# Daffodil International University

## Machine Learning Driven Data Analysis I and Communicating Data Insights Lab

Course Code: DS422

# LAB REPORT

Submitted To:

**Musabbir Hasan Sammak**

Lecturer
Daffodil International University

Submitted By:

**Meher Durdana Khan Raisa**

ID: 192-35-2818

# SUPERVISED LEARNING

In [7]:

```python
import pandas as pd
import numpy as np
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score,
from sklearn.utils import resample

import seaborn as sns
import matplotlib.pyplot as plt

# Step I: Download and extract the Iris dataset
iris = load_iris()
features = pd.DataFrame(iris.data, columns=iris.feature_names)
targets = pd.DataFrame(iris.target, columns=["Species"])
```

Explore the datasets

```
In [8]:  1
         2
         3  print("Features shape:", features.shape)
         4  print("Features data types:\n", features.dtypes)
         5  print("Targets shape:", targets.shape)
         6  print("Targets data types:\n", targets.dtypes)
         7  print("Features dimensions:\n", features.head())
         8
         9  # Check for missing data
        10  print("Missing data:\n", features.isnull().sum())
        11
```

```
Features shape: (150, 4)
Features data types:
 sepal length (cm)    float64
sepal width (cm)     float64
petal length (cm)    float64
petal width (cm)     float64
dtype: object
Targets shape: (150, 1)
Targets data types:
 Species    int64
dtype: object
Features dimensions:
    sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)
0                5.1               3.5                1.4               0.2
1                4.9               3.0                1.4               0.2
2                4.7               3.2                1.3               0.2
3                4.6               3.1                1.5               0.2
4                5.0               3.6                1.4               0.2
Missing data:
 sepal length (cm)    0
sepal width (cm)     0
petal length (cm)    0
petal width (cm)     0
dtype: int64
```
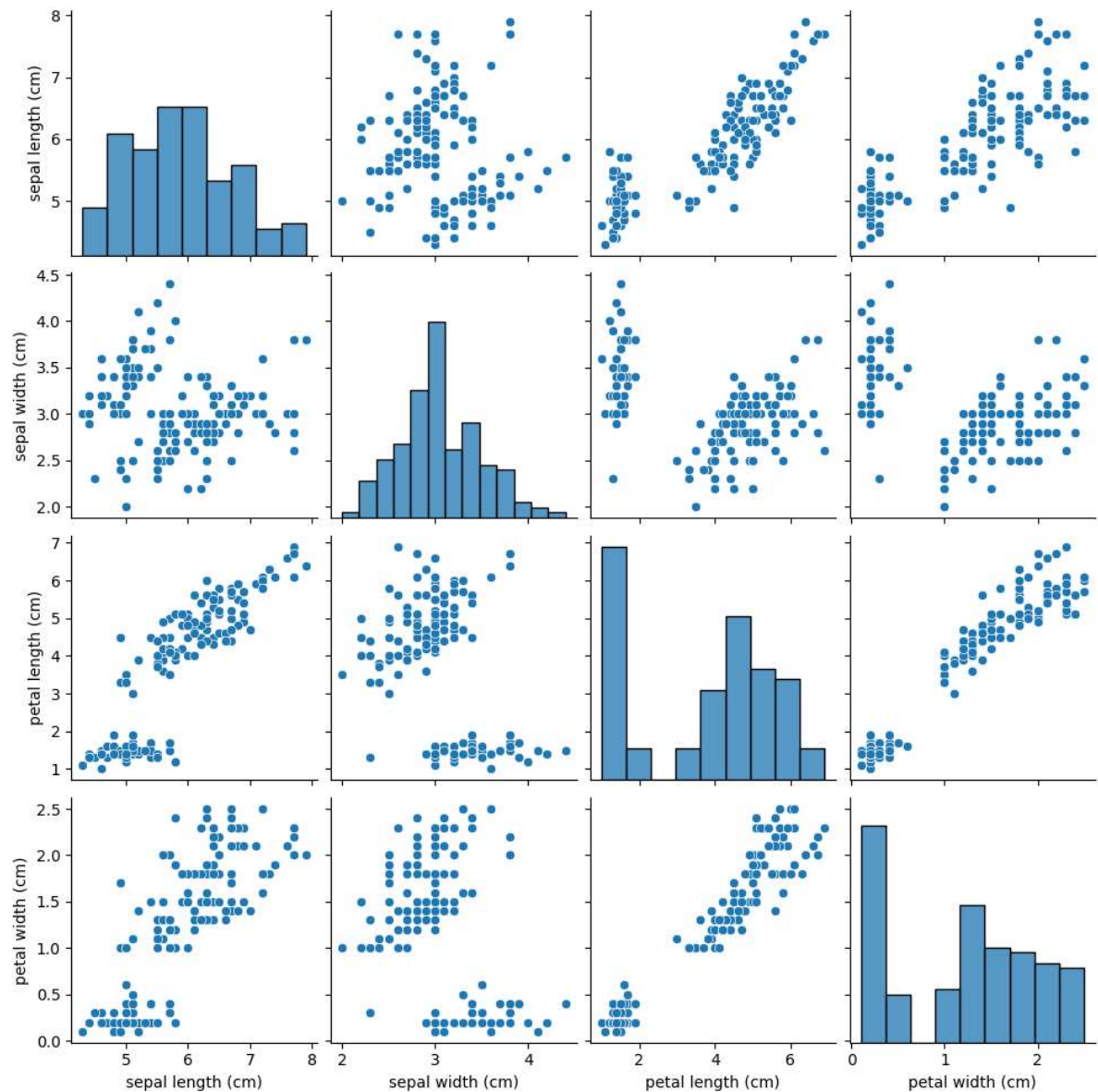
Visualize the features

```
1
2  sns.pairplot(features)
3  plt.show()
4
```



