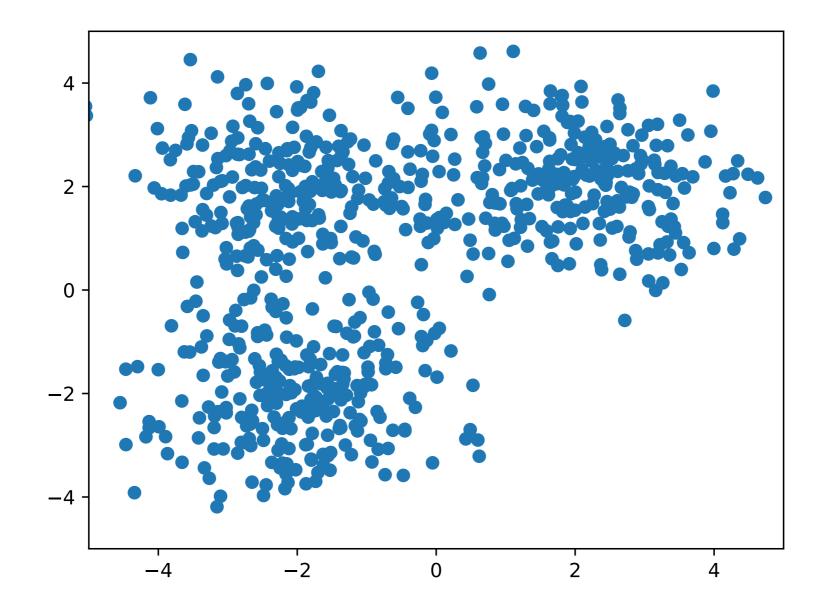
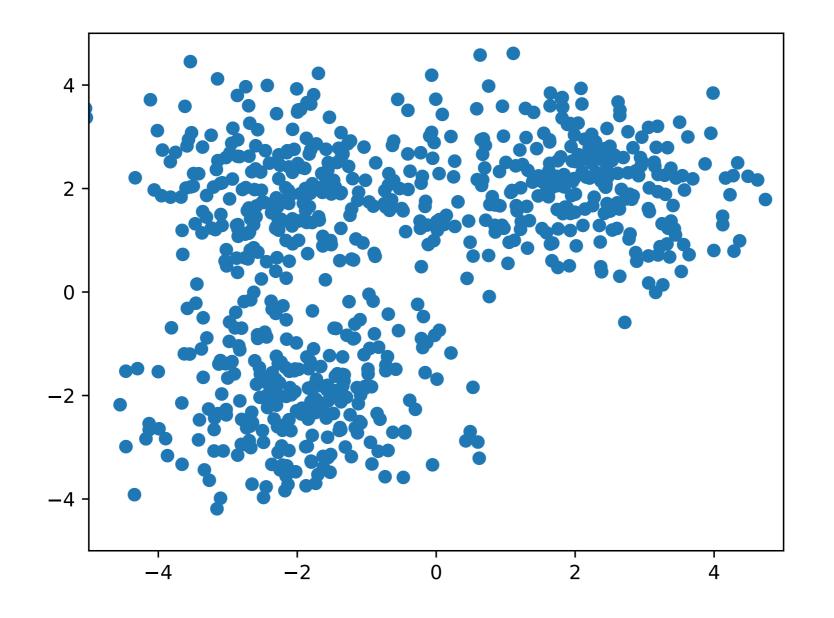
CS 453X: Class 17

Jacob Whitehill

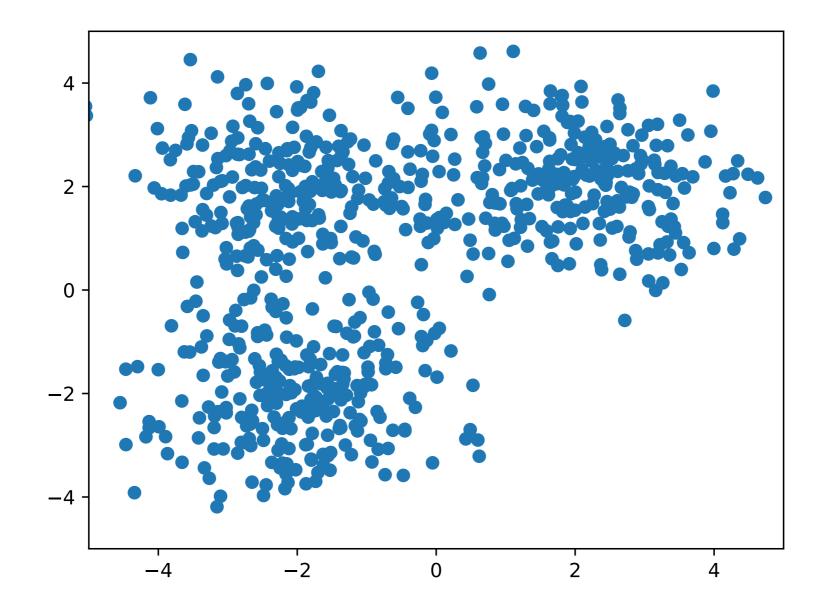
- What do you see in this figure?
 - 1.A candelabra.



