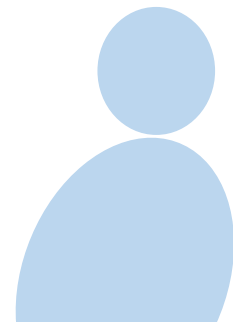


Chapter 6

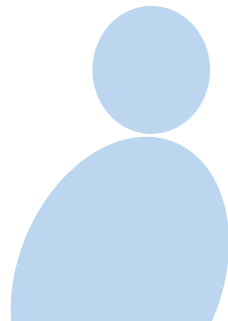
K Means Clustering



Clustering

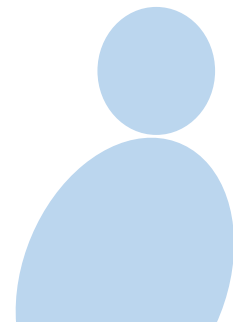
- Unsupervised learning

- Requires data, but _____



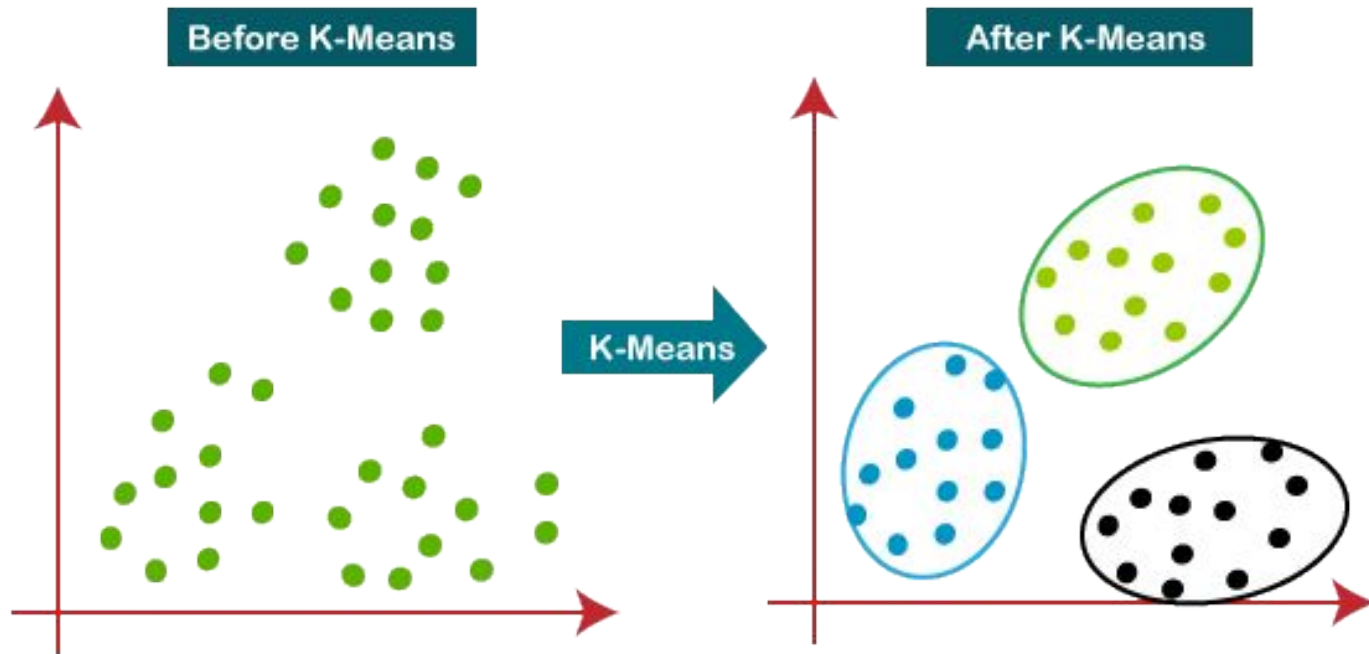
Clustering

- Unsupervised learning
- Requires data, but **NO** labels



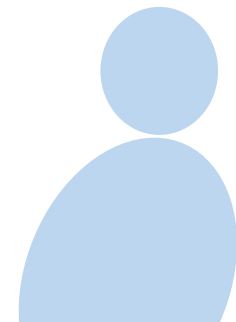
Clustering

- Basic idea: Group together *similar instances*



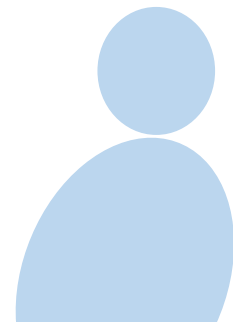
Clustering

- Image segmentation
 - Goal: Partition an image into perceptually similar regions



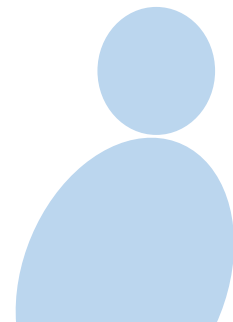
Other examples

- Group emails or search results
- Customer shopping patterns



K-Means Algorithm

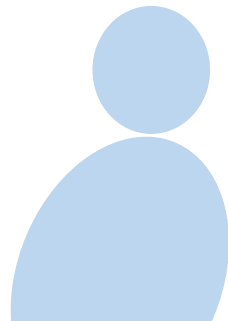
- K –means is an iterative algorithm
- Consider the *dataset* of n –dimensional vectors



