# K-means Clustering: Choosing k

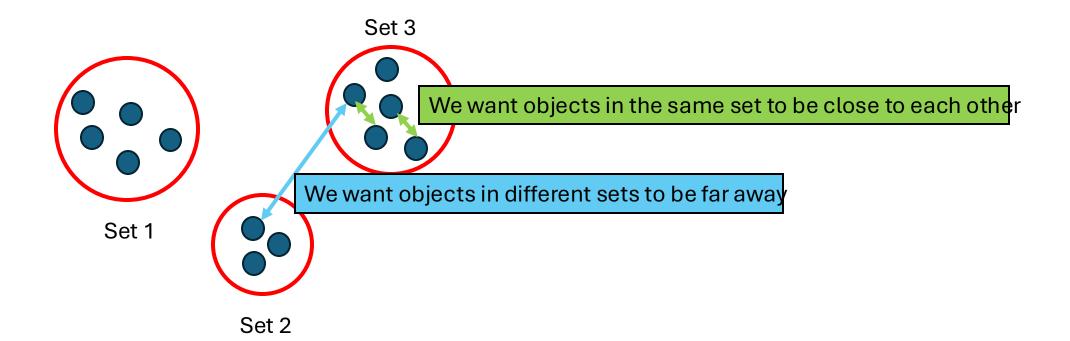
El Kindi Rezig

- In order to do clustering, we need to have two things:
  - The input dataset to be clustered  $X = \{x_1, x_2, ..., x_n\} \subset \mathbb{R}^d$
  - A distance function  $d(x_i, x_j) \ge 0$  that can tell us how similar two points are
  - In this class,  $d(x_i, x_i) = ||x_i x_i||$  (Euclidean distance)

- Goal of clustering is to group objects into k sets (or clusters)
  - Objects in the same set (or cluster) are close to each other

- Goal of clustering is to group objects into k sets (or clusters)
  - Objects in the same set (or cluster) are close to each other
  - Objects in different sets (or clusters) are far away from each other

- Goal of clustering is to group objects into k sets (or clusters)
  - Objects in the same set (or cluster) are close to each other
  - · Objects in different sets (or clusters) are far away from each other



- In mathematical terms, for a set of sites  $S = \{s_1, s_2, ..., s_k\} \subset \mathbb{R}^d$  and a dataset  $X \subset \mathbb{R}^d$  we want to minimize:
  - $cost_2(X,S) = \sum_{i=1}^n (x_i \phi_S(x_i))^2$  Can be thought of as the projection of  $x_i$  onto closest site  $x_j$
  - $\phi_s(x_i) = argmin_{s_i \in S} ||x_i s_j||$  is the site  $s_j$  that is closest to  $x_i$

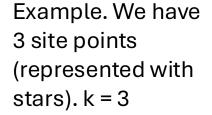
• In mathematical terms, for a set of sites  $S = \{s_1, s_2, ..., s_k\} \subset \mathbb{R}^d$  and a dataset  $X \subset \mathbb{R}^d$  we want to minimize:

• 
$$cost_2(X,S) = \sum_{i=1}^{n} (x_i - \phi_S(x_i))^2$$

•  $\phi_s(x_i) = argmin_{s_i \in S} ||x_i - s_j||$  is the site  $s_j$  that is closest to  $x_i$ 

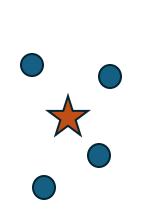








- In mathematical terms, for a set of sites  $S = \{s_1, s_2, ..., s_k\} \subset \mathbb{R}^d$  and a dataset  $X \subset \mathbb{R}^d$  we want to minimize:
  - $cost_2(X,S) = \sum_{i=1}^{n} (x_i \phi_S(x_i))^2$
  - $\phi_s(x_i) = argmin_{s_i \in S} ||x_i s_j||$  is the site  $s_j$  that is closest to  $x_i$

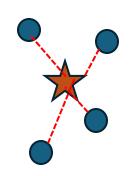


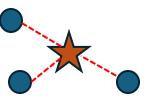


Data points (blue circles)

