# Introduction to Artificial Intelligence (AI) and Machine Learning (ML)

Zhenhua He   |   Ridham Patoliya

# Learning objectives

Terminology of Machine Learning

The difference between AI and ML

The different types of machine learning techniques

Applications of machine learning techniques

# Data Exploration



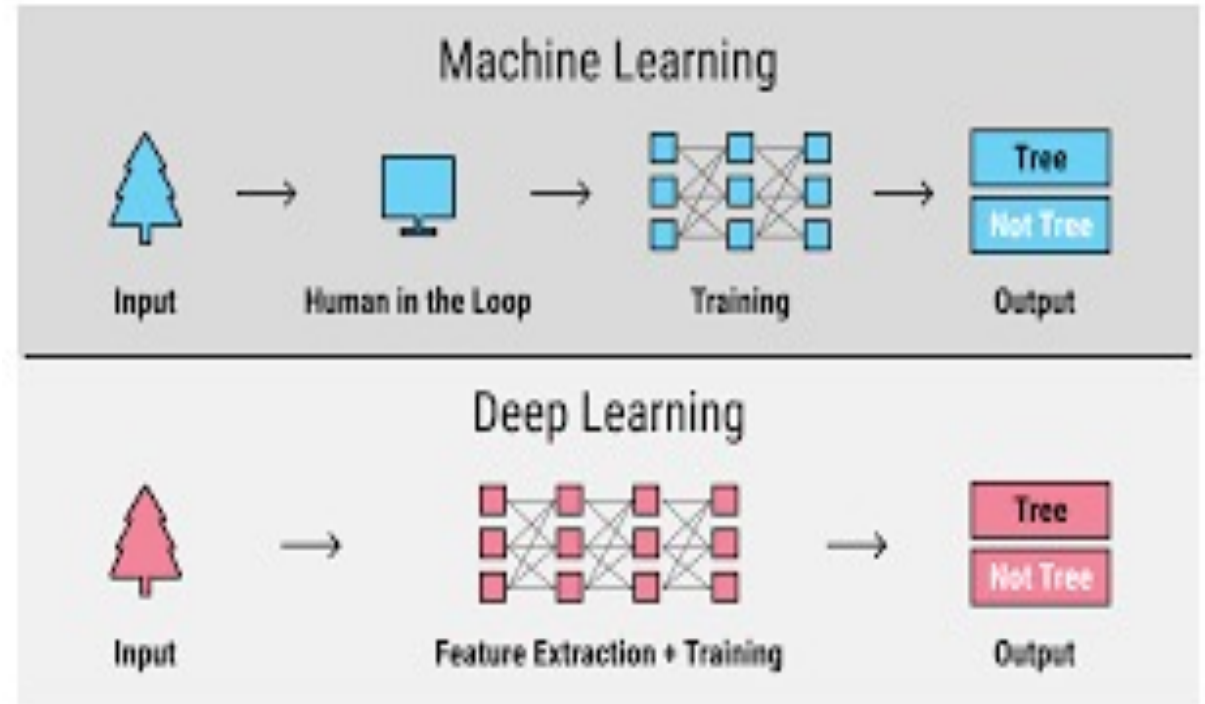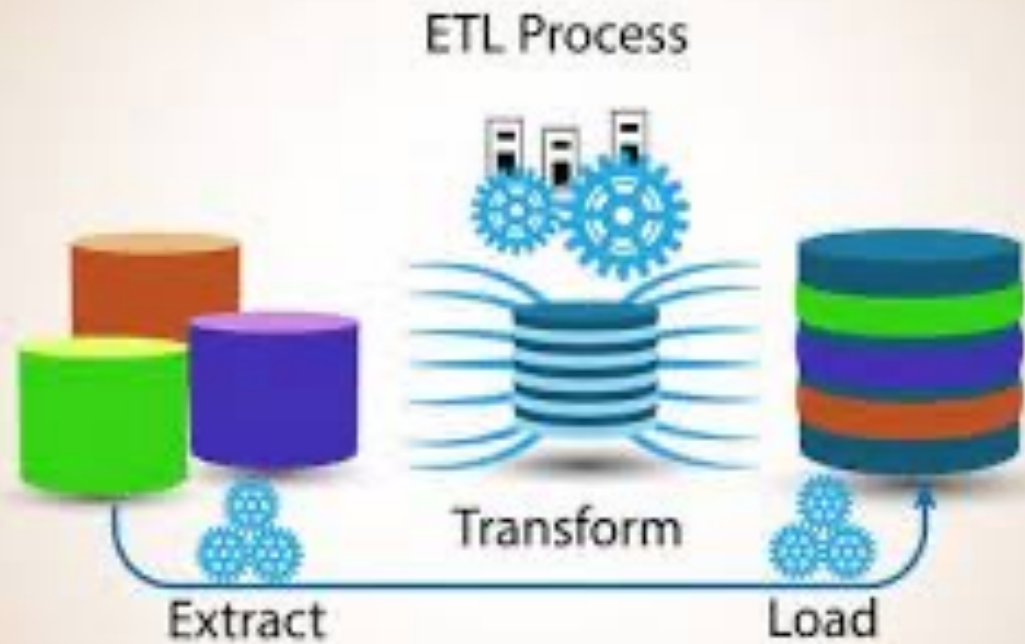Data manipulation/analysis library
- Exercises