K-Means Algorithm

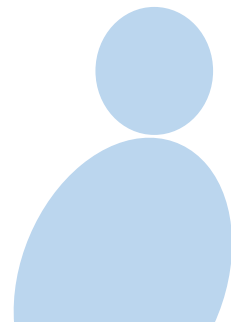
- K –means is an iterative algorithm
- Consider the *dataset* of n –dimensional vectors

$$\bar{\mathbf{x}}(1), \bar{\mathbf{x}}(2), \dots, \bar{\mathbf{x}}(M)$$



Clusters

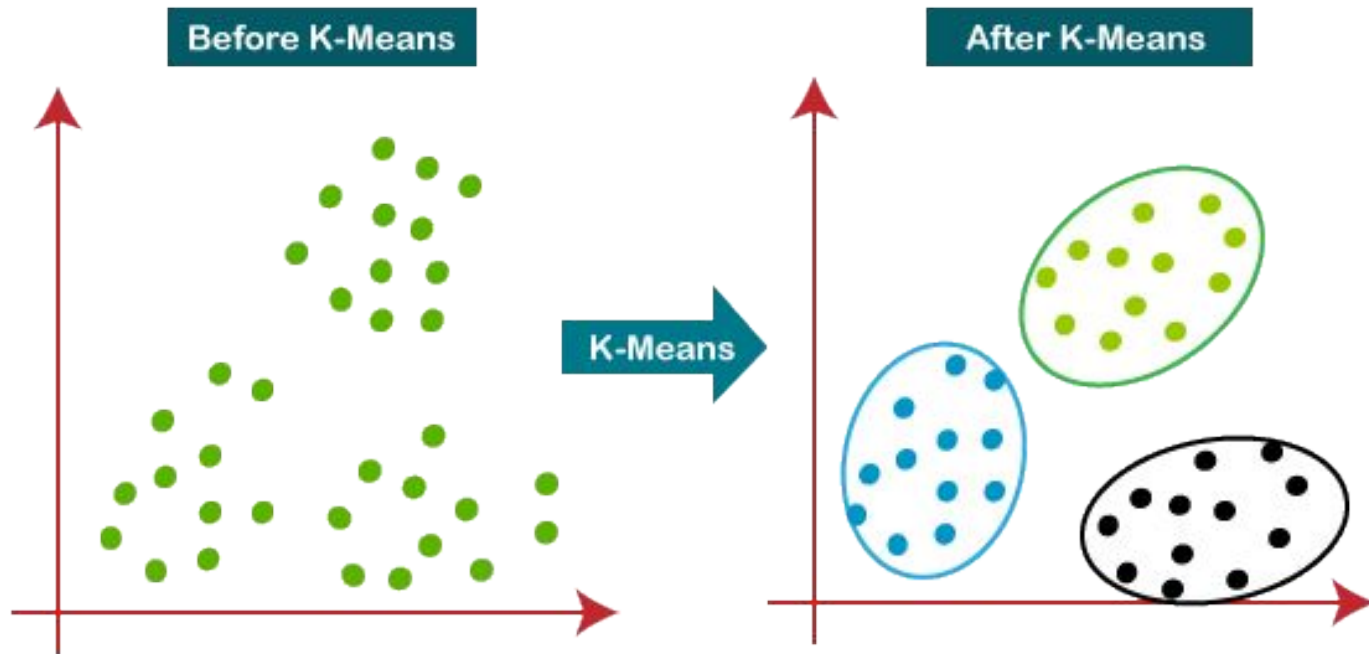
- Organize the data into K *clusters*



Clusters

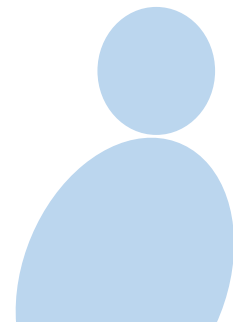
- Organize the data into K *clusters*

$$\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K$$



Clusters

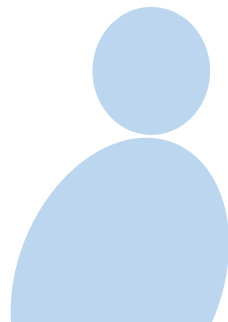
- The centroids for the clusters are



Clusters

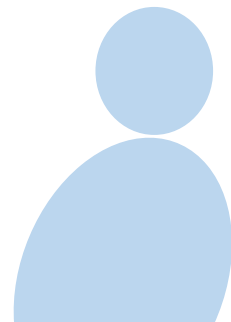
- The centroids for the clusters are

$$\bar{\mu}_1, \bar{\mu}_2, \dots, \bar{\mu}_K$$

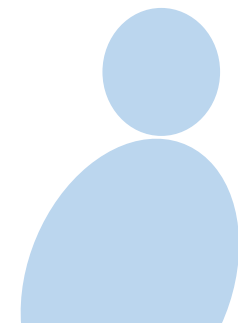


Cluster assignment

- Let $\alpha_i(j)$ denote the *cluster assignment indicator*



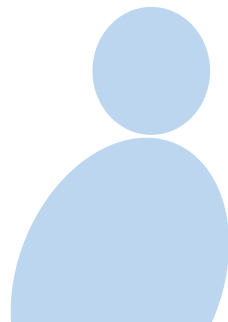
Cluster assignment



Cluster assignment

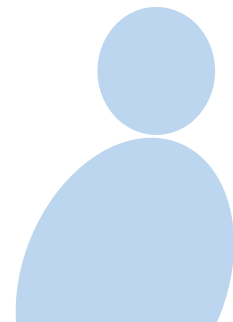
- Let $\alpha_i(j)$ denote the *cluster assignment indicator*

$$\alpha_i(j) = \begin{cases} 1 & \bar{\mathbf{x}}(j) \in \mathcal{C}_i \\ 0 & \bar{\mathbf{x}}(j) \notin \mathcal{C}_i \end{cases}$$



K-Means procedure

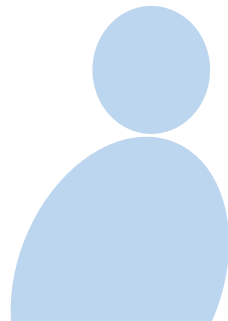
- *Initialize centroids* randomly
- $\bar{\mu}_i^{(l-1)}$ denotes centroids in iteration $l - 1$



K-Means procedure

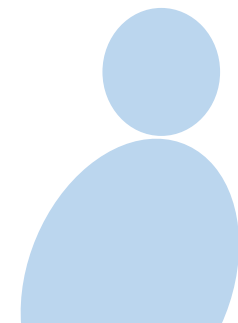
- *Initialize centroids* randomly
 $\bar{\mu}_1^{(0)}, \bar{\mu}_2^{(0)}, \dots, \bar{\mu}_K^{(0)}$