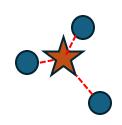


- In mathematical terms, for a set of sites  $S = \{s_1, s_2, ..., s_k\} \subset \mathbb{R}^d$  and a dataset  $X \subset \mathbb{R}^d$  we want to minimize:
  - $cost_2(X,S) = \sum_{i=1}^{n} (x_i \phi_S(x_i))^2$
  - $\phi_S(x_i) = argmin_{S_i \in S} ||x_i s_j||$  is the site  $s_j$  that is closest to  $x_i$

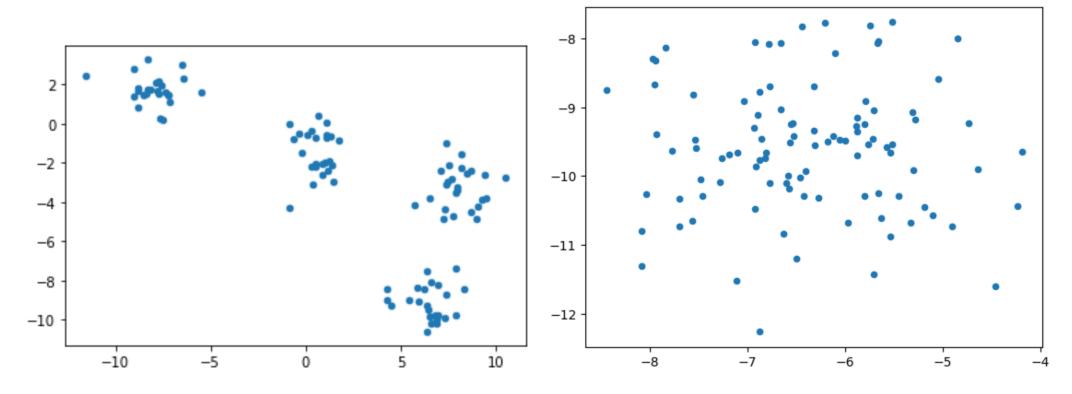




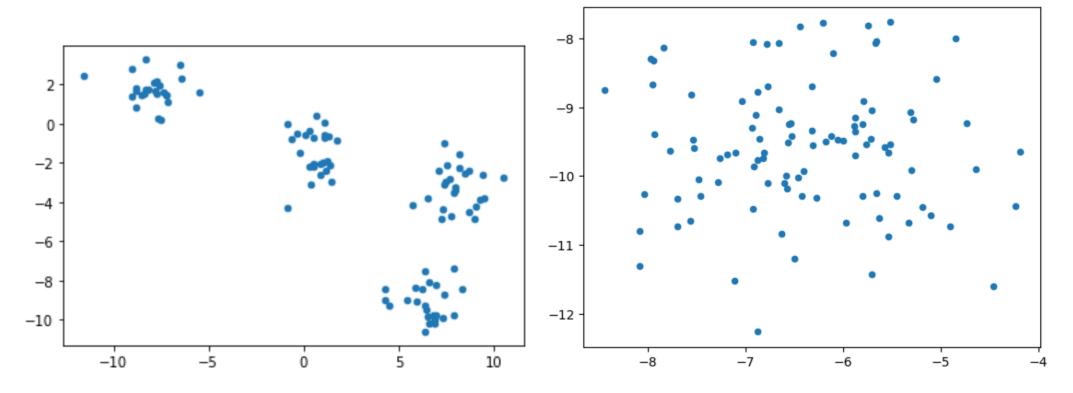
Projecting data points onto closest site which is applying the projection function  $\phi_s(x_i)$  to every data point



Whenever you can, especially when data  $X \subset \mathbb{R}^2$ , you should plot the data! Consider the two data sets shown:

