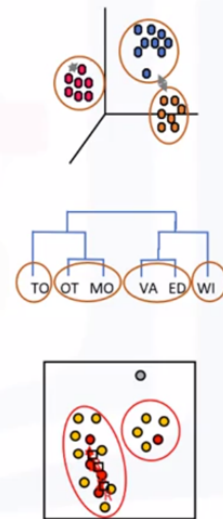


# Clustering

A cluster is a group of data points or objects in a dataset that are similar to other objects in the group, and dissimilar to datapoints in other clusters

## Clustering algorithms

- Partitioned-based Clustering
  - Relatively efficient
  - E.g. k-Means, k-Median, Fuzzy c-Means
- Hierarchical Clustering
  - Produces trees of clusters
  - E.g. Agglomerative, Divisive
- Density-based Clustering
  - Produces arbitrary shaped clusters
  - E.g. DBSCAN



### 1. Partitioned-based Clustering

#### KMeans

**iris setosa**



petal      sepal

**iris versicolor**



petal      sepal

**iris virginica**



petal      sepal

In [1]:



```
from sklearn.datasets import load_iris
iris = load_iris()
```

In [4]:



```
print(iris.keys())
```

```
dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_name', 'filename'])
```

In [5]:



```
print(iris.DESCR)
```

```
.. _iris_dataset:
```

```
Iris plants dataset  
-----
```

```
**Data Set Characteristics:**
```

```
:Number of Instances: 150 (50 in each of three classes)  
:Number of Attributes: 4 numeric, predictive attributes and the class  
:Attribute Information:  
  - sepal length in cm  
  - sepal width in cm  
  - petal length in cm  
  - petal width in cm  
  - class:  
    - Iris-Setosa  
    - Iris-Versicolour  
    - Iris-Virginica
```

```
:Summary Statistics:
```

```
=====  =====  
              Min  Max   Mean    SD   Class Correlation  
=====  =====  
sepal length:  4.3  7.9   5.84   0.83    0.7826  
sepal width:   2.0  4.4   3.05   0.43   -0.4194  
petal length:  1.0  6.9   3.76   1.76    0.9490 (high!)  
petal width:   0.1  2.5   1.20   0.76    0.9565 (high!)  
=====  =====
```

```
:Missing Attribute Values: None  
:Class Distribution: 33.3% for each of 3 classes.  
:Creator: R.A. Fisher  
:Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)  
:Date: July, 1988
```

The famous Iris database, first used by Sir R.A. Fisher. The dataset is taken from Fisher's paper. Note that it's the same as in R, but not as in the UCI Machine Learning Repository, which has two wrong data points.

This is perhaps the best known database to be found in the pattern recognition literature. Fisher's paper is a classic in the field and is referenced frequently to this day. (See Duda & Hart, for example.) The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other.

```
.. topic:: References
```

- Fisher, R.A. "The use of multiple measurements in taxonomic problems" Annual Eugenics, 7, Part II, 179-188 (1936); also in "Contributions to Mathematical Statistics" (John Wiley, NY, 1950).
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis

is.

(Q327.D83) John Wiley & Sons. ISBN 0-471-22361-1. See page 218.

- Dasarathy, B.V. (1980) "Nosing Around the Neighborhood: A New System Structure and Classification Rule for Recognition in Partially Exposed Environments". IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. PAMI-2, No. 1, 67-71.
- Gates, G.W. (1972) "The Reduced Nearest Neighbor Rule". IEEE Transactions on Information Theory, May 1972, 431-433.
- See also: 1988 MLC Proceedings, 54-64. Cheeseman et al's AUTOCLASS II conceptual clustering system finds 3 classes in the data.
- Many, many more ...