Preprocess the dataset

```
1
2
3  X_train, X_test, y_train, y_test = train_test_split(features, targets, tes
4
5  scaler = StandardScaler()
6  X_train_scaled = scaler.fit_transform(X_train)
7  X_test_scaled = scaler.transform(X_test)
8
```

Develop an Ensemble Model

```
1
2
3  logreg = LogisticRegression()
4  dt = DecisionTreeClassifier()
5  svc = SVC()
6  knn = KNeighborsClassifier()
7  nb = GaussianNB()
8
9  ensemble_model = VotingClassifier(estimators=[('lr', logreg), ('dt', dt),
10
```

Hyperparameter tuning

```
1
2  parameters = {'lr__C': [0.1, 1, 10],
3                'dt__max_depth': [None, 5, 10],
4                'svc__C': [0.1, 1, 10],
5                'knn__n_neighbors': [3, 5, 7],
6                'nb__var_smoothing': [1e-09, 1e-08, 1e-07]}
7
8  grid_search = GridSearchCV(estimator=ensemble_model, param_grid=parameters
9  grid_search.fit(X_train_scaled, y_train.values.ravel())
10
11  best_model = grid_search.best_estimator_
12
```

Evaluate training performance

```python
def evaluate_model(model, X, y):
    y_pred = model.predict(X)
    acc = accuracy_score(y, y_pred)
    precision = precision_score(y, y_pred, average='macro')
    recall = recall_score(y, y_pred, average='macro')
    f1 = f1_score(y, y_pred, average='macro')
    cm = confusion_matrix(y, y_pred)

    print("Confusion Matrix:\n", cm)
    print("Accuracy:", acc)
    print("Precision:", precision)
    print("Recall:", recall)
    print("F1 Score:", f1)

evaluate_model(best_model, X_train_scaled, y_train)
```

```
Confusion Matrix:
 [[40  0  0]
 [ 0 39  2]
 [ 0  1 38]]
Accuracy: 0.975
Precision: 0.975
Recall: 0.9751928288513655
F1 Score: 0.9749960931395529
```