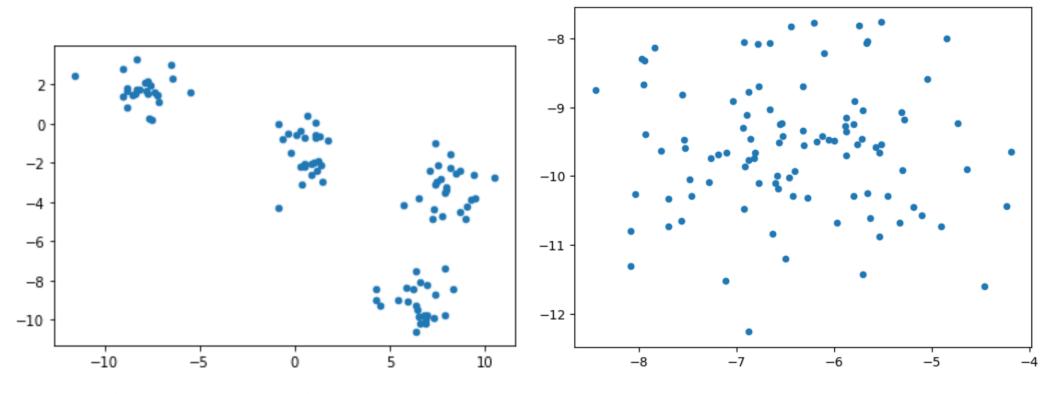


Whenever you can, especially when data  $X \subset \mathbb{R}^2$ , you should plot the data! Consider the two data sets shown:



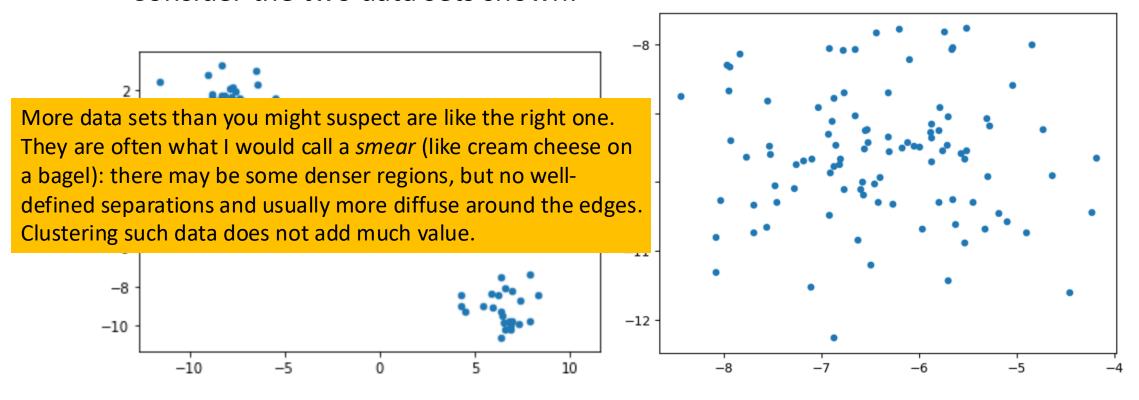
Which one shows "better" clusters?

Whenever you can, especially when data  $X \subset \mathbb{R}^2$ , you should plot the data! Consider the two data sets shown:



The left one has 4 well-defined blobs, the right one does not.

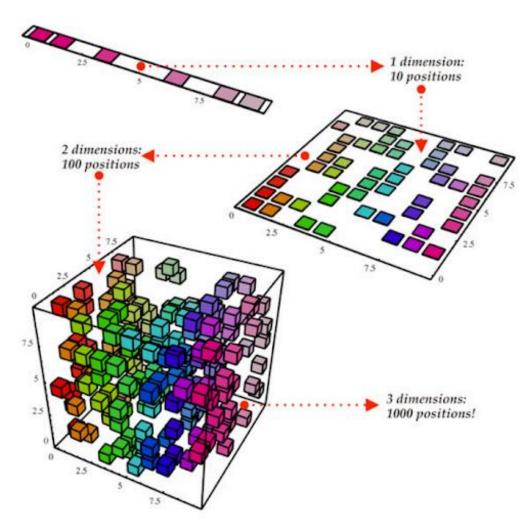
Whenever you can, especially when data  $X \subset \mathbb{R}^2$ , you should plot the data! Consider the two data sets shown:



The left one has 4 well-defined blobs, the right one does not.

- If your data is not naturally in  $\mathbb{R}^2$ , you may need to apply dimensionality reduction to visualize it.
- One method already seen: Laplacian Eigenmaps (covered in Lecture 10 as part of spectral clustering).

- If your data is not naturally in  $\mathbb{R}^2$ , you may need to apply dimensionality reduction to visualize it.
- One method already seen: Laplacian Eigenmaps (covered in Lecture 10 as part of spectral clustering).
- Additional methods that will be covered in this course:
  - PCA (Principal Component Analysis)
  - MDS (Multidimensional Scaling)
  - Distance metric learning techniques
- Choosing a different distance metric can change the visualization:
  - It may transform a diffuse "smear" into more distinct clusters/blobs.