- What do you see in this figure?
 - 2.A young woman frowning.

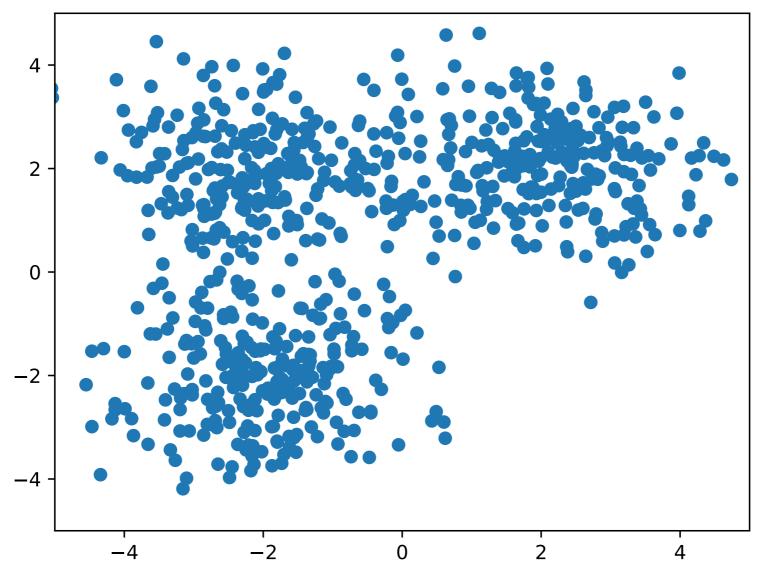


- What do you see in this figure?
 - 3. Three somewhat distinct clusters of data points.

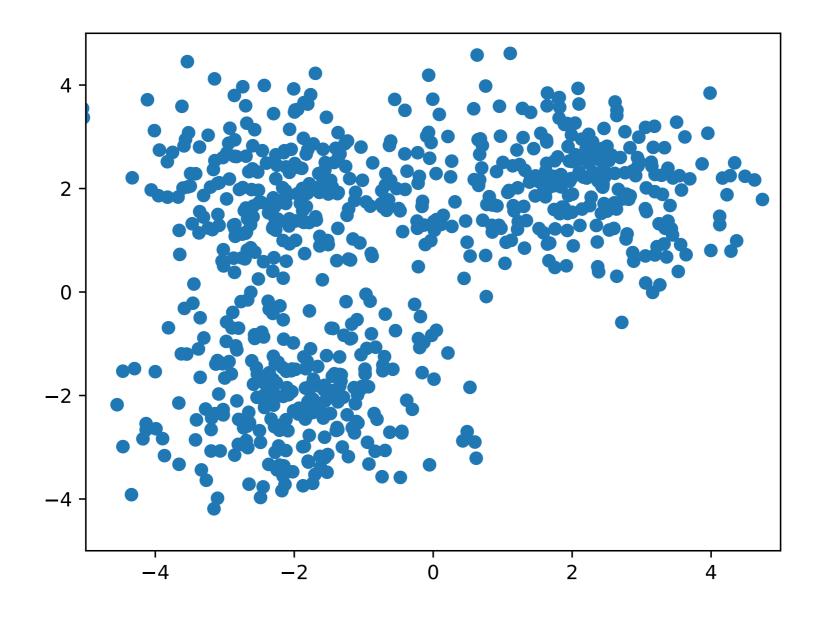


- What do you see in this figure?
 - 3. Three somewhat distinct clusters of data points.