Evaluate test performance

```python
In [14]:
 1
 2
 3  evaluate_model(best_model, X_test_scaled, y_test)
 4  n_iterations = 1000
 5  n_size = int(len(X_train_scaled))
 6
 7  accuracy_scores = []
 8  precision_scores = []
 9  recall_scores = []
10
11  for _ in range(n_iterations):
12      X_train_resampled, y_train_resampled = resample(X_train_scaled, y_trai
13      best_model.fit(X_train_resampled, y_train_resampled.values.ravel())
14      y_pred = best_model.predict(X_test_scaled)
15
16      accuracy_scores.append(accuracy_score(y_test, y_pred))
17      precision_scores.append(precision_score(y_test, y_pred, average='macro
18      recall_scores.append(recall_score(y_test, y_pred, average='macro'))
19
20  accuracy_mean = np.mean(accuracy_scores)
21  precision_mean = np.mean(precision_scores)
22  recall_mean = np.mean(recall_scores)
23
24  accuracy_ci = np.percentile(accuracy_scores, [2.5, 97.5])
25  precision_ci = np.percentile(precision_scores, [2.5, 97.5])
26  recall_ci = np.percentile(recall_scores, [2.5, 97.5])
27
28  print("Bootstrapping results:")
29  print("Accuracy Mean:", accuracy_mean)
30  print("Accuracy 95% CI:", accuracy_ci)
31  print("Precision Mean:", precision_mean)
32  print("Precision 95% CI:", precision_ci)
33  print("Recall Mean:", recall_mean)
34  print("Recall 95% CI:", recall_ci)
35
36
```

```
Confusion Matrix:
 [[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
Accuracy: 1.0
Precision: 1.0
Recall: 1.0
F1 Score: 1.0
Bootstrapping results:
Accuracy Mean: 1.0
Accuracy 95% CI: [1. 1.]
Precision Mean: 1.0
Precision 95% CI: [1. 1.]
Recall Mean: 1.0
Recall 95% CI: [1. 1.]
```

Compare training and test results

```
In [15]:    1
            2  train_acc = accuracy_score(y_train, best_model.predict(X_train_scaled))
            3  test_acc = accuracy_score(y_test, best_model.predict(X_test_scaled))
            4
            5  print("Training Accuracy:", train_acc)
            6  print("Test Accuracy:", test_acc)
```

```
Training Accuracy: 0.9583333333333334
Test Accuracy: 1.0
```

# UNSUPERVISED LEARNING

```
In [1]:     1  import pandas as pd
            2  import numpy as np
            3  import matplotlib.pyplot as plt
            4  from sklearn.datasets import load_iris
            5  from sklearn.cluster import KMeans
            6  from scipy.cluster.hierarchy import dendrogram, linkage
            7
            8  # I. Download the Iris dataset and extract features and targets
            9  iris = load_iris()
           10  df_features = pd.DataFrame(data=iris.data, columns=iris.feature_names)
           11  df_targets = pd.DataFrame(data=iris.target, columns=['target'])
           12
```

```python
# II. Explore the datasets
print("Features Shape:", df_features.shape)
print("Targets Shape:", df_targets.shape)
print("\nFeatures Data Types:")
print(df_features.dtypes)
print("\nFeatures Dimensions:")
print(df_features.head())

# Check for missing data
print("\nMissing Data:")
print(df_features.isnull().sum())
```

```
Features Shape: (150, 4)
Targets Shape: (150, 1)

Features Data Types:
sepal length (cm)    float64
sepal width (cm)     float64
petal length (cm)    float64
petal width (cm)     float64
dtype: object

Features Dimensions:
   sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)
0                5.1               3.5                1.4               0.2
1                4.9               3.0                1.4               0.2
2                4.7               3.2                1.3               0.2
3                4.6               3.1                1.5               0.2
4                5.0               3.6                1.4               0.2

Missing Data:
sepal length (cm)    0
sepal width (cm)     0
petal length (cm)    0
petal width (cm)     0
dtype: int64
```

```python
# IIIa. K-Means clustering using Petal Length and Petal Width
petal_data = df_features[['petal length (cm)', 'petal width (cm)']]
k_values = range(1, 10)
inertia = []

for k in k_values:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(petal_data)
    inertia.append(kmeans.inertia_)

# Plotting elbow curve
plt.plot(k_values, inertia, 'bx-')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Inertia')
plt.title('Elbow Method - Petal Length & Petal Width')
plt.show()

# Perform K-Means clustering with optimal k=3
kmeans_petal = KMeans(n_clusters=3, random_state=42)
kmeans_petal.fit(petal_data)
petal_labels = kmeans_petal.labels_

# Visualize clusters
plt.scatter(petal_data['petal length (cm)'], petal_data['petal width (cm)'
plt.xlabel('Petal Length (cm)')
plt.ylabel('Petal Width (cm)')
plt.title('K-Means Clustering - Petal Length & Petal Width')
plt.show()
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
```