Data visualization library
- Exercises

# ENTERPRISE DATA   vs   MACHINE LEARNING DATA
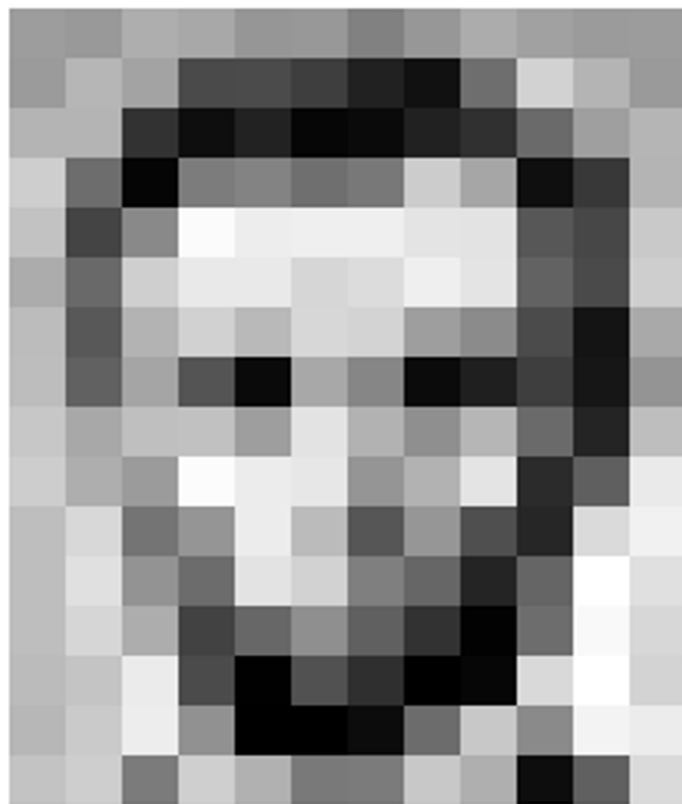
# Formats of data

NUMBERS

TEXTS

IMAGES

# Image (as seen by computers)

# Types of data:

- Labelled
- Unlabelled

# Labelled data



Dog

Cat

Dog

Dog

Cat

Cat

Dog

Cat

# Unlabelled Data

# Neural Networks and Biology



Synapses: connection to the next neuron (release of neurotransmitter)

Axon: delivers output from neuron to other neurons

Dendrites: to facilitate the connection with other neurons

Cell Body

Nucleus

Figure: Structure of a typical neuron

# Neural Network

# What is training?

The process used to create our ML model.

Find a set of weights and biases that have high accuracy.

# What is testing?

The process used to test our ML model.

Run the model against known outcomes

# What is inference?

Running our model on live data to produce actionable output.