Dimensionality reduction captures all the "significant" directions In which the data is changing

- Dimensionality reduction and visualization approaches are not limited to assignment-based clustering.
- They can also work for non-centrally symmetric clusters (e.g., irregular shapes).
- Regarding the two-moons example:
  - If you think it's relevant, you should demonstrate real-world data that resembles it.
  - Artificial examples generated with t-SNE or UMAP don't count as evidence.
- Regardless, you can still plot your data and circle clusters to highlight structure.

• In clustering methods, we define a **cost function** (Cost of a clustering).

- In clustering methods, we define a **cost function** (Cost of a clustering).
- The goal is to find the clustering that **minimizes this cost** (e.g., using Lloyd's algorithm).

- In clustering methods, we define a **cost function** (Cost of a clustering).
- The goal is to find the clustering that **minimizes this cost** (e.g., using Lloyd's algorithm).
- One might think:
  - Evaluate the cost for each choice of k.
  - Return the k with the smallest cost.

- In clustering methods, we define a **cost function** (Cost of a clustering).
- The goal is to find the clustering that **minimizes this cost** (e.g., using Lloyd's algorithm).
- One might think:
  - Evaluate the cost for each choice of k.
  - Return the k with the smallest cost.
- Problem: For most formulations, the cost will always decrease as k increases.

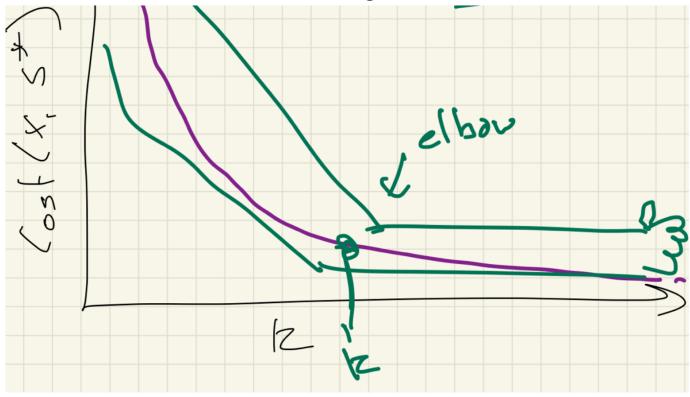
- In clustering methods, we define a cost function (Cost of a clustering).
- The goal is to find the clustering that **minimizes this cost** (e.g., using Lloyd's algorithm).
- One might think:
  - Evaluate the cost for each choice of k.
  - Return the k with the smallest cost.
- Problem: For most formulations, the cost will always decrease as k increases.
- Simply minimizing cost would always favor larger k, which is not meaningful.

• For a clustering defined by sites  $s = \{s_1, s_2, ..., s_k\}$  for a data set X, consider the cost

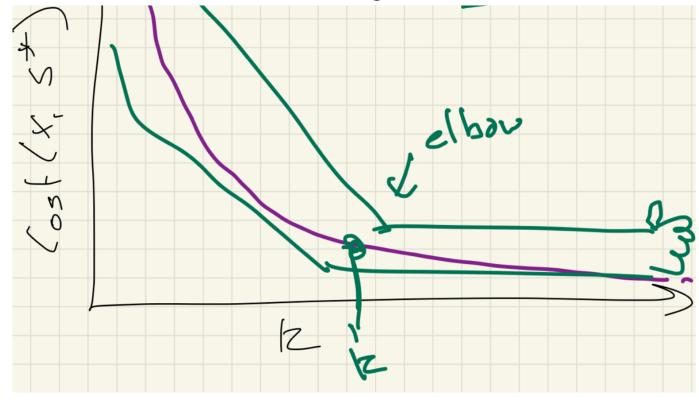
$$\mathsf{Cost}_2(S, X) = \frac{1}{|X|} \sum_{x \in X} (x - \phi_S(x))^2$$

Recall that:  $\phi_S(x) = \arg\min_{s_j \in S} \|x - s\|$ 

For some choice of k, if we let  $S^*$  be the optimal clustering in terms of  $Cost_2(S^*, X)$ , then we can plot the score for each k. It will look something like this:

