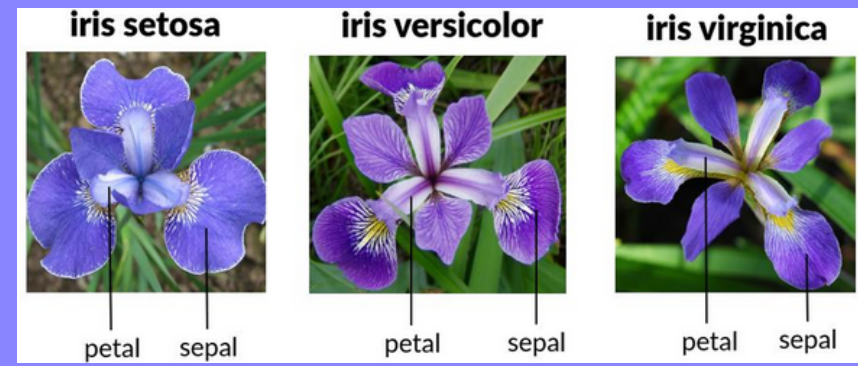


# UNSUPERVISED MACHINE LEARNING WITH IRIS DATASET

K-Means Clustering

by: Raniah Mufidah Admayana

Code in Google Colab = <https://colab.research.google.com/drive/1uipfdSnZYJrdK4kztmKuFOIZhvaVdRIr?usp=sharing>.



Three classes of IRIS dataset (Mijwil & Abttan, 2021)

## INTRODUCTION

### IRIS DATASET

Iris dataset is a classic dataset in machine learning and statistics. It classifies three types of iris flowers based on sepal and petal measurements and was originally analyzed by R.A. Fisher in the 1930s (Marsland, 2014).

### UNSUPERVISED MACHINE LEARNING

Unsupervised machine learning finds patterns and relationships in data without labeled outcomes. It identifies hidden structures, trends, or groups, making it useful when labeled data is unavailable. This approach helps uncover insights and can complement supervised learning to enhance models. Example of unsupervised machine learning methods:

- Clustering
  - Anomaly detection
  - Dimensionality reduction
  - Association rule-mining
- (Sarkar et al., 2017).

### RESEARCH PURPOSE

Evaluate the performance of the K-Means clustering model on the Iris dataset.

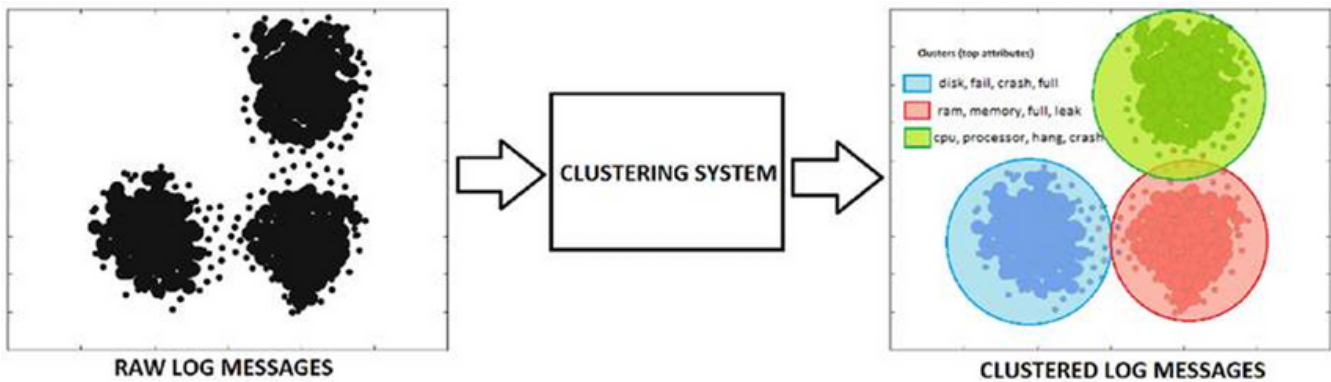
## LITERATURE REVIEW

### CLUSTERING

Clustering is grouping similar data points into different clusters based on their features. It identifies patterns and relationships in the data without prior labels or supervision.