In [6]:

```
data = iris['data']
```

In [7]:

```
data = iris.data
```

In [8]:

```
data
```

Out[8]:

```
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],
       [5.4, 3.7, 1.5, 0.2],
       [4.8, 3.4, 1.6, 0.2],
       [4.8, 3. , 1.4, 0.1],
       [4.3, 3. , 1.1, 0.1],
       [5.8, 4. , 1.2, 0.2],
       [5.7, 4.4, 1.5, 0.4],
       [5.4, 3.9, 1.3, 0.4],
```

In [9]:

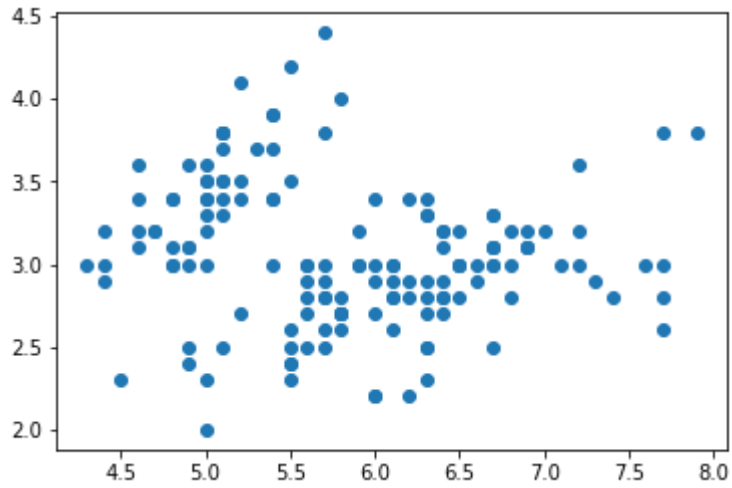
```
import matplotlib.pyplot as plt
```

In [10]:

```
plt.scatter(data[:,0],data[:,1])
```

Out[10]:

<matplotlib.collections.PathCollection at 0x24dab61bfa0>



In [11]:

```
req = data[:,0:2]  
req
```

Out[11]:

```
array([[5.1, 3.5],  
       [4.9, 3. ],  
       [4.7, 3.2],  
       [4.6, 3.1],  
       [5. , 3.6],  
       [5.4, 3.9],  
       [4.6, 3.4],  
       [5. , 3.4],  
       [4.4, 2.9],  
       [4.9, 3.1],  
       [5.4, 3.7],  
       [4.8, 3.4],  
       [4.8, 3. ],  
       [4.3, 3. ],  
       [5.8, 4. ],  
       [5.7, 4.4],  
       [5.4, 3.9],
```

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

In [12]:



```
from sklearn.model_selection import train_test_split
cl_tr,cl_ts = train_test_split(req,test_size = 0.3,random_state = 7)

from sklearn.cluster import KMeans

model = KMeans(n_clusters = 3)

model.fit(cl_tr)
```

Out[12]:

KMeans(n\_clusters=3)

In [9]:



```
y_pred = model.predict(cl_ts)
y_pred
```

Out[9]:

```
array([0, 0, 1, 0, 0, 1, 2, 0, 1, 2, 2, 0, 1, 2, 1, 0, 2, 2, 1, 1, 0, 2,
       0, 0, 2, 0, 0, 0, 0, 2, 0, 0, 1, 2, 2, 1, 1, 1, 1, 2, 2, 2, 2, 2,
       0])
```

In [10]:



```
# Import pyplot
import matplotlib.pyplot as plt

# Assign the columns of new_points: xs and ys
xs = cl_ts[:,0]
ys = cl_ts[:,1]

# Make a scatter plot of xs and ys, using labels to define the colors
plt.scatter(xs,ys,c = y_pred,alpha=0.5)

# Assign the cluster centers: centroids
centroids = model.cluster_centers_

# Assign the columns of centroids: centroids_x, centroids_y
centroids_x = centroids[:,0]
centroids_y = centroids[:,1]

# Make a scatter plot of centroids_x and centroids_y
plt.scatter(centroids_x,centroids_y,marker = 'D',s = 50)
plt.show()
```