Elbow Method - Petal Length & Petal Width

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
```
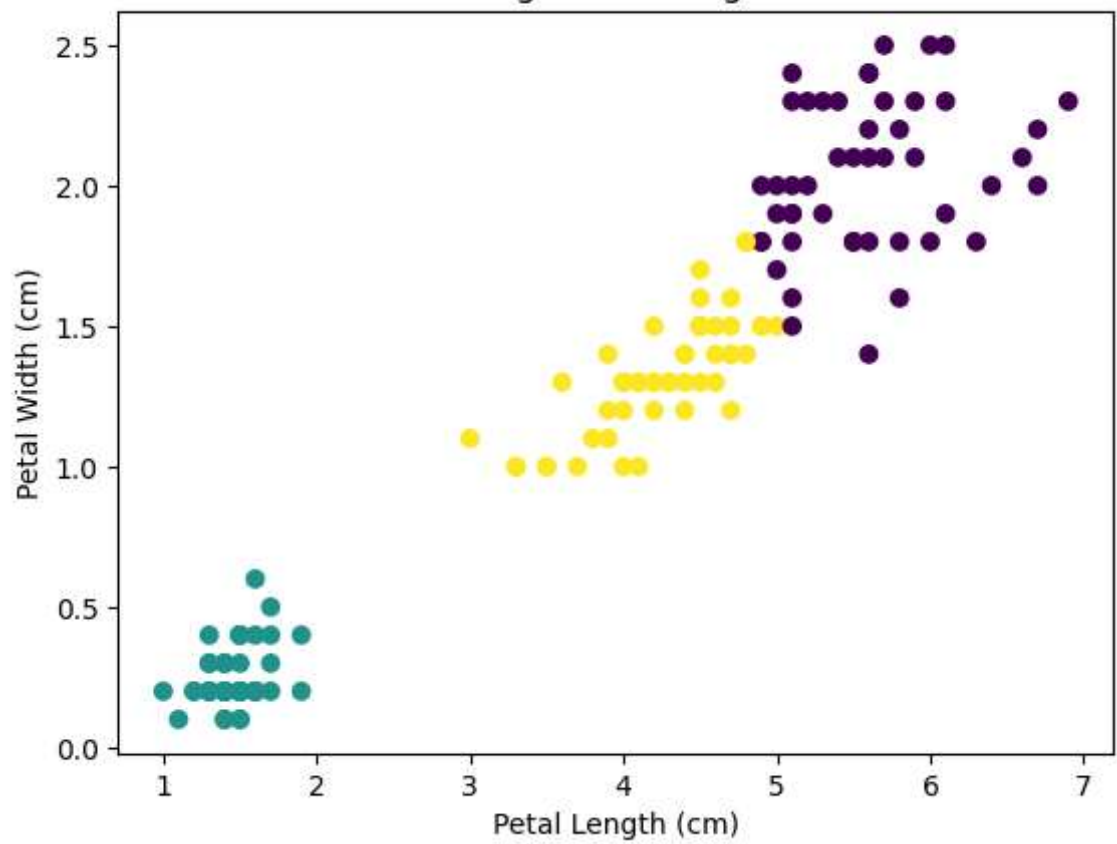
K-Means Clustering - Petal Length & Petal Width