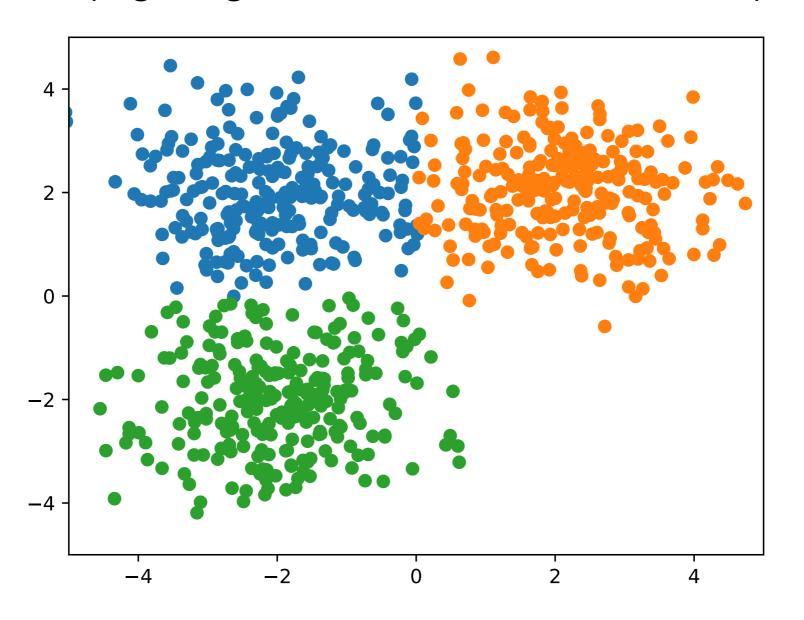
Bingo.



 Intuitively, we can define clustering as putting data into groups such that: data within each group are more similar than data between groups.

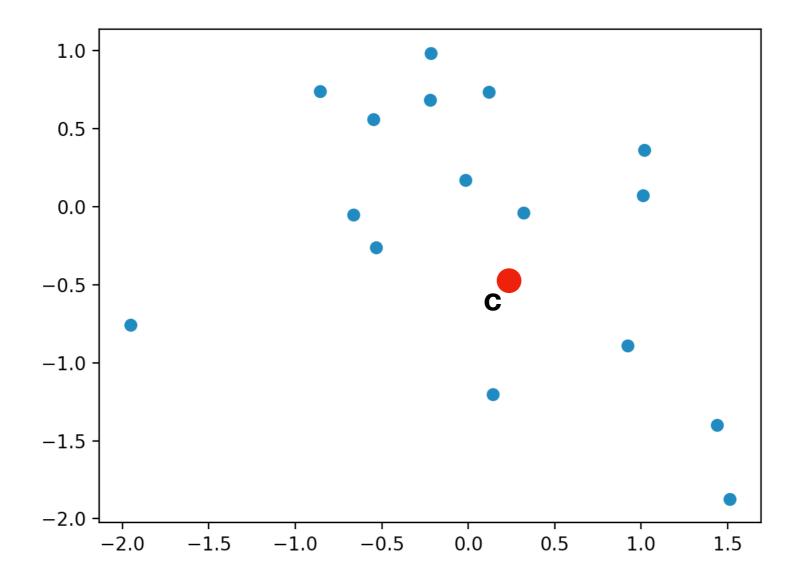


- Wouldn't it be nice to be able to cluster data automatically?
 - Maybe the clusters align with certain natural structure in the data (e.g., digit classes in MNIST dataset)?



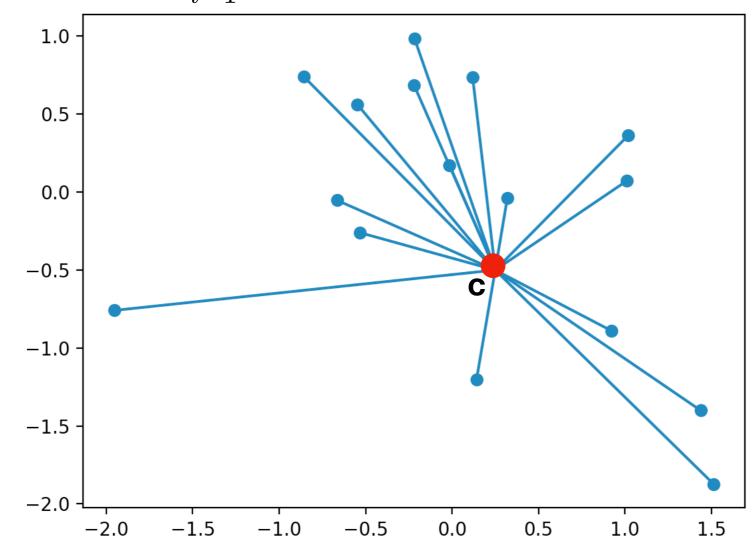
Show demo.