Types of Clustering Methods:

- Centroid-based (e.g., K-Means, K-Medoids)
  - Distribution-based (e.g., Gaussian Mixture Models)
  - Hierarchical (e.g., Agglomerative, Divisive)
  - Density-based (e.g., DBSCAN, OPTICS)
- (Sarkar et al., 2017).



An overview of the working process of Clustering (Sarkar et al., 2017)

### K-Means Clustering

K-Means clustering works by grouping data into k clusters (k is chosen by the user).

The algorithm works as follows:

1. k cluster “center” points are created at random locations.
2. For each observation:
  - a. The distance between each observation and the k center points is calculated.
  - b. The observation is assigned to the cluster of the nearest center point.
3. The center points are moved to the means (i.e., centers) of their respective clusters.
4. Steps 2 and 3 are repeated until no observation changes in cluster membership

K-Means assumes that clusters are roughly circular and of similar size. The key measure used is the distance between each data point and the centroids (Albon, 2018).

### MODEL EVALUATION

#### Sillhoutte Score

Measures how well data points fit within their assigned cluster compared to other clusters.

The score ranges from -1 to 1:

- Close to 1 → The data point is well-clustered.
  - Close to 0 → The data point is between two clusters.
  - Close to -1 → The data point may belong to a different cluster.
- (Rousseeuw, 1987).

#### Adjusted Rand Index (ARI)

ARI measures the similarity between the clustering results generated by the algorithm and the original labels, taking into account the possibility of random clustering.

ARI values range from -1 to 1:

- 1 → Perfect clustering according to the original labels.
  - 0 → Clustering occurs randomly.
  - < 0 → Clustering is worse than random clustering.
- (Scikit-learn developers, n.d.)

## RESEARCH METHODOLOGY

**DATA SOURCE** → secondary data from the website <https://www.kaggle.com/datasets/uciml/iris>, accessed on December 17, 2024.

### RESEARCH VARIABLES

**Independent Variables (Features):**

represented by the symbol “X”

- Sepal Length (cm)
- Petal Length (cm)
- Sepal Width (cm)
- Petal Width (cm)

**Dependent Variable (for ARI):**

represented by the symbol “y”

Species (Iris-setosa, Iris-versicolor, Iris-virginica)

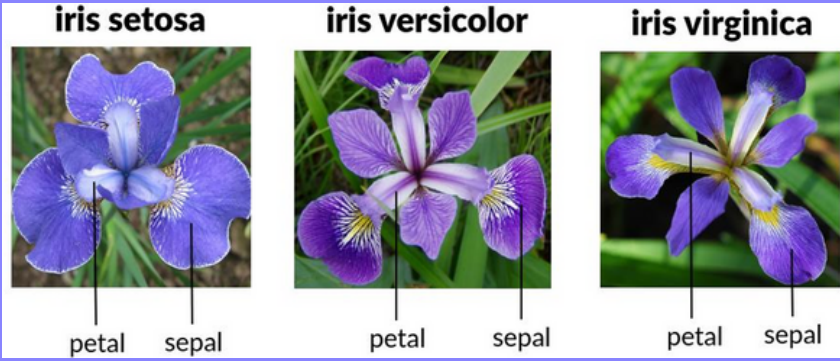
### ANALYSIS STEPS

1. Identify the problem and input the dataset.
2. Exploratory Data Analysis.
3. Train the models and evaluate their performance.
4. Summarize the results and draw conclusions.



# UNSUPERVISED MACHINE LEARNING WITH IRIS DATASET

K-Means Clustering  
by: Raniah Mufidah Admayana  
Code in Google Colab = <https://colab.research.google.com/drive/1uipfdSnZYJrdK4kztmKuFOIZhvaVdRIr?usp=sharing>



Three classes of IRIS dataset (Mijwil & Abttan, 2021)

## DISCUSSION

### 1. INPUT THE DATASET

change all column names into lowercase

```
df = pd.read_csv('/content/iris.csv')
```

```
df.columns = df.columns.str.lower().str.replace(' ', '_')
```

```
df.head()
```

	id	sepalengthcm	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

```
# Divide the dataset into X and y
X = df[['sepalengthcm', 'sepalwidthcm', 'petallengthcm', 'petalwidthcm']]
y = df['species']
```

