# Clustering in Python
# *k*-means clustering

Dr. Chun-Hao Chen

# Outline

# 1. Import Library

✓import matplotlib.pyplot as plt          # Draw figures

✓from sklearn.datasets import make_blobs          # make_blobs: generate clustering instances

✓from sklearn.cluster import KMeans          # Kmeans: k-means clustering algorithm
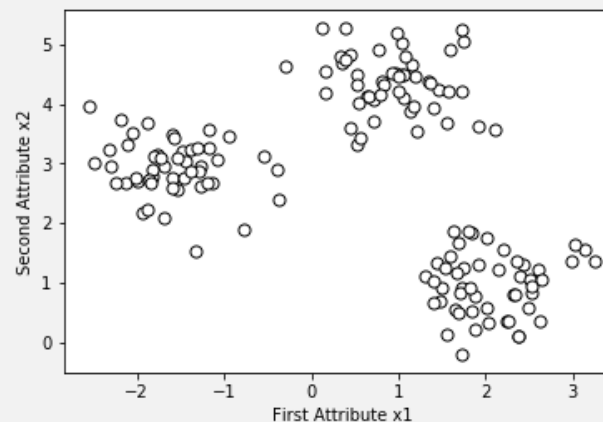
# 2. Generate Clustering Instances

```
1.   # create instances
2.   # X: attribute values
3.   # y: group label
4.   X, y = make_blobs(
5.     n_samples=150, n_features=2,
6.     centers=3, cluster_std=0.5,
7.     shuffle=True, random_state=0
8.   )
9.   print(X)
10.  print(X[:, 0])
11.  print(X[:, 1])

12.  # Draw scatter of the instances
13.  plt.scatter(
14.    X[:, 0], X[:, 1],
15.    c='white', marker='o',
16.    edgecolor='black', s=50
17.  )
18.  plt.xlabel('First Attribute x1')
19.  plt.ylabel('Second Attribute x2')
20.  plt.show()
```

```
#n_sample = 5, list X is shown as follows:
[
 [-1.3049724 3.08471943]
 [ 0.92466065 4.50908658]
 [ 1.45131429 4.22810872]
 [ 2.43578638 0.95850117]
 [ 2.12728931 1.62480041]
]
#centers = 3, list y is shown as follows:
[2 0 0 1 1]
# X[:, 0]
[-1.3049724  0.92466065 1.45131429 2.43578638 2.12728931]
# X[:, 1]
[ 3.08471943 4.50908658 4.22810872 0.95850117 1.62480041]
```
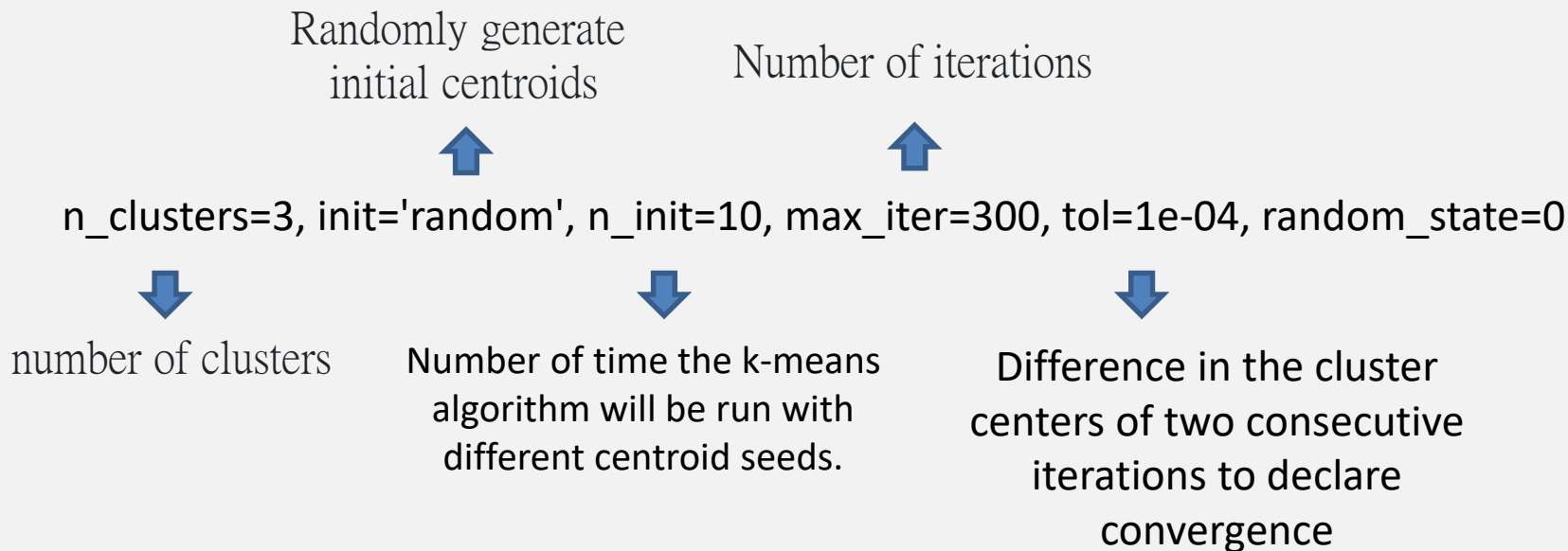