- $\bar{\mu}_i^{(l-1)}$ denotes centroids in iteration $l - 1$



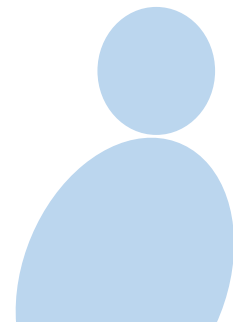
Cluster determination

- Assign $\bar{\mathbf{x}}(j)$ to the cluster \tilde{i} with *closest centroid* $\bar{\boldsymbol{\mu}}_{\tilde{i}}^{(l-1)}$



Cluster determination

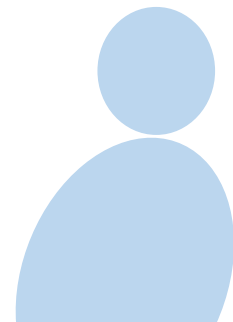
- This is *minimized* when $\alpha_{\tilde{i}}^{(l)}(j) = 1$,
where



Cluster determination

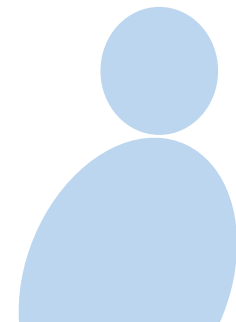
- This is *minimized* when $\alpha_{\tilde{i}}^{(l)}(j) = 1$,
where

$$\tilde{i} = \arg \min \left\| \bar{\mathbf{x}}(j) - \bar{\boldsymbol{\mu}}_i^{(l-1)} \right\|^2$$



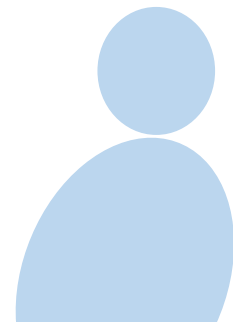
Centroid determination

- The centroid of cluster i can be determined as below
- The **average** of all points assigned to cluster i in iteration l



Centroid determination

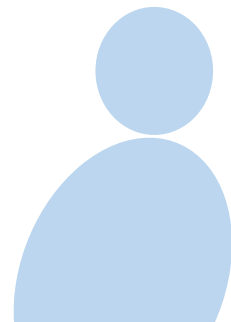
- This can be mathematically expressed as



Centroid determination

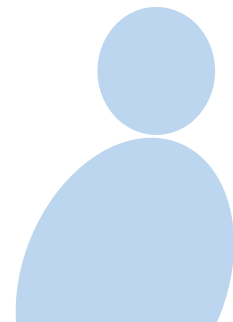
- This can be mathematically expressed as

$$\begin{aligned}\bar{\mu}_i^{(l)} &= \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{\mathbf{x}}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)} \\ &= \frac{\sum_{j: \bar{\mathbf{x}}(j) \in \mathcal{C}_i} \bar{\mathbf{x}}(j)}{\sum_{j: \bar{\mathbf{x}}(j) \in \mathcal{C}_i} 1}\end{aligned}$$



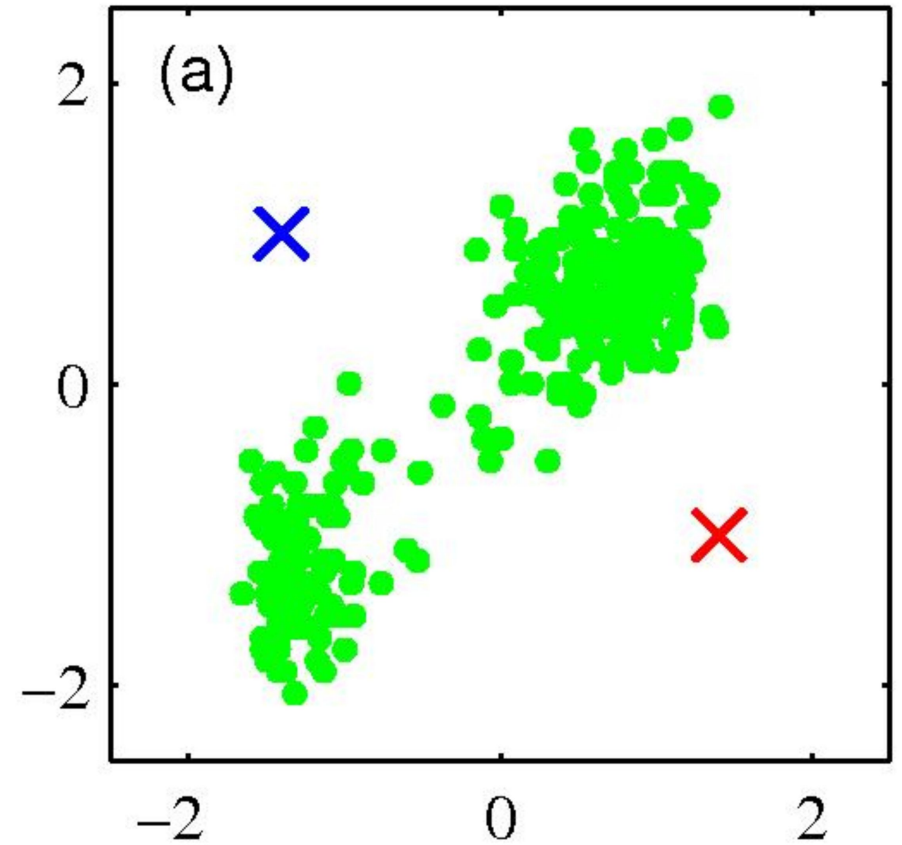
Convergence

- The cluster and centroid computation steps are repeated until convergence
- i.e., when cluster assignments do NOT change