# Common types of Learning

**Supervised learning**

**We have <mark>labelled data,</mark> and we want to make some prediction**

- Regression
- Classification

**Unsupervised learning**

**We have <mark>unlabeled data,</mark> and we want to make some prediction**

- Clustering

# Supervised learning

# Regression

# Regression



Y-axis
Housing Price

$0.1 million

1000 square feet

X-axis – area in sq. feet

# Regression

# Quiz

- Which of the following CANNOT be an example of regression?

  - A) Using past data of weather in college station to predict future's weather.

  - B) Predicting prices of stocks using previous month's price data

  - C) Determining if an email is spam or not

  - D) Determining network traffic for today using previous month's data

Classification

# Classification



Decision Boundary

Chipmunk · Gopher

**Rodent Classification**

Length in cm

20 cm
15 cm

75 g · 170 g

Weight (mass in grams)

# Quiz

- Which of the following CANNOT be an example of classification?

  - A) Using blood pressure and weight data to determine if a patient is diabetic or not

  - B) Estimating amount of annual rain from previous year's data

  - C) Classifying Pokémon in different types (e.g., fire, ice, poison, electric)

  - D) Determining if an email is spam or not

# Unsupervised learning

clustering

# Clustering



Species 1

Species 2

Species 3

**Rodent Clustering**

Length in cm

20 cm

15 cm

75 g

170 g

Weight (mass in grams)

# Quiz

- Which of the following CANNOT be an example of clustering?

  - A) Sorting and making groups of research papers having similar content

  - B) Determining whether a news article is about politics or sports

  - C) Identifying clusters of stars having similar characteristics

  - D) Sorting through subjects of emails and grouping them accordingly

# Quiz

- Which of the following CANNOT be an example of machine learning? Select all that apply.

    - A) Manually trying out different passwords on your amazon account to check if it works

    - B) Your virtual assistant starts recognizing your voice after first few tries

    - C) Fire alarm goes off when smoke level is more than a specific level

    - D) Sorting through subjects of emails and grouping them accordingly

# What is Artificial Intelligence

- **Wikipedia**: intelligence demonstrated by machines as opposed to natural intelligence displayed by animals including humans.

- **Oxford**: the theory and development of computer systems able to perform tasks that normally require human intelligence.

- **IBM**: leverages computers and machines to mimic the problem-solving and decision-making capabilities of the human mind

# Train a linear regression model

Jupyter Notebooks on UArizona HPC with Python.

ood.hpc.arizona.edu

# Jupyter Notebooks on UArizona HPC with Python

# Jupyter Notebooks on UArizona HPC with Python

**Interactive Apps**

Desktops

🖵 Interactive Desktop

GUIs

⚖ ABAQUS GUI

Ⓐ ANSYS Workbench GUI

◣ MATLAB GUI

✳ Mathematica GUI

Servers

☕ Jupyter Notebook

⬢ RStudio Server

---

## Jupyter Notebook

This app will launch a Jupyter server using Python on a UAz cluster.

**Cluster**

| Ocelote Cluster | ⏶ |
|---|---|

**Run Time**

| 1 | ⏳ |
|---|---|

Enter maximum number of wall clock hours the job is allowed to run.

**Core count on a single node**

| 1 | ⏳ |
|---|---|

Enter the number of cores on a single node that the job is allowed to use.

**Memory per core**

| 6 | ⏳ |
|---|---|

Enter the number of Gigabytes of RAM needed per core.

**Special Options**

| |
|---|

Enter node specific requirements, if any.

**PI Group**

| chrisreidy |
|---|

Enter an HPC PI group to be charged for time used.

# Jupyter Notebooks on UArizona HPC with Python

# Jupyter Notebooks on UArizona HPC with Python

# Jupyter Notebooks on UArizona HPC with Python

# Jupyter Notebooks on HPC

ood.hpc.arizona.edu

ocelote / 2 hours / 1 core / 6 mem / standard queue / chrisreidy

**Accessing files for the exercises**

ssh netid@hpc.arizona.edu

shell

ocelote

mkdir intro-to-hpc

cd intro-to-hpc

https://ua-researchcomputing-hpc.github.io/Intro-to-HPC/

Then Accessing Workshop Files and cut / paste the section starting "wget"

(old method:  cp /xdisk/chrisreidy/workshops/* .)