n_samples = 150
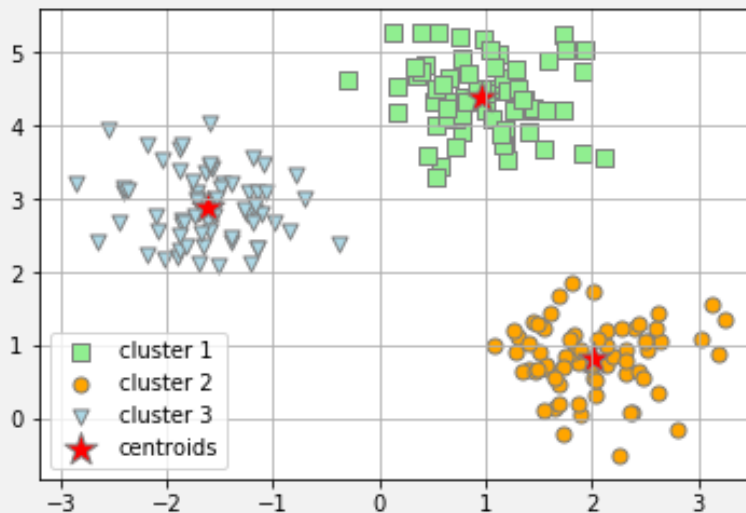centers = 3

# 3、Call KMeans() for Clustering

1. km = KMeans(
2.    n_clusters=3, init='random', n_init=10, max_iter=300, tol=1e-04, random_state=0
3. )
4. y_km = km.fit_predict(X) # Compute cluster centers and predict cluster index for each sample.
5. print(y_km) # Continue previous example: y_km = [1 2 2 0 0]

Randomly generate
initial centroids

Number of iterations

⬆                 ⬆

n_clusters=3, init='random', n_init=10, max_iter=300, tol=1e-04, random_state=0

⬇              ⬇                          ⬇

number of clusters

Number of time the k-means
algorithm will be run with
different centroid seeds.

Difference in the cluster
centers of two consecutive
iterations to declare
convergence

# 4、Show Clustering Result

1. # Draw the 3 clusters
2. plt.scatter( X[y_km == 0, 0], X[y_km == 0, 1], s=50, c='lightgreen', marker='s', edgecolor='gray', label='cluster 1')
3. plt.scatter( X[y_km == 1, 0], X[y_km == 1, 1], s=50, c='orange', marker='o', edgecolor='gray', label='cluster 2')
4. plt.scatter( X[y_km == 2, 0], X[y_km == 2, 1], s=50, c='lightblue', marker='v', edgecolor='gray', label='cluster 3')
5. # Draw the centroids
6. plt.scatter(
7. km.cluster_centers_[:, 0], km.cluster_centers_[:, 1], s=250, marker='*', c='red', edgecolor='gray', label='centroids'
8. )
9. plt.legend(scatterpoints=1)
10. plt.grid()
11. plt.show()

# 5、 Find Suitable Number of Clusters

```
1.  distortions = []
2.  for i in range(1, 15):
3.      km = KMeans( n_clusters=i, init='random', n_init=10, max_iter=300, tol=1e-04, random_state=0 )
4.      km.fit(X)
5.      distortions.append(km.inertia_)
```

➡ Inertia_ : Sum of squared distances of samples to their closest cluster center.

```
6.  # Draw figures
7.  plt.plot(range(1, 15), distortions, marker='o')
8.  plt.xlabel('Number of clusters')
9.  plt.ylabel('Distortion')
10. plt.show()
```

k = 3 is suitable