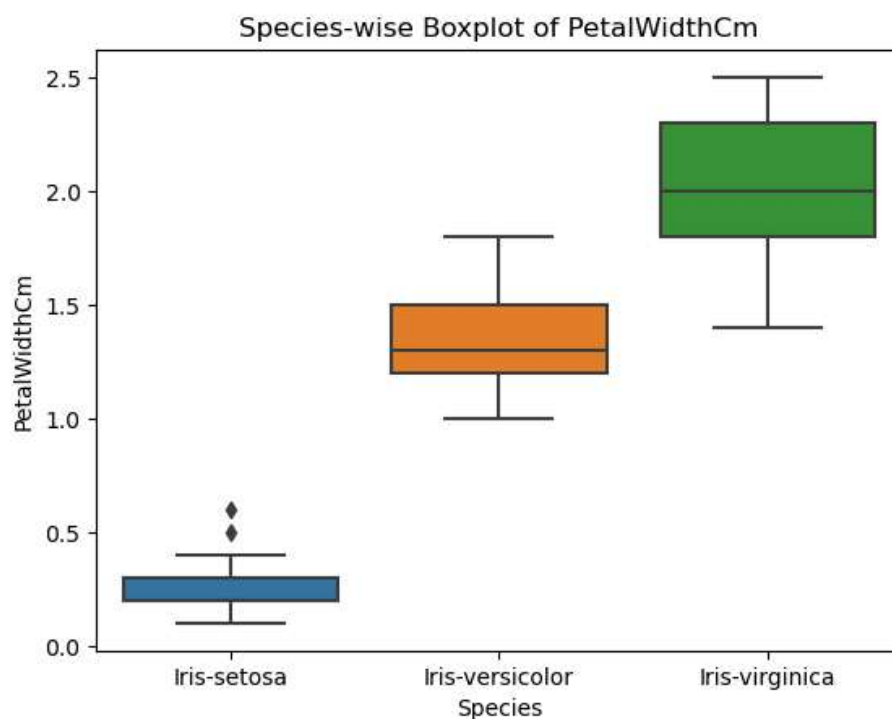
In [120]:

*#OPTIONAL*

*#With above commands, the practical problem statements are finished. But to study the the Iris Flower in the perspective of Biologists, the morphological (structural) differences between the species are to be st  
#Below is the extended code (not for exam)*

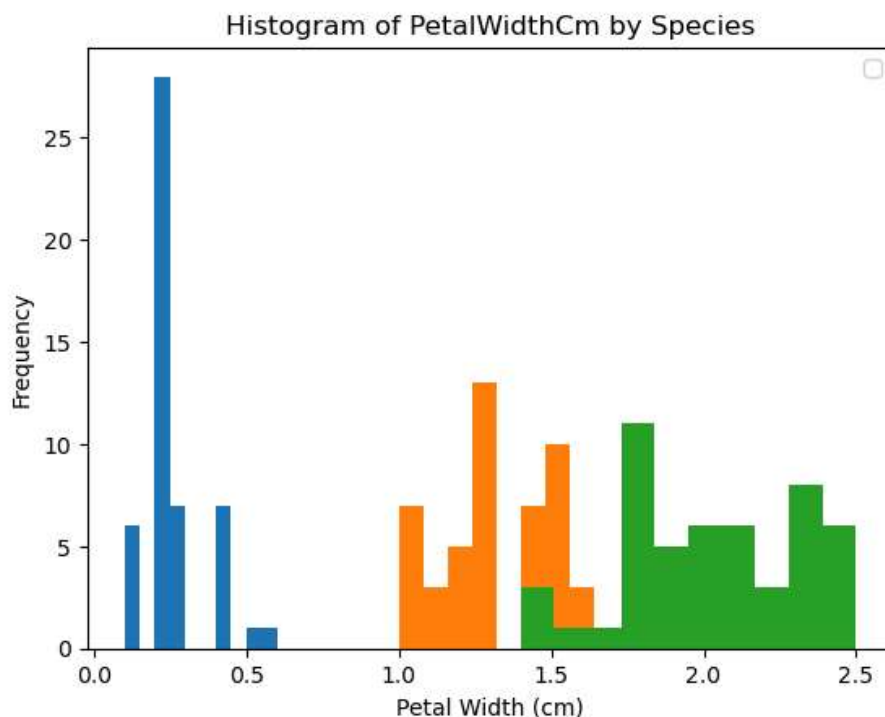
```
In [121]: # Lets do a comparative analysis of all species on PetalWidthCm
import seaborn as sns
sns.boxplot(x='Species', y='PetalWidthCm', data=df)
plt.title('Species-wise Boxplot of PetalWidthCm')
# Draw Specieswise BoxPlot for PetalWidthCm

plt.show()
```



```
In [122]: # Create a histogram of PetalWidthCm for each species
for species in df['Species'].unique():
    species_data = df[df['Species'] == species]
    plt.hist(species_data['PetalWidthCm'], bins=10) # Adjust bins as needed
plt.title('Histogram of PetalWidthCm by Species')
plt.xlabel('Petal Width (cm)')
plt.ylabel('Frequency')
plt.legend()
plt.show()
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



```
In [123]: # Calculate IQR and identify outliers for each species for PetalWidthCm
def find_outliers_iqr(data):
    Q1 = data.quantile(0.25)
    Q3 = data.quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = data[(data < lower_bound) | (data > upper_bound)]
    return outliers
```

```
In [124]: upper_bound
```

```
Out[124]: 4.05
```

```
In [125]: lower_bound
```

```
Out[125]: 2.05
```

```
In [126]: print("Printing Outlier for each species for PetalWidthCm")
for species in df['Species'].unique():
    species_data = df[df['Species'] == species]['PetalWidthCm']
    outliers = find_outliers_iqr(species_data)
    print(f"Outliers for {species}: {outliers.values}")
```

```
Printing Outlier for each species for PetalWidthCm
Outliers for Iris-setosa: [0.5 0.6]
Outliers for Iris-versicolor: []
Outliers for Iris-virginica: []
```

```
In [128]: # Return row indices for outliers for Iris-setosa
def find_outlier_indices(df, species_name, column_name):
    species_data = df[df['Species'] == species_name][column_name]
    outliers = find_outliers_iqr(species_data)
    outlier_indices = outliers.index
```

```
In [131]: outlier_indices
```

```
Out[131]: Int64Index([15, 32, 33, 60], dtype='int64')
```

```
In [133]: setosa_outlier_indices = find_outlier_indices(df, 'Iris-setosa', 'PetalWidthCm')
print(f"\nRow indices of outliers for Iris-setosa in 'PetalWidthCm': {setosa_outlier_indices.tolist()}")
```

```
Row indices of outliers for Iris-setosa in 'PetalWidthCm': [23, 43]
```