Consider a set of n data points { x⁽ⁱ⁾ }, and another point c:



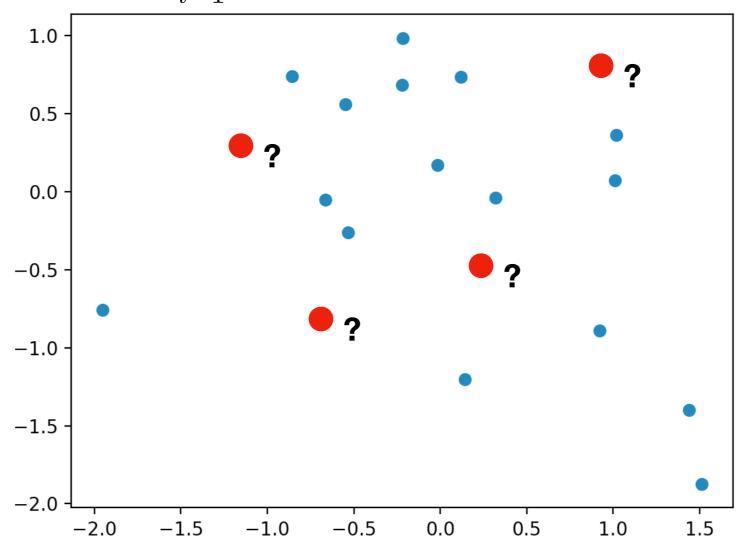
 Let's define a cost function f_{SSD} as the sum of squared distances between each data point x⁽ⁱ⁾ and c.

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^2$$



• Which point **c** minimizes f_{SSD} for $\{ \mathbf{x}^{(i)} \}$?

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^2$$



$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^2$$

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$\nabla_{\mathbf{c}} f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} \nabla_{\mathbf{c}} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$\nabla_{\mathbf{c}} f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} \nabla_{\mathbf{c}} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$= \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})$$

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$\nabla_{\mathbf{c}} f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} \nabla_{\mathbf{c}} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$= \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})$$

$$= \sum_{i=1}^{n} \mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$\nabla_{\mathbf{c}} f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} \nabla_{\mathbf{c}} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$= \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})$$

$$= \sum_{i=1}^{n} \mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

$$= n\mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)} = 0$$

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$\nabla_{\mathbf{c}} f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} \nabla_{\mathbf{c}} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$= \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})$$

$$= \sum_{i=1}^{n} \mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

$$= n\mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)} = 0$$

$$n\mathbf{c} = \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

$$f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$\nabla_{\mathbf{c}} f_{\text{SSD}}(\mathbf{c}) = \frac{1}{2} \sum_{i=1}^{n} \nabla_{\mathbf{c}} (\mathbf{c} - \mathbf{x}^{(i)})^{2}$$

$$= \sum_{i=1}^{n} (\mathbf{c} - \mathbf{x}^{(i)})$$

$$= \sum_{i=1}^{n} \mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

$$= n\mathbf{c} - \sum_{i=1}^{n} \mathbf{x}^{(i)} = 0$$

$$n\mathbf{c} = \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

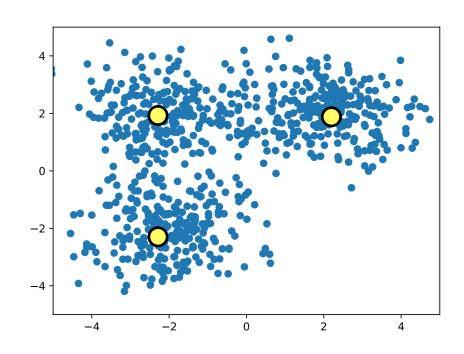
$$\mathbf{c} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

- In other words, the point c that minimizes f_{SSD} is the mean of the n points { x⁽ⁱ⁾ }.
 - Every other point c'≠ c must have a higher f_{SSD}.