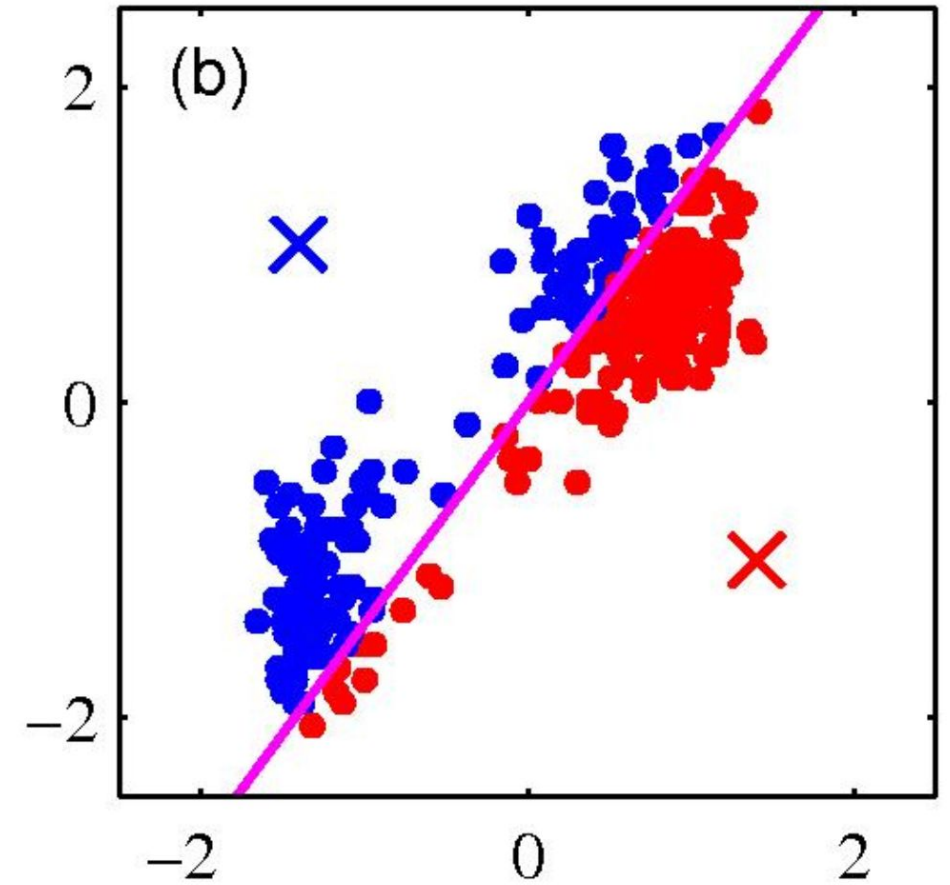
K-Means Example

- Pick K *random points* as cluster *centroids*
- Shown here for $K = 2$



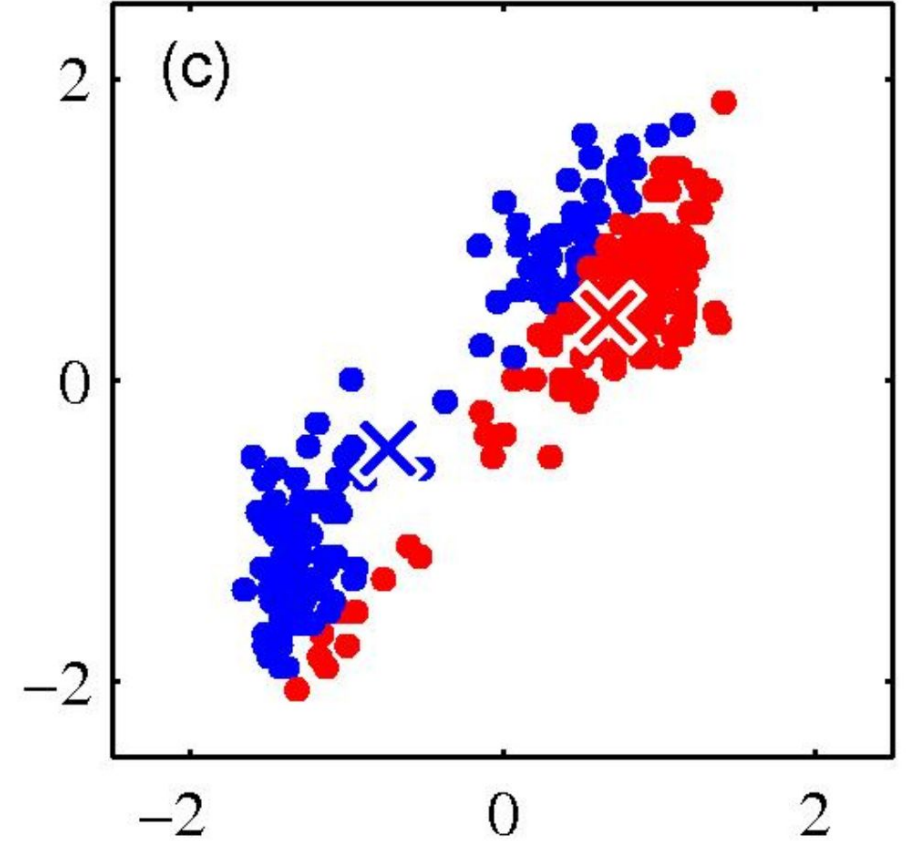
Iteration 1

- Assign data points to *closest centroid*



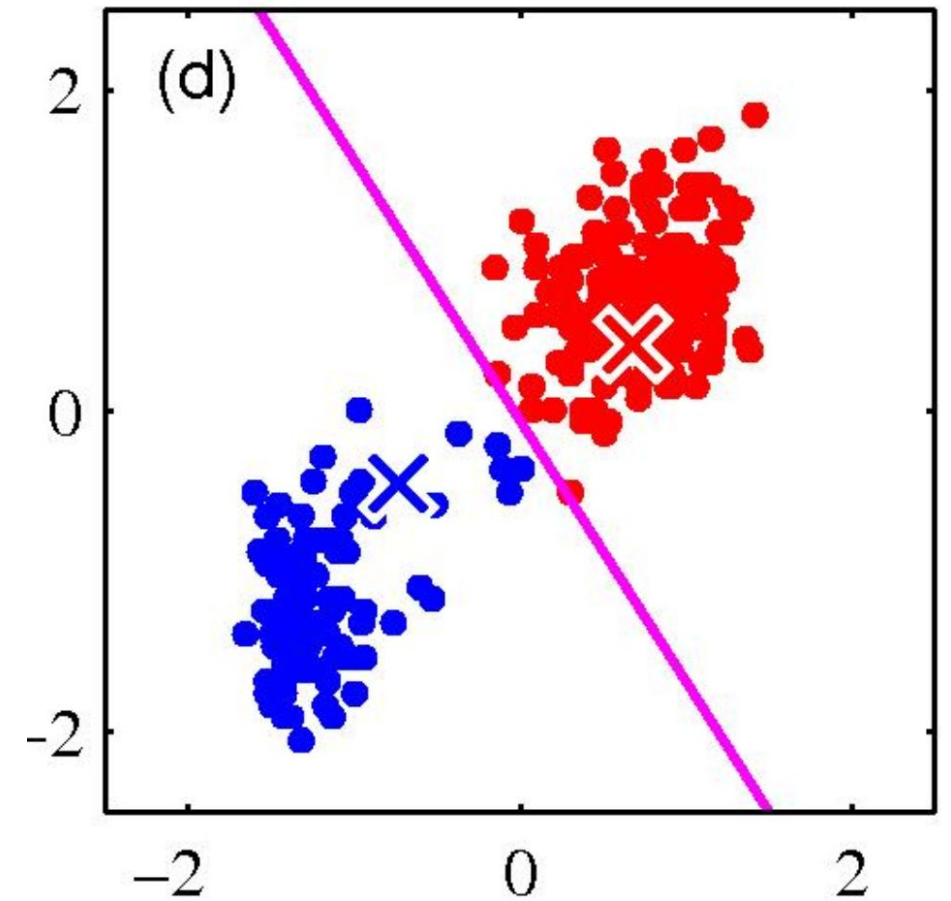
Iteration 1

- Change each *centroid* to the average of the assigned points



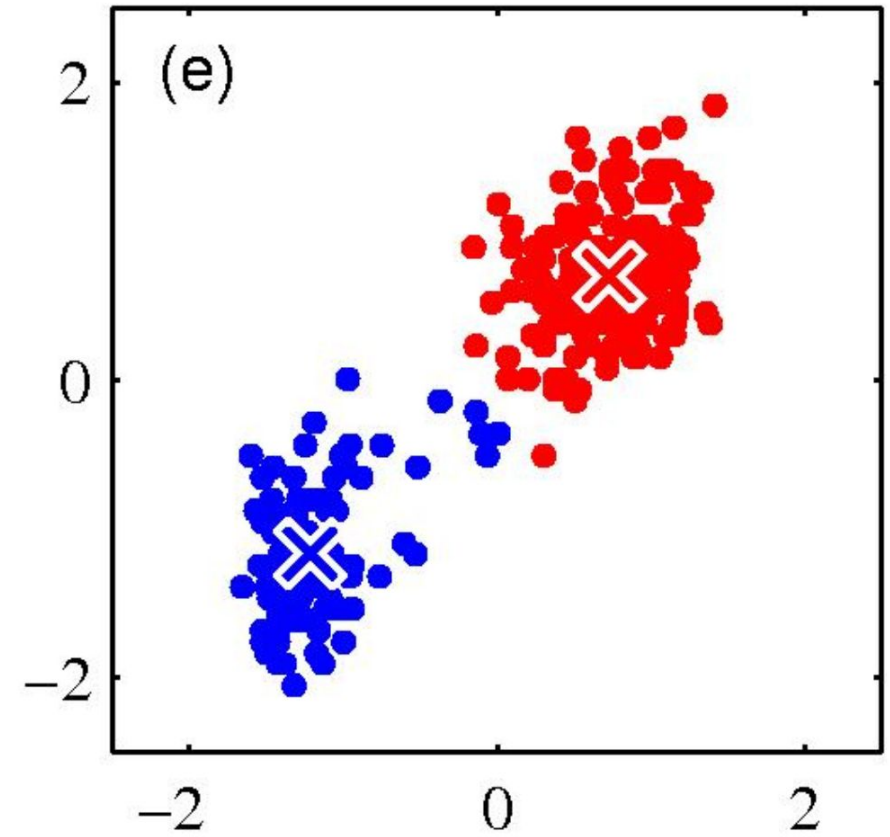
Iteration 2

- Assign data points to *closest centroid*