```python
# IIIb. K-Means clustering using Sepal Length and Sepal Width
sepal_data = df_features[['sepal length (cm)', 'sepal width (cm)']]
inertia = []

for k in k_values:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(sepal_data)
    inertia.append(kmeans.inertia_)

# Plotting elbow curve
plt.plot(k_values, inertia, 'bx-')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Inertia')
plt.title('Elbow Method - Sepal Length & Sepal Width')
plt.show()

# Perform K-Means clustering with optimal k=2
kmeans_sepal = KMeans(n_clusters=2, random_state=42)
kmeans_sepal.fit(sepal_data)
sepal_labels = kmeans_sepal.labels_

# Visualize clusters
plt.scatter(sepal_data['sepal length (cm)'], sepal_data['sepal width (cm)'
plt.xlabel('Sepal Length (cm)')
plt.ylabel('Sepal Width (cm)')
plt.title('K-Means Clustering - Sepal Length & Sepal Width')
plt.show()
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
```

Elbow Method - Sepal Length & Sepal Width

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
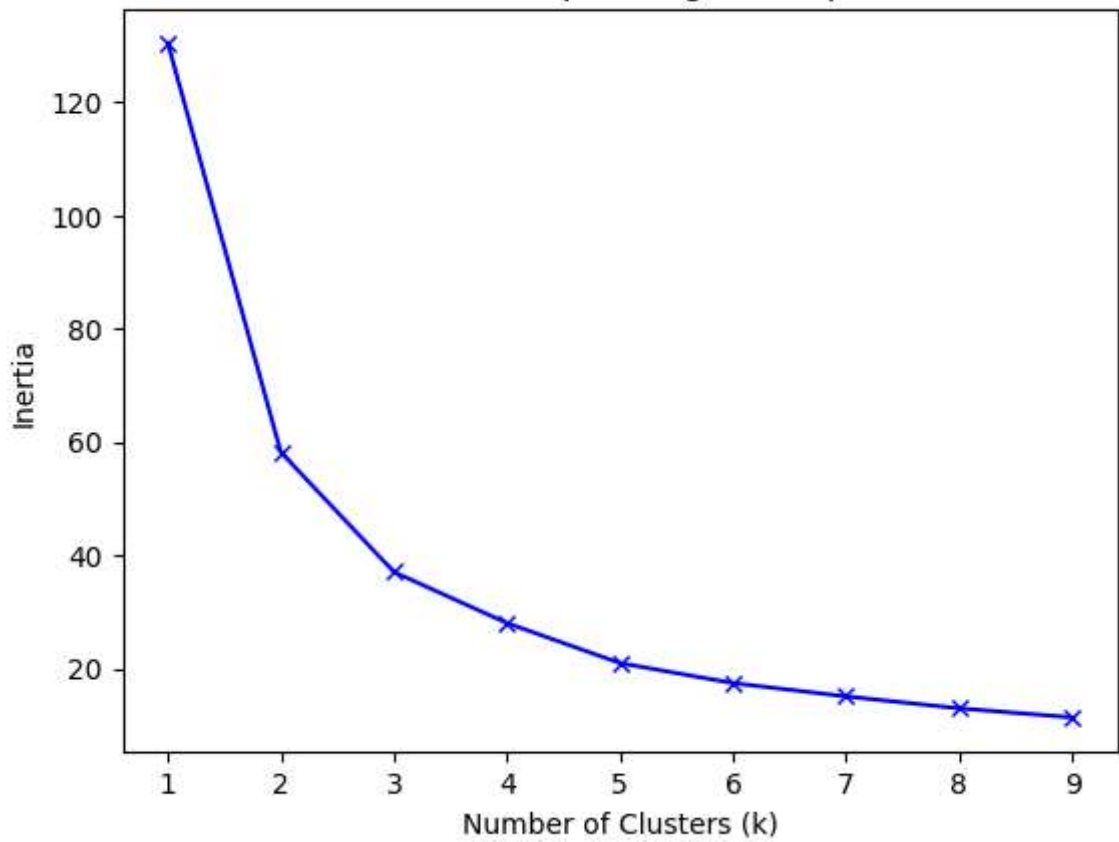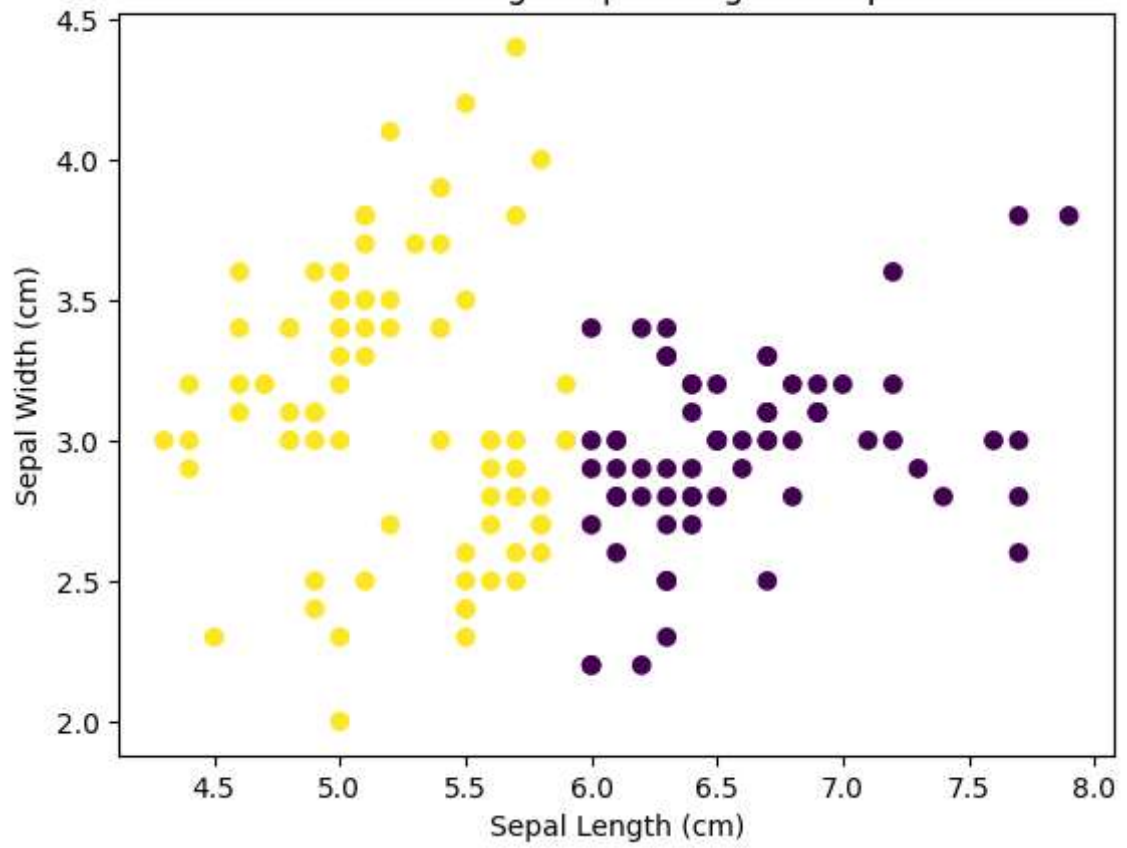  warnings.warn(

K-Means Clustering - Sepal Length & Sepal Width

```python
In [5]:
# IIIc. K-Means clustering using all features
inertia = []

for k in k_values:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(df_features)
    inertia.append(kmeans.inertia_)