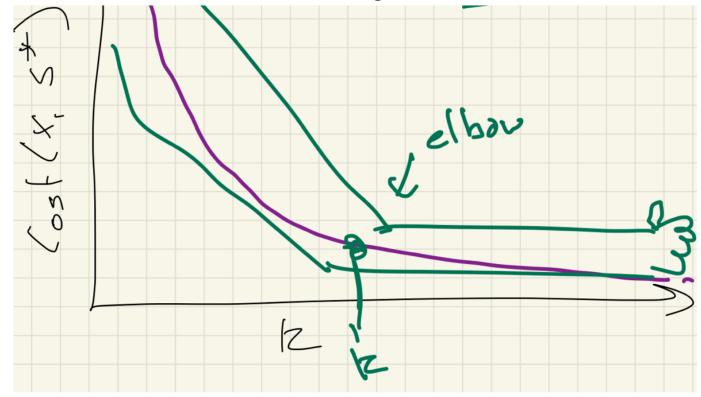


For some choice of k, if we let  $S^*$  be the optimal clustering in terms of  $Cost_2(S^*, X)$ , then we can plot the score for each k. It will look something like this:



What happens when k = n?

For some choice of k, if we let  $S^*$  be the optimal clustering in terms of  $Cost_2(S^*, X)$ , then we can plot the score for each k. It will look something like this:



Now that as k increases the curve (in purple) decreases. It should start high, and then decrease towards 0. When k = n it will be exactly 0, since each  $x \in X$  can be a center, and the distance to the closest one is 0.

- Suppose there are **well-defined clusters** for some choice of  $k^*$ .
- When  $k < k^*$ :

- Suppose there are well-defined clusters for some choice of  $k^*$ .
- When  $k < k^*$ :
  - Multiple true clusters get merged into one site  $s_i$ .
  - This causes the cost to be high.

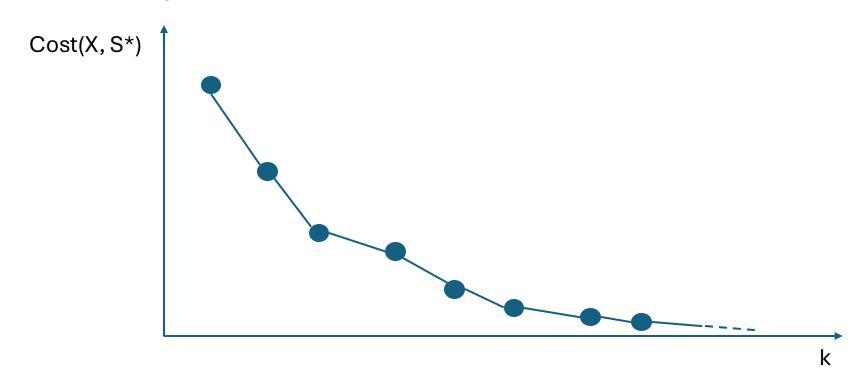
- Suppose there are well-defined clusters for some choice of  $k^*$ .
- When  $k < k^*$ :
  - Multiple true clusters get merged into one site  $s_i$ .
  - This causes the cost to be high.
- When  $k > k^*$ :

- Suppose there are **well-defined clusters** for some choice of  $k^*$ .
- When  $k < k^*$ :
  - Multiple true clusters get merged into one site  $s_i$ .
  - This causes the cost to be high.
- When  $k > k^*$ :
  - True clusters get split across multiple sites.
  - The cost decreases only **slightly**, since points were already in compact clusters.