Choice #1: Cut and paste commands into Jupyter from .txt file

Choice #2: Run the Notebook .ipynb file

Choice #3: Type in the commands. Syntax is very important

# Train a linear regression model

- Import libraries

```
# Import libraries

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
```

- Use Pandas to load the data and view the first 5 rows

```
# Load data and view the first 5 rows
data = pd.read_excel("king_county_house_data.xlsx")

data.head(5)
```

| | id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors |
|---|---|---|---|---|---|---|---|---|
| 0 | 7129300520 | 20141013T000000 | 221900 | 3 | 1.00 | 1180 | 5650 | 1.0 |
| 1 | 6414100192 | 20141209T000000 | 538000 | 3 | 2.25 | 2570 | 7242 | 2.0 |
| 2 | 5631500400 | 20150225T000000 | 180000 | 2 | 1.00 | 770 | 10000 | 1.0 |
| 3 | 2487200875 | 20141209T000000 | 604000 | 4 | 3.00 | 1960 | 5000 | 1.0 |
| 4 | 1954400510 | 20150218T000000 | 510000 | 3 | 2.00 | 1680 | 8080 | 1.0 |

# Train a linear regression model

- Choose the columns from the data
- Split the data into train and test sets

```python
space = data['sqft_living']
price = data['price']

# Change X into 2D array
X = np.array(space).reshape(-1, 1)
Y = np.array(price)

# Split data into train sets and test sets
X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=1/3,random_state=0)
```

- Visualize the train set

```python
# Visualize training set
plt.scatter(X_train,Y_train,color="red",label="Living Area")
plt.title("Housing Prices in King County, WA")
plt.xlabel("Area (sq-ft)")
plt.ylabel("Price (USD)")
plt.legend()
plt.show()
```

# Train a linear regression model

- Train the model with train set

- Predict on test set

```python
# Train
regressor = LinearRegression()
regressor.fit(X_train, Y_train)

# Prediction
y_pred = regressor.predict(X_test)
```

- Visualize the train data and the best fit line

```python
# Visualize the data and the bestfit line
plt.scatter(X_train,Y_train,color="red",label="Living Area")
plt.title("Housing Prices in King County, WA")
plt.plot(X_train,regressor.predict(X_train),color="blue",label="Price")
plt.xlabel("Area (sq-ft)")
plt.ylabel("Price (USD)")
plt.legend()
plt.show()
```

# Train a linear regression model

- Predict the price of a house with a certain area

```
area = 5000

price = regressor.predict([[area]])

print('House of %d sq-ft costs about $%d' % (area, price))

House of 5000 sq-ft costs about $1339969
```

- Visualize the test data

```
# Visualize test set
plt.scatter(X_test,Y_test,color='red',label="Living Area")
plt.plot(X_test,regressor.predict(X_test),color="blue",label="Price")
plt.xlabel("Area (sq-ft)")
plt.ylabel("Price (USD)")
plt.legend()
plt.show()
```

# Build a clustering model for Iris Dataset

# Build a clustering model – Iris dataset



Iris setosa    Iris versicolor    Iris virginica

- Import libraries

```
[1] #import libraries
    import numpy as np
    import pandas as pd
    import matplotlib.pyplot as plt
    import seaborn as sns
    %matplotlib inline
    from sklearn.cluster import KMeans
    from sklearn.datasets import load_iris
```

# Build a clustering model – Iris dataset

- Load the data

```
[2] iris=load_iris()
    iris

    {'DESCR': '.. _iris_dataset:\n\nIris plants dataset\n--
     'data': array([[5.1, 3.5, 1.4, 0.2],
            [4.9, 3. , 1.4, 0.2],
            [4.7, 3.2, 1.3, 0.2],
            [4.6, 3.1, 1.5, 0.2],
            [5. , 3.6, 1.4, 0.2],
            [5.4, 3.9, 1.7, 0.4],
            [4.6, 3.4, 1.4, 0.3],
            [5. , 3.4, 1.5, 0.2],
            [4.4, 2.9, 1.4, 0.2],
            [4.9, 3.1, 1.5, 0.1],
```

```
[3] df=pd.DataFrame(data=iris.data, columns=['sepal length','sepal width','petal length','petal width'])
    df['target']=pd.Series(iris.target)
    df
```