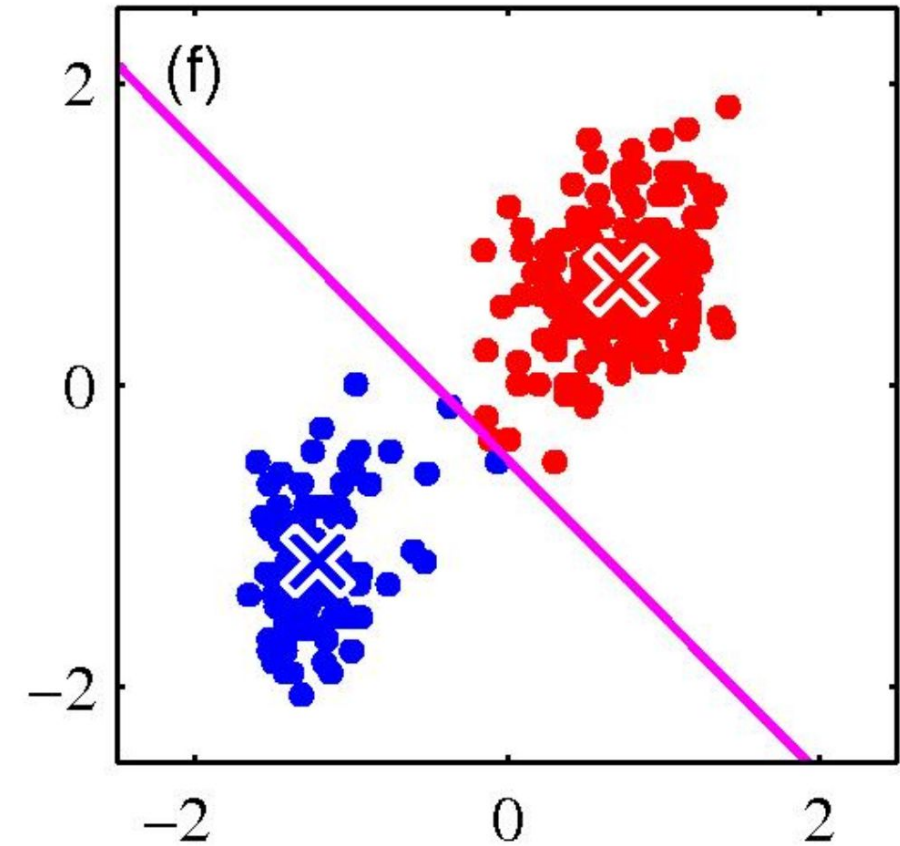
Iteration 2

- Change each *centroid* to the average of the assigned points



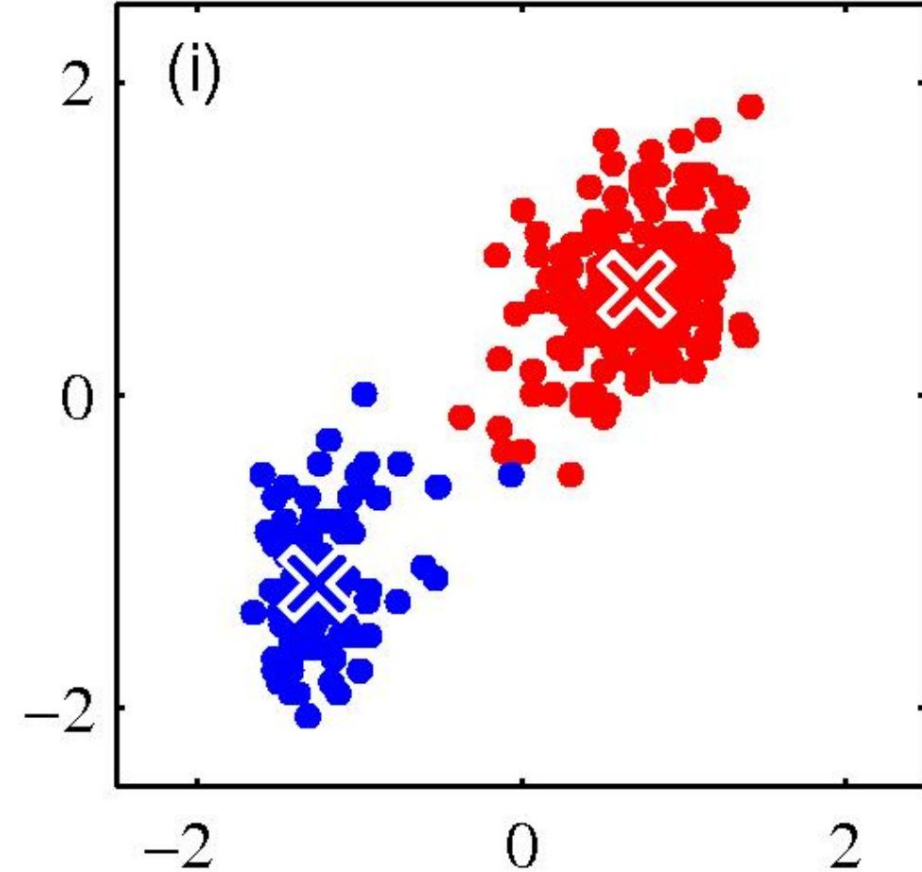
Iteration 3

- Assign data points to *closest centroid*

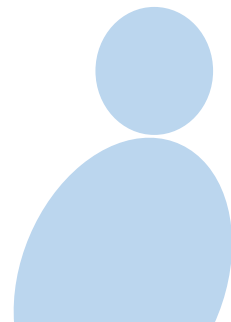


Iteration 3

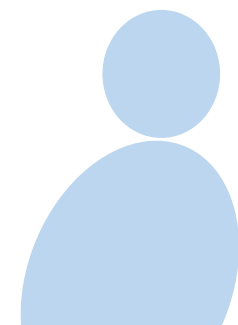
- Change each *centroid* to the average of the assigned points
- *Convergence achieved!*



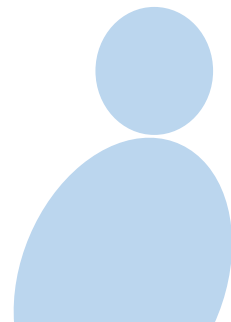
```
1 import matplotlib.pyplot as plt