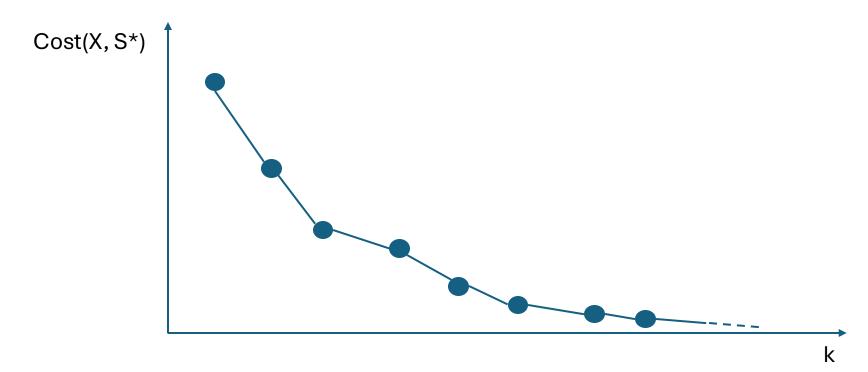
- Suppose there are **well-defined clusters** for some choice of  $k^*$ .
- When  $k < k^*$ :
  - Multiple true clusters get merged into one site  $s_i$ .
  - This causes the cost to be high.
- When  $k > k^*$ :
  - True clusters get split across multiple sites.
  - The cost decreases only **slightly**, since points were already in compact clusters.
- The **elbow point** occurs where:

- Suppose there are well-defined clusters for some choice of  $k^*$ .
- When  $k < k^*$ :
  - Multiple true clusters get merged into one site  $s_i$ .
  - This causes the cost to be high.
- When  $k > k^*$ :
  - True clusters get split across multiple sites.
  - The cost decreases only **slightly**, since points were already in compact clusters.
- The **elbow point** occurs where:
  - The cost transitions from rapidly decreasing to slowly decreasing.
  - This bend in the curve (like an arm) indicates a **good choice of** k.

Which elbow point is best?



Which elbow point is best?



Caveat: Multiple elbow points are possible

#### Caveat: Elbow point is not always clear-cut

- Even in well-clustered data, some points are very well separated → splitting them lowers cost a lot.
- Some clusters are close together → splitting has only a moderate effect on cost.
- Some clusters are spread out → splitting them still reduces cost noticeably.
- $\Rightarrow$  Even with good clustering, the elbow point may not be obvious.

## Caveats of the Elbow method

### Caveat: Elbow point is not always clear-cut

- Even in well-clustered data, some points are very well separated → splitting them lowers cost a lot.
- Some clusters are close together → splitting has only a moderate effect on cost.
- Some clusters are spread out → splitting them still reduces cost noticeably.
- $\Rightarrow$  Even with good clustering, the elbow point may not be obvious.

#### Positive side

• Getting k slightly wrong may still be acceptable.

## Caveats of the Elbow method

#### Caveat: Elbow point is not always clear-cut

- Even in well-clustered data, some points are very well separated → splitting them lowers cost a lot.
- Some clusters are close together → splitting has only a moderate effect on cost.
- Some clusters are spread out → splitting them still reduces cost noticeably.
- $\Rightarrow$  Even with good clustering, the elbow point may not be obvious.

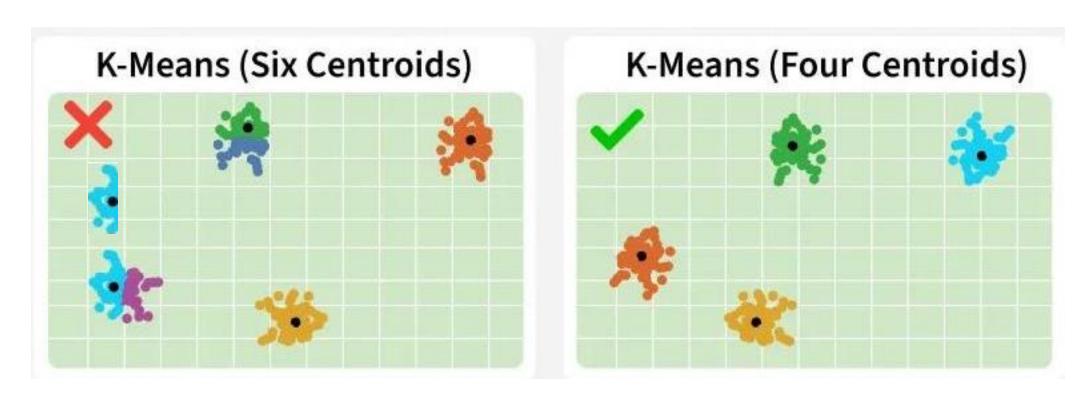
#### Positive side

• Getting k slightly wrong may still be acceptable.

#### Additional complication: Hierarchical structure