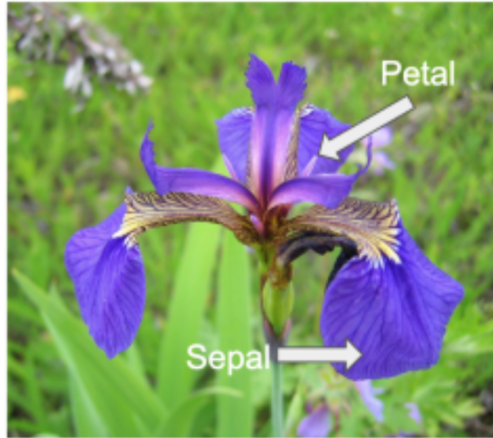|     | sepal length | sepal width | petal length | petal width | target |
|-----|--------------|-------------|--------------|-------------|--------|
| 0   | 5.1          | 3.5         | 1.4          | 0.2         | 0      |
| 1   | 4.9          | 3.0         | 1.4          | 0.2         | 0      |
| 2   | 4.7          | 3.2         | 1.3          | 0.2         | 0      |
| 3   | 4.6          | 3.1         | 1.5          | 0.2         | 0      |
| 4   | 5.0          | 3.6         | 1.4          | 0.2         | 0      |
| ... | ...          | ...         | ...          | ...         | ...    |
| 145 | 6.7          | 3.0         | 5.2          | 2.3         | 2      |

# Build a clustering model – Iris dataset

- Visualize the data

# Build a clustering model – Iris dataset

- Estimate k with elbow method- first try k = 5

```
[5]  # Let's first try k = 5
     x = iris.data
     kmeans5 = KMeans(n_clusters=5,init = 'k-means++', random_state = 0)
     y = kmeans5.fit_predict(x)
     print(y)

     [0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
      0 0 0 0 0 0 0 0 0 0 0 0 0 4 4 4 3 4 4 4 3 4 3 3 4 3 4 3 4 4 3 4 3 4 3 4 4
      4 4 4 4 4 3 3 3 3 4 3 4 4 4 3 3 3 4 3 3 3 3 3 4 3 3 1 4 2 1 1 2 3 2 1 2 1
      1 1 4 1 1 1 2 2 4 1 4 2 4 1 2 4 4 1 2 2 2 1 4 4 2 1 1 4 1 1 1 4 1 1 1 4 1
      1 4]


[6]  kmeans5.cluster_centers_

     array([[5.006     , 3.428     , 1.462     , 0.246     ],
            [6.52916667, 3.05833333, 5.50833333, 2.1625    ],
            [7.475     , 3.125     , 6.3       , 2.05      ],
            [5.508     , 2.6       , 3.908     , 1.204     ],
            [6.20769231, 2.85384615, 4.74615385, 1.56410256]])
```

# Build a clustering model – Iris dataset

- Estimate k with elbow method

```
[7] plt.scatter(x[y == 0,0], x[y==0,1], s = 15, c= 'red', label = 'Cluster_1')
    plt.scatter(x[y == 1,0], x[y==1,1], s = 15, c= 'blue', label = 'Cluster_2')
    plt.scatter(x[y == 2,0], x[y==2,1], s = 15, c= 'green', label = 'Cluster_3')
    plt.scatter(x[y == 3,0], x[y==3,1], s = 15, c= 'cyan', label = 'Cluster_4')
    plt.scatter(x[y == 4,0], x[y==4,1], s = 15, c= 'magenta', label = 'Cluster_5')

    plt.scatter(kmeans5.cluster_centers_[:,0], kmeans5.cluster_centers_[:,1], s = 25, c = 'yellow', label = 'Centroids')
    plt.legend()
    plt.show()
```
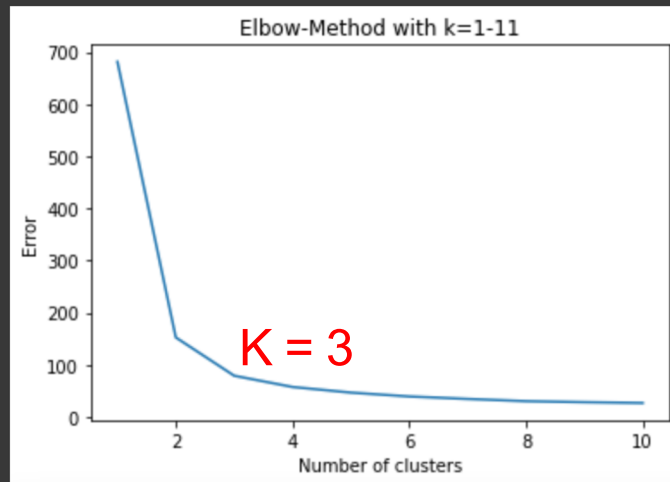
# Build a clustering model – Iris dataset