2 from sklearn.cluster import KMeans
3 from sklearn.datasets import make_blobs
4
```



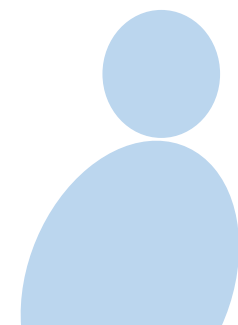
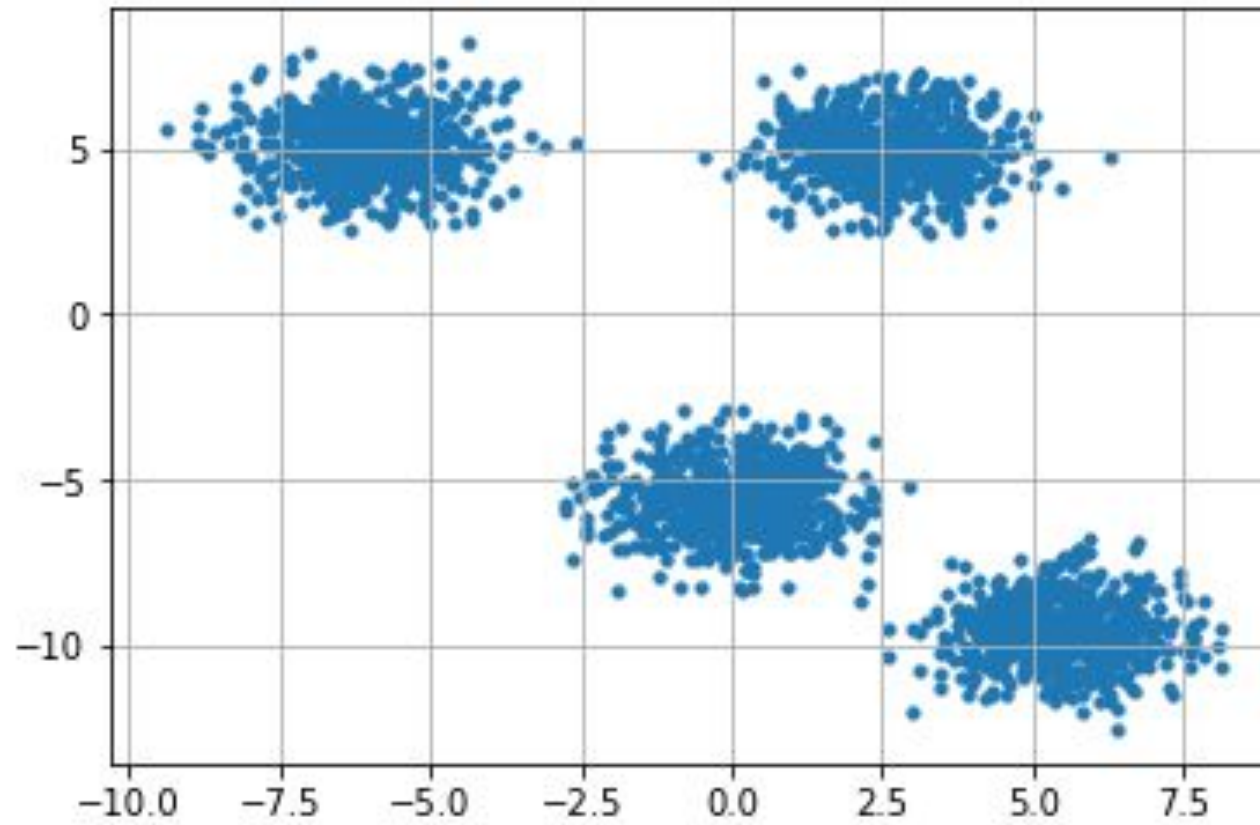

```
5  
6  
7 X, y = make_blobs(n_samples=2500, centers=4, n_features=2, random_state = 10)  
8
```



```
9 plt.figure()
10 plt.scatter(X[:, 0], X[:, 1], c=None, cmap='jet', s=10)
11 plt.suptitle('Original Data')
12 plt.grid(1, which='both')
13 plt.axis('tight')
14 plt.show()
15
```