- Probably the simplest and most commonly used clustering algorithm is called *k*-means.
- It partitions a set of n data $\{ \mathbf{x}^{(i)} \}$ into k clusters.

- Chicken-and-the-egg problem:
 - If we knew the *mean* $\mathbf{c}^{(j)}$ of each cluster j, we could assign each data point \mathbf{x} to the closest cluster center.

$$a(\mathbf{x}) = \arg\min_{j} \left(\mathbf{x} - \mathbf{c}^{(j)}\right)^{2}$$

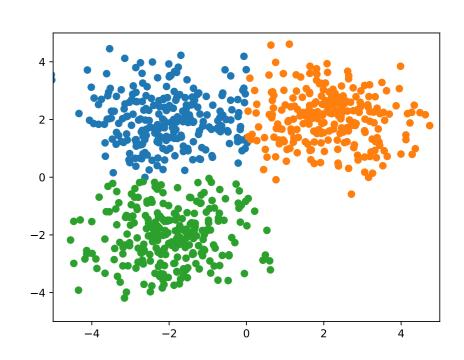


- Chicken-and-the-egg problem:
 - If we knew the *mean* $\mathbf{c}^{(j)}$ of each cluster j, we could assign each data point \mathbf{x} to the closest cluster center.

$$a(\mathbf{x}) = \arg\min_{j} \left(\mathbf{x} - \mathbf{c}^{(j)}\right)^{2}$$

 If we knew which data were in which cluster, we could compute the mean of each cluster.

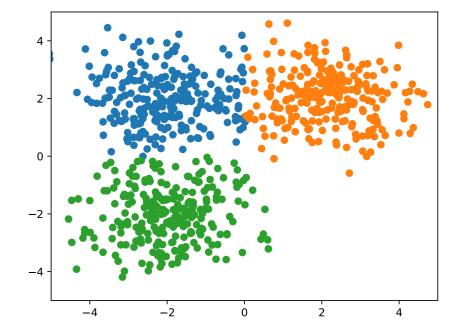
$$\mathbf{c}_j = \frac{1}{n_j} \sum_{i: a(\mathbf{x}^{(i)}) = j} \mathbf{x}^{(i)}$$



- Chicken-and-the-egg problem:
 - If we knew the *mean* $\mathbf{c}^{(j)}$ of each cluster j, we could assign each data point \mathbf{x} to the closest cluster center.

$$a(\mathbf{x}) = \arg\min_{j} \left(\mathbf{x} - \mathbf{c}^{(j)}\right)^{2}$$

 If we knew which data were in which cluster, we could compute the mean of each cluster.



$$\mathbf{c}_j = rac{1}{n_j}\sum_{i:a(\mathbf{x}^{(i)})=j}\mathbf{x}^{(i)}$$
 # data in cluster j

• The *k*-means algorithm seeks to optimize the assignment of data points to clusters so as to minimize:

$$\sum_{j=1}^{k} \sum_{i:a(\mathbf{x}^{(i)})=j} \left(\mathbf{x}^{(i)} - \mathbf{c}^{(j)}\right)^{2}$$

- Algorithm:
 - 1.Randomly assign the data points to clusters.
 - 2. Repeat until the cost does not change:
 - A.Compute $\mathbf{c}^{(j)}$ of each cluster j as the mean of the points assigned to it.
 - B.Assign each point $\mathbf{x}^{(i)}$ to the nearest cluster mean $\mathbf{c}^{(i)}$.

 Why is the "Repeat until the cost does not change" loop guaranteed to converge?

- Why is the "Repeat until the cost does not change" loop guaranteed to converge? $\sum_{j=1}^k \sum_{i: a(\mathbf{x}^{(i)}) = j} \left(\mathbf{x}^{(i)} \mathbf{c}^{(j)}\right)^2$
 - The *k*-means cost function has a *lower bound* (0) since the sum can never be negative.
 - Each step within the loop can only lower the cost:

- Why is the "Repeat until the cost does not change" loop guaranteed to converge? $\sum_{i=1}^{k} \sum_{i: a(\mathbf{x}^{(i)}) = i} (\mathbf{x}^{(i)} \mathbf{c}^{(j)})^2$
 - The *k*-means cost function has a *lower bound* (0) since the sum can never be negative.
 - Each step within the loop can only lower the cost:
 - A.Compute $\mathbf{c}^{(j)}$ of each cluster j as the mean of the points assigned to it.