- Estimate k with elbow method

```
[8] Error =[]
    for i in range(1, 11):
        kmeans11 = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0).fit(x)
        kmeans11.fit(x)
        Error.append(kmeans11.inertia_)
    import matplotlib.pyplot as plt
    plt.plot(range(1, 11), Error)
    plt.title('Elbow-Method with k=1-11') #within cluster sum of squares
    plt.xlabel('Number of clusters')
    plt.ylabel('Error')
    plt.show()
```
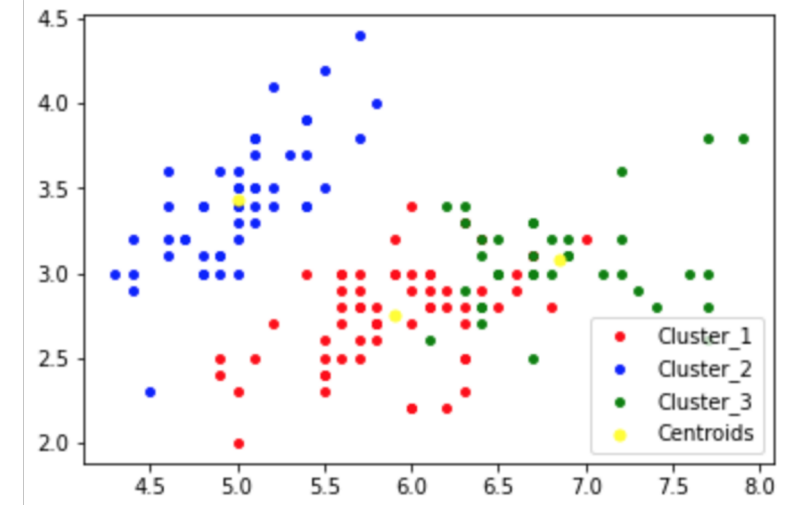