- X is for k-means clustering
- y is for model evaluation on Adjusted Rand Index (ARI)

### 2. EXPLORATORY DATA ANALYSIS

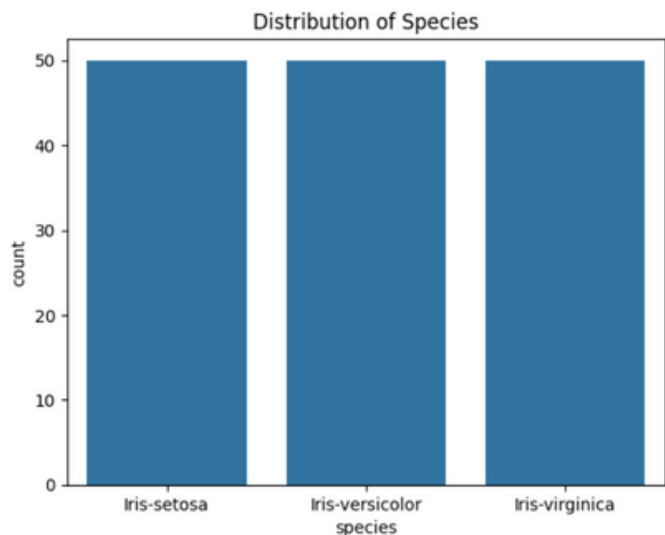
Clean, explore, and visualize the data

```
df.isnull().sum()
# To see the missing values
```

	0
id	0
sepalengthcm	0
sepalwidthcm	0
petallengthcm	0
petalwidthcm	0
species	0

No missing value

```
# Species distribution
sns.countplot(x='species', data=df)
plt.title('Distribution of Species')
plt.show()
```



the amount of data between species is evenly distributed

```
# Species distribution
print(df['species'].value_counts())

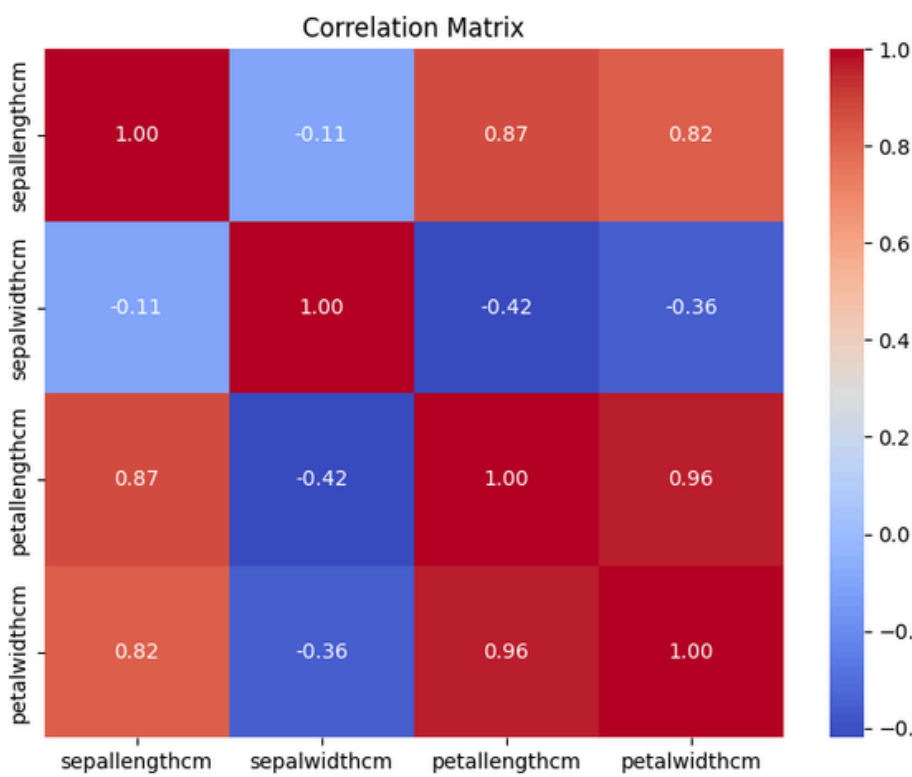
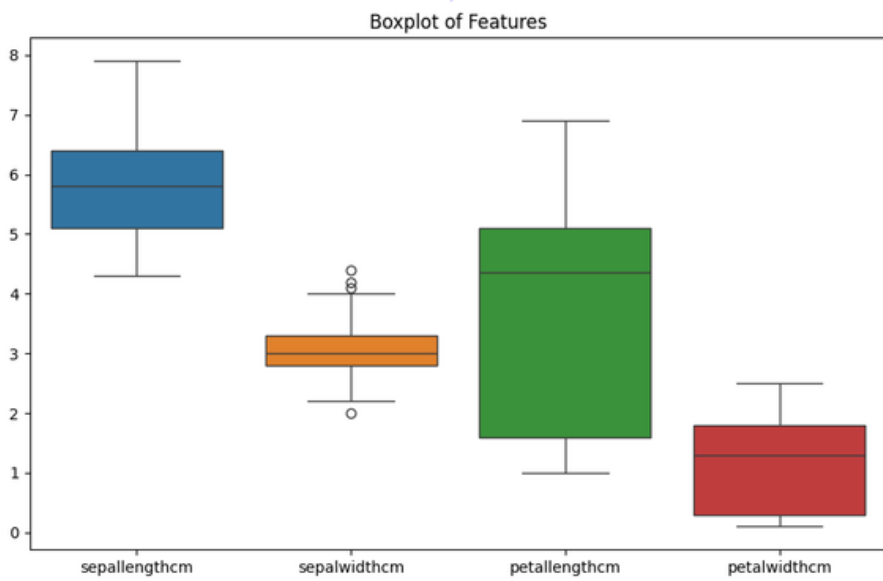
species
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
Name: count, dtype: int64
```