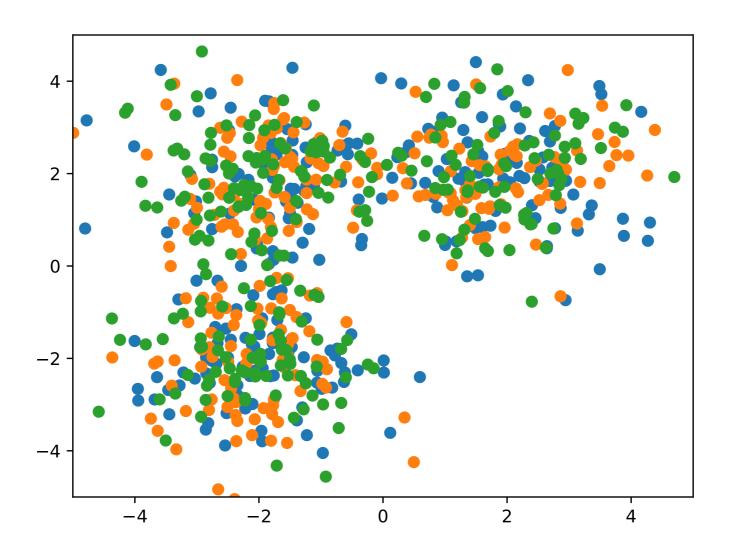
As shown earlier, the mean of the data $\{x^{(i)}\}$ minimizes f_{SSD} .

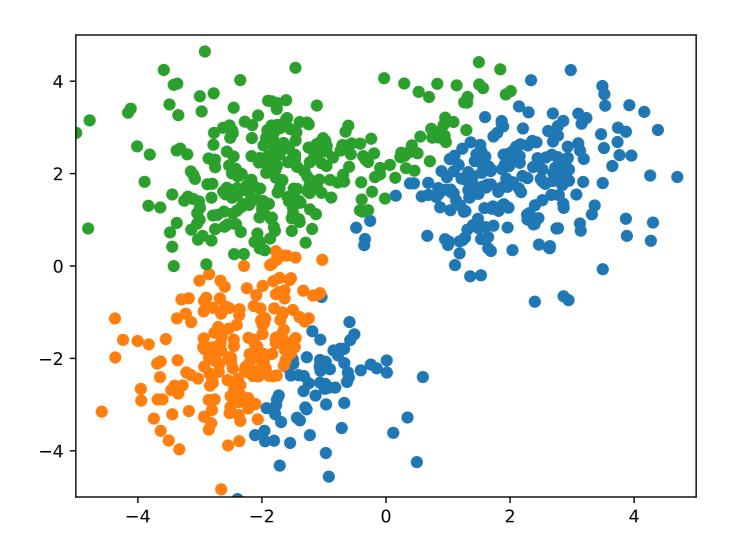
- Why is the "Repeat until the cost does not change" loop guaranteed to converge? $\sum_{j=1}^k \sum_{i: a(\mathbf{x}^{(i)}) = j} \left(\mathbf{x}^{(i)} \mathbf{c}^{(j)}\right)^2$
 - The *k*-means cost function has a *lower bound* (0) since the sum can never be negative.
 - Each step within the loop can only lower the cost:
 - A.Compute $\mathbf{c}^{(j)}$ of each cluster j as the mean of the points assigned to it.

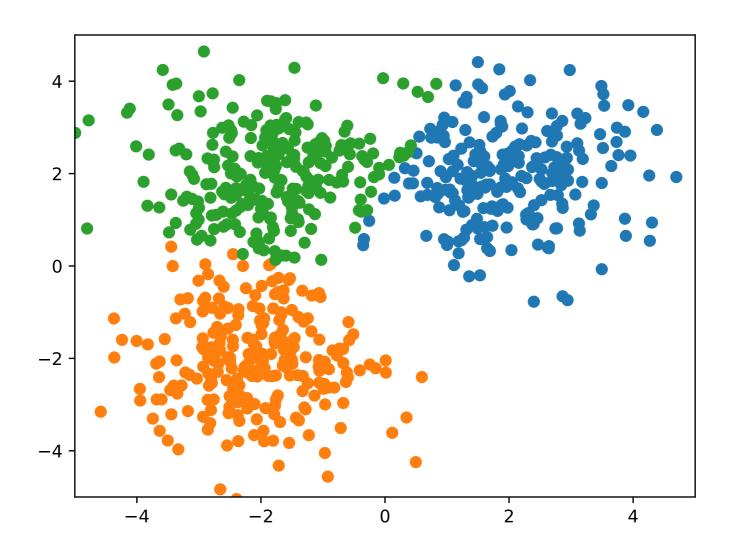
As shown earlier, the mean of the data $\{x^{(i)}\}$ minimizes f_{SSD} .