K = 3

# Build a clustering model – Iris dataset

- Get the optimal k = 3 from the elbow method. Cluster centers

```
[9] kmeans3 = KMeans(n_clusters=3, random_state=21)
    y = kmeans3.fit_predict(x)
    kmeans3.cluster_centers_

    array([[5.9016129 , 2.7483871 , 4.39354839, 1.43387097],
           [5.006     , 3.428     , 1.462     , 0.246     ],
           [6.85      , 3.07368421, 5.74210526, 2.07105263]])
```
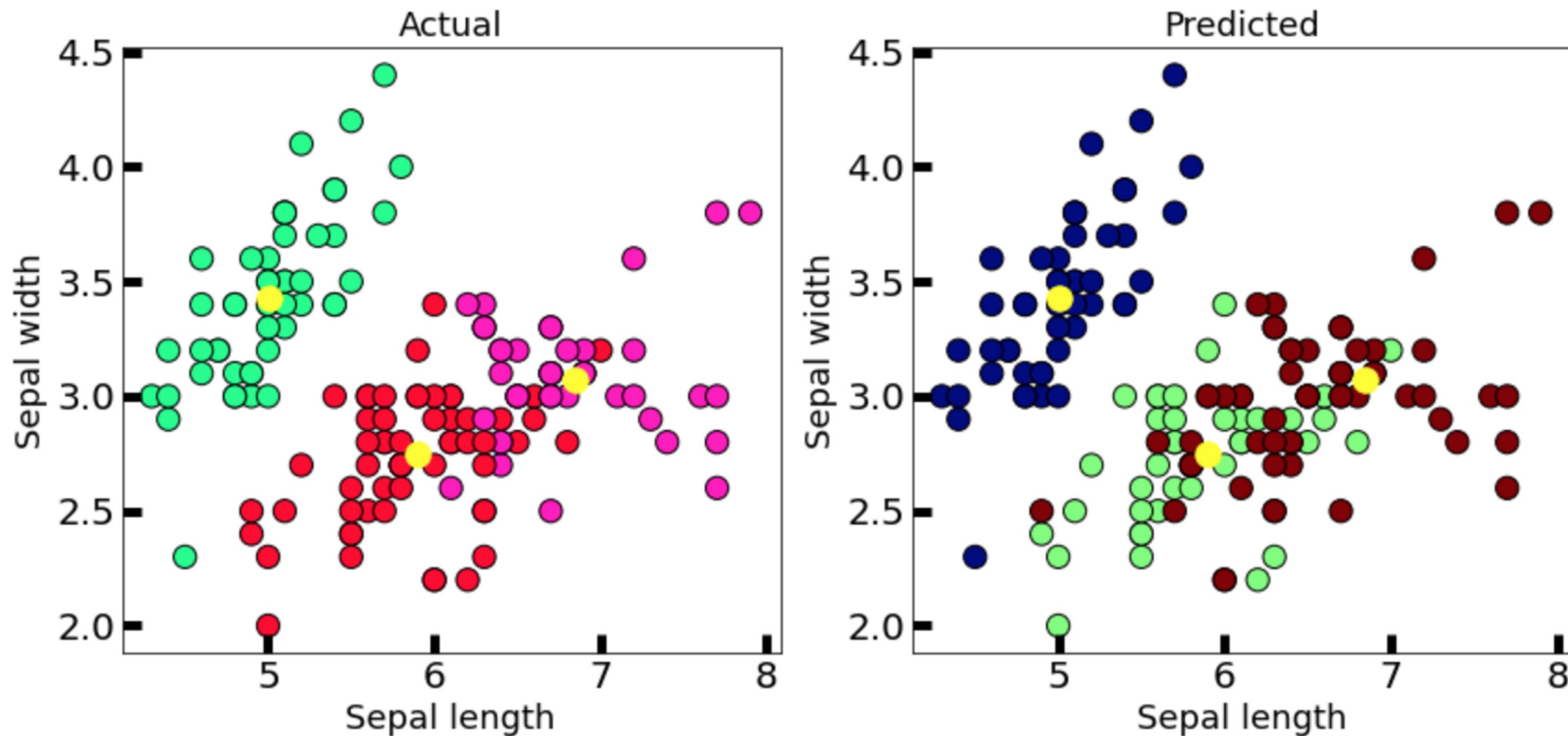


```
[11] plt.scatter(x[y == 0,0], x[y==0,1], s = 15, c= 'red', label = 'Cluster_1')
     plt.scatter(x[y == 1,0], x[y==1,1], s = 15, c= 'blue', label = 'Cluster_2')
     plt.scatter(x[y == 2,0], x[y==2,1], s = 15, c= 'green', label = 'Cluster_3')
     plt.scatter(kmeans3.cluster_centers_[:,0], kmeans3.cluster_centers_[:,1], s = 25, c = 'yellow', label = 'Centroids')
     plt.legend()
     plt.show()
```

# Build a clustering model – Iris dataset

- Compared the actual and predicted clusters

# Getting help

- HPC documentation docs.hpc.arizona.edu

- Support ticket
  https://uarizona.service-now.com/sp?id=sc_cat_item&sys_id=2983102adbd23c109627d90d689619c6&sysparm_category=84d3d1acdbc8f4109627d90d6896191f

- Office Hours – Wednesday 2-4 PM
  https://gather.town/app/dVsAprPNBVmI9NpL/hpc-office-hours

- HPC consulting
  hpc-consult@list.arizona.edu

- Visualization consulting
  vislab-consult@list.arizona.edu

- Statistics consulting
  stat-consult@list.arizona.edu