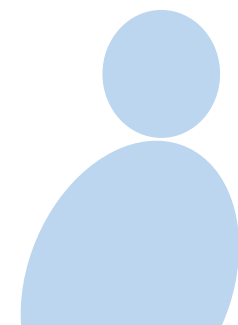
Original Data



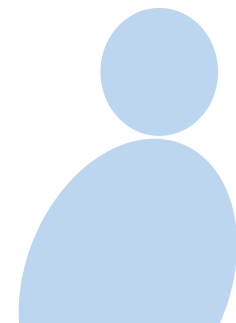
15

```
16 y_pred = KMeans(n_clusters=4, random_state=0).fit_predict(X)
```

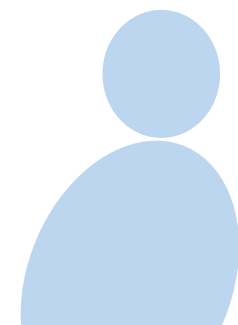
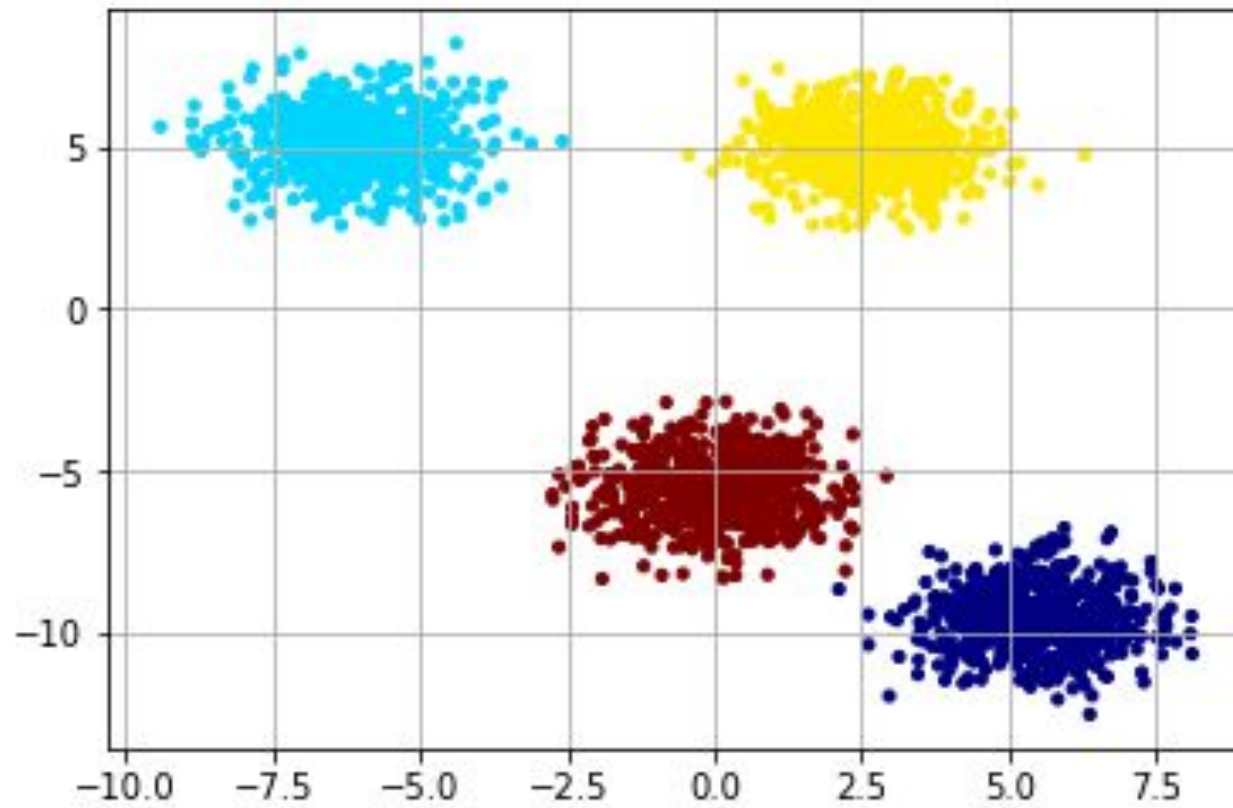
17



```
18 plt.figure()
19 plt.scatter(X[:, 0], X[:, 1], c=y_pred, cmap='jet', s=10)
20 plt.suptitle('K-Means Clusters')
21 plt.grid(1, which='both')
22 plt.axis('tight')
23 plt.show()
24
```