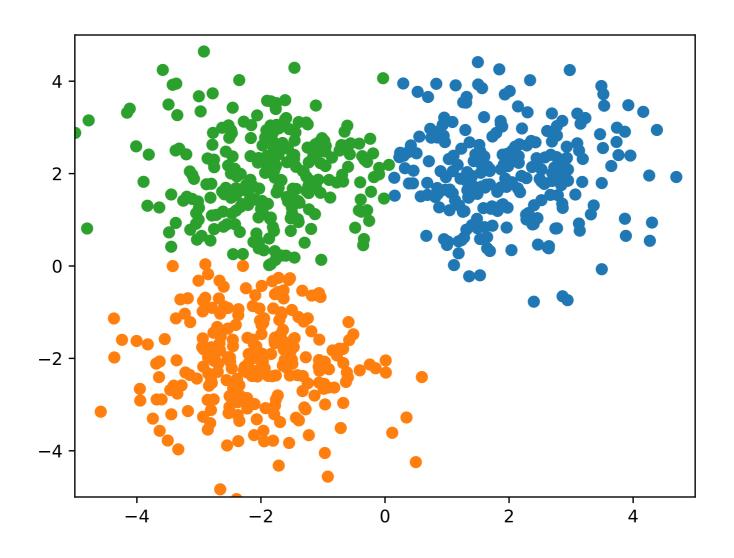
B.Assign each point $\mathbf{x}^{(i)}$ to the nearest cluster mean $\mathbf{c}^{(i)}$.

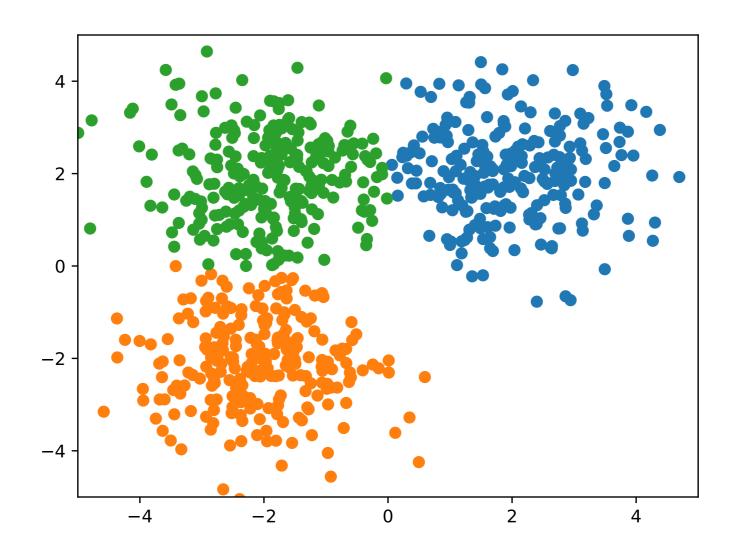
Assigning $x^{(i)}$ to the closest cluster mean $c^{(i)}$ can only decrease the cost due to $x^{(i)}$.