# Plotting elbow curve
plt.plot(k_values, inertia, 'bx-')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Inertia')
plt.title('Elbow Method - All Features')
plt.show()

# Perform K-Means clustering with optimal k=3
kmeans_all = KMeans(n_clusters=3, random_state=42)
kmeans_all.fit(df_features)
all_labels = kmeans_all.labels_
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
```
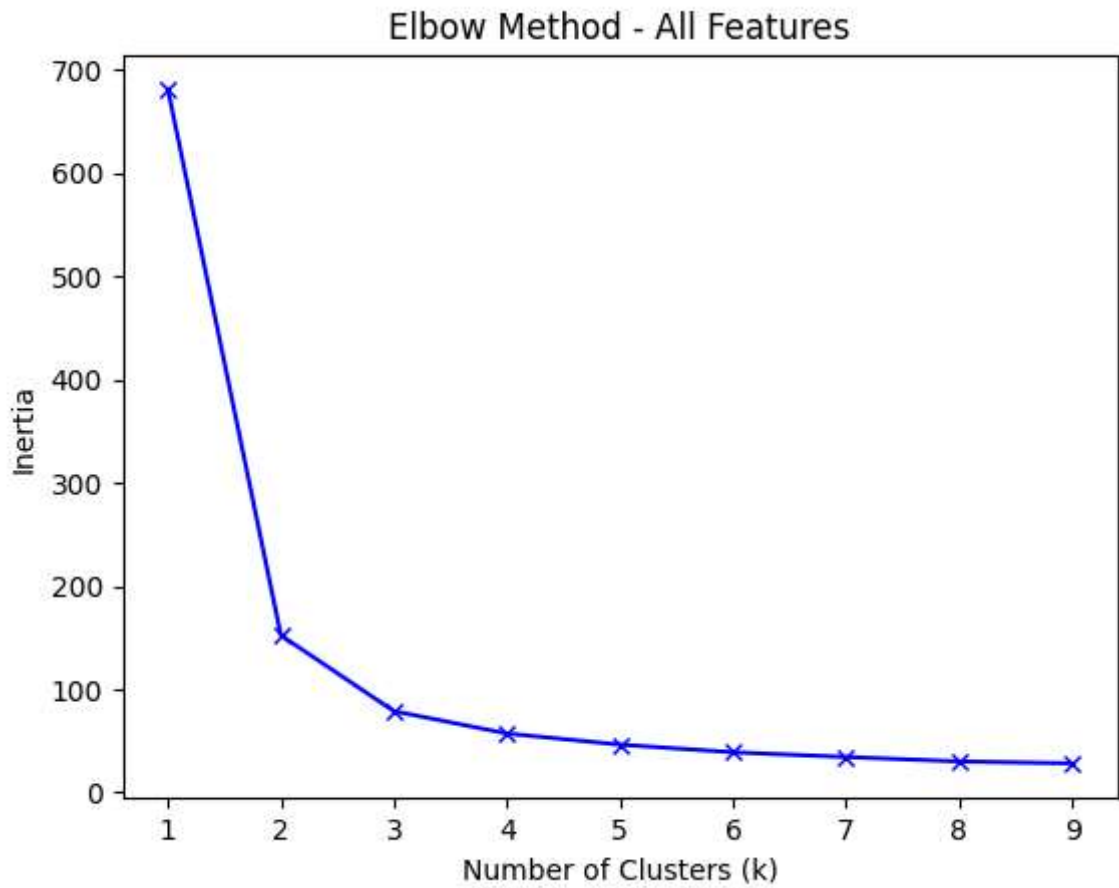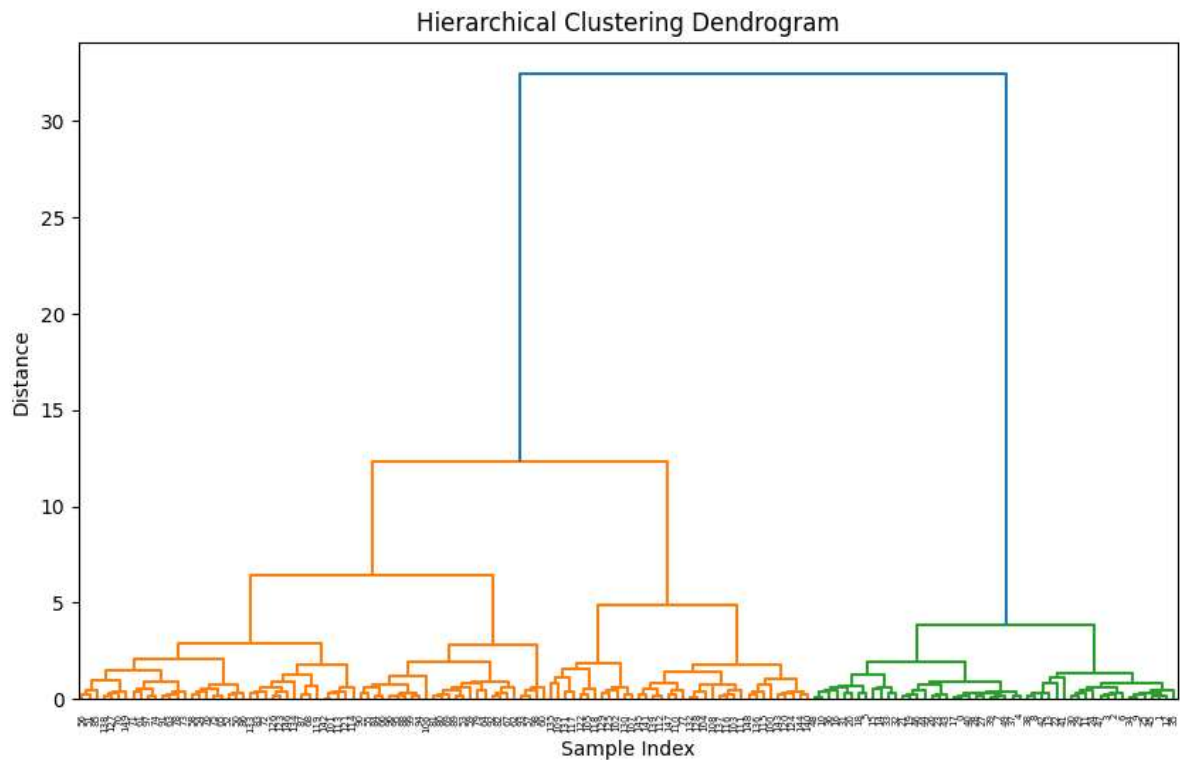
Elbow Method - All Features

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Futur
eWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
  warnings.warn(
```

```
1
2  # IV. Hierarchical clustering
3  linked = linkage(df_features, 'ward')
4
5  plt.figure(figsize=(10, 6))
6  dendrogram(linked, orientation='top', distance_sort='descending', show_lea
7  plt.title('Hierarchical Clustering Dendrogram')
8  plt.xlabel('Sample Index')
9  plt.ylabel('Distance')
10 plt.show()
11
```



Hierarchical Clustering Dendrogram

V. Comparison between K-Means and Hierarchical clustering:

K-Means and Hierarchical clustering have their own strengths and weaknesses. K-Means is a partition-based clustering algorithm that assigns each data point to a single cluster, and it assumes equal-sized and spherical clusters. Hierarchical clustering, on the other hand, creates a hierarchy of clusters and does not assume equal-sized clusters.

The choice between the two depends on the nature of the dataset and the desired outcome. K-Means is computationally efficient and works well when the clusters are relatively well-separated and have a spherical shape. Hierarchical clustering is more flexible and can handle different cluster shapes, but it can be computationally expensive for large datasets.

In this particular case, based on the plots and the elbow method, K-Means with k=3 seems to perform reasonably well for both the Petal and Sepal features. However, the clusters based on the Petal features seem to have a clearer separation and correlation with the target variable, as observed in the scatter plot.