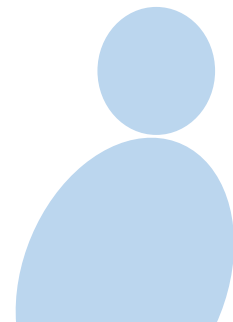


K-Means Clusters

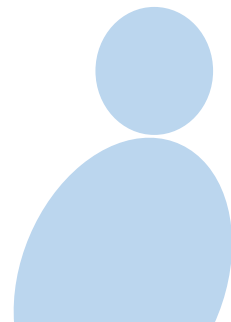


Number of clusters

- How to choose number of clusters?
- Use the Sum Square Error (SSE) as a metric
- Also termed **inertia**

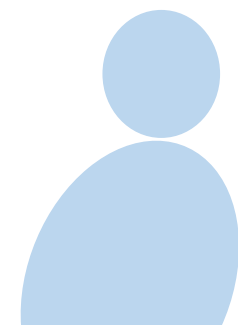


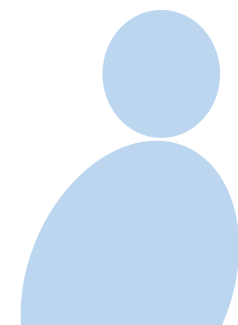
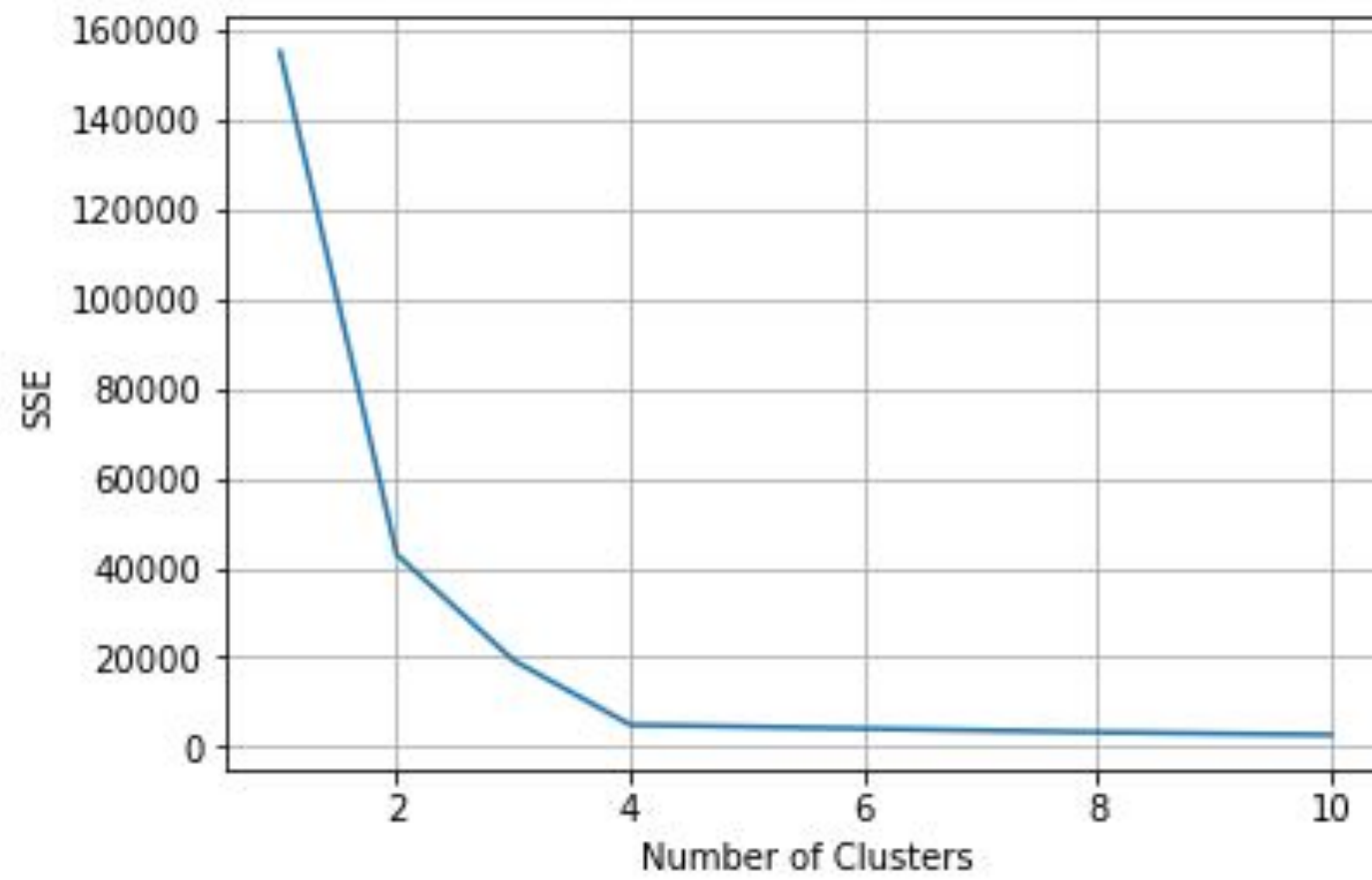
```
25 #Sum squared error
26 SSE = []
27
28 for k in range(1, 11):
29     kmeans = KMeans(n_clusters=k, random_state=0)
30     kmeans.fit(X)
31     SSE.append(kmeans.inertia_)
32
```




```
32  
33 plt.plot(range(1, 11), SSE)  
34 plt.xlabel("Number of Clusters")  
35 plt.ylabel("SSE")  
36 plt.grid(1, which='both')  
37 plt.axis('tight')  
38 plt.show()
```

```
39  
40
```





Instructors may use this white area (14.5 cm / 25.4 cm) for the text. Three options provided below for the font size.

Font: Avenir (Book), Size: 32, Colour: Dark Grey

Font: Avenir (Book), Size: 28, Colour: Dark Grey

Font: Avenir (Book), Size: 24, Colour: Dark Grey

Do not use the space below.