### Data Distribution & Outlier Detection

```
X.describe()
```

	sepalengthcm	sepalwidthcm	petallengthcm	petalwidthcm
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

Outliers can be ignored



### Correlation between features

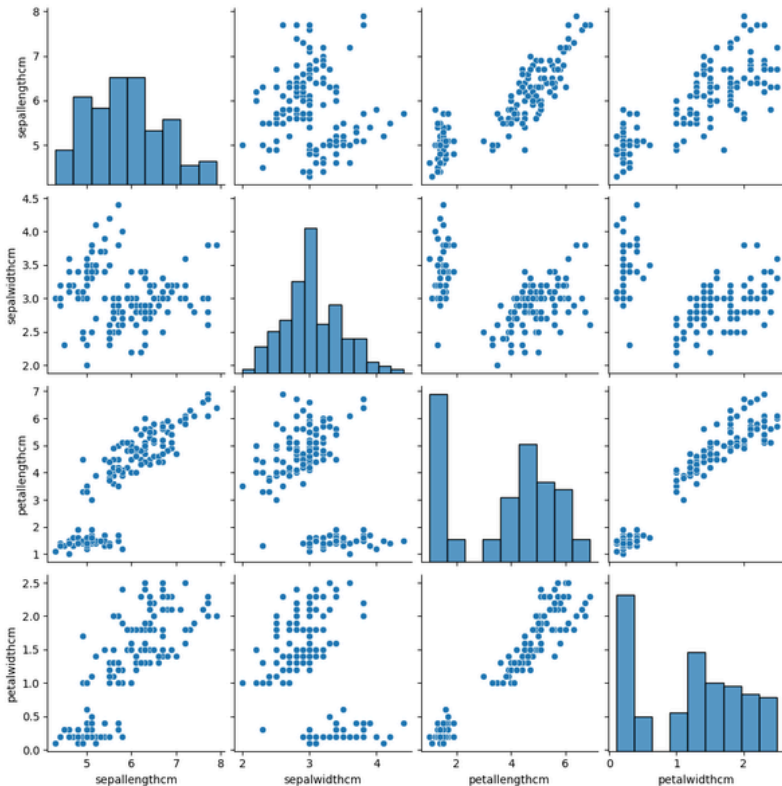
Correlation Matrix shows the strength of the relationship between features.

- a number close to 1 or -1 → have strong correlation
- 0 → don't have a correlation

(Sarkar et al., 2017)

Features that have strong correlation

- petallengthcm & petalwidthcm (0.96)
- sepalengthcm & petallengthcm (0.87)
- sepalengthcm & petalwidthcm (0.82)







## REFERENCE