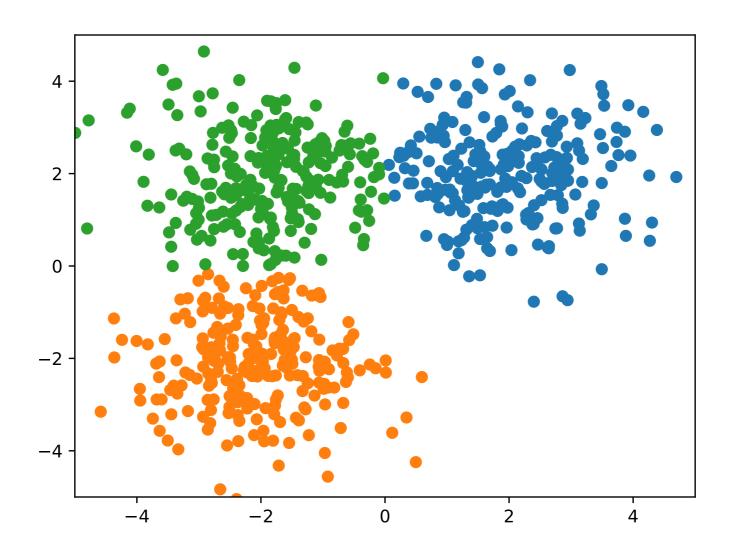






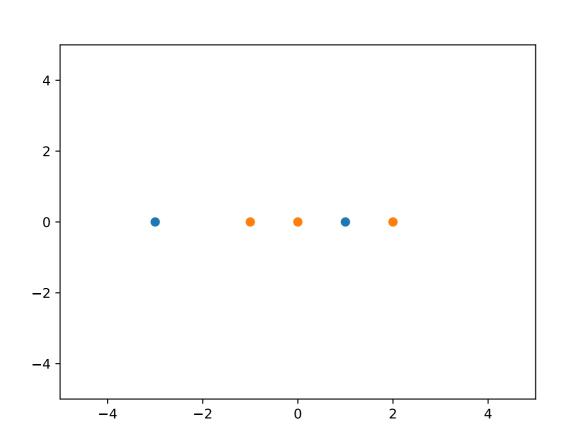






Exercise: *k*-means on 1-D data

- Suppose the (1-D) data consist of { -3, -1, 0, 1, 2 }.
- Suppose we initialize the clusters as:



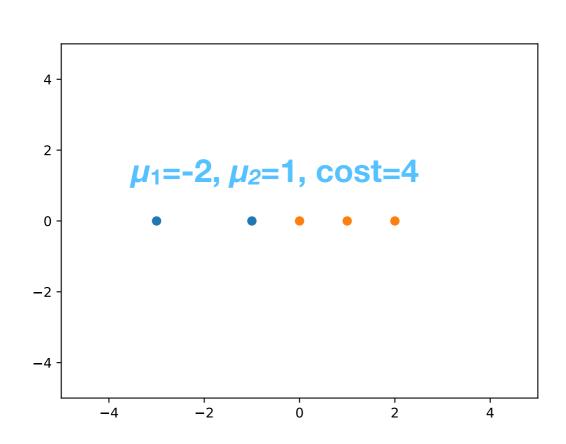
$$\sum_{j=1}^{k} \sum_{i: a(\mathbf{x}^{(i)}) = j} \left(\mathbf{x}^{(i)} - \mathbf{c}^{(j)} \right)^{2}$$

- Repeat until the cost does not change:
 - A.Compute **c**(*i*) of each cluster *j* as the mean of the points assigned to it.
 - B.Assign each point $\mathbf{x}^{(j)}$ to the nearest cluster mean $\mathbf{c}^{(j)}$.

What will k-means output (cluster assignments & means)?

Exercise: *k*-means on 1-D data

- Suppose the (1-D) data consist of { -3, -1, 0, 1, 2 }.
- Suppose we initialize the clusters as:



$$\sum_{j=1}^{k} \sum_{i: a(\mathbf{x}^{(i)}) = j} \left(\mathbf{x}^{(i)} - \mathbf{c}^{(j)} \right)^{2}$$

- Repeat until the cost does not change:
 - A.Compute **c**(*i*) of each cluster *j* as the mean of the points assigned to it.
 - B.Assign each point $\mathbf{x}^{(j)}$ to the nearest cluster mean $\mathbf{c}^{(j)}$.

• What will *k*-means output (cluster assignments & means)?

Caveats

• The output of *k*-means can differ depending on how the algorithm was initialized.

Caveats

- The output of *k*-means can differ depending on how the algorithm was initialized.
- There is no guarantee that the algorithm will converge to a global minimum.

Caveats

- The output of *k*-means can differ depending on how the algorithm was initialized.
- There is no guarantee that the algorithm will converge to a global minimum.
- Sometimes, clusters can become empty.
 - In this case, we have to eliminate one (or more) of the clusters to avoid dividing by 0:

$$\mathbf{c}_j = \frac{1}{n_j} \sum_{i: a(\mathbf{x}^{(i)}) = j} \mathbf{x}^{(i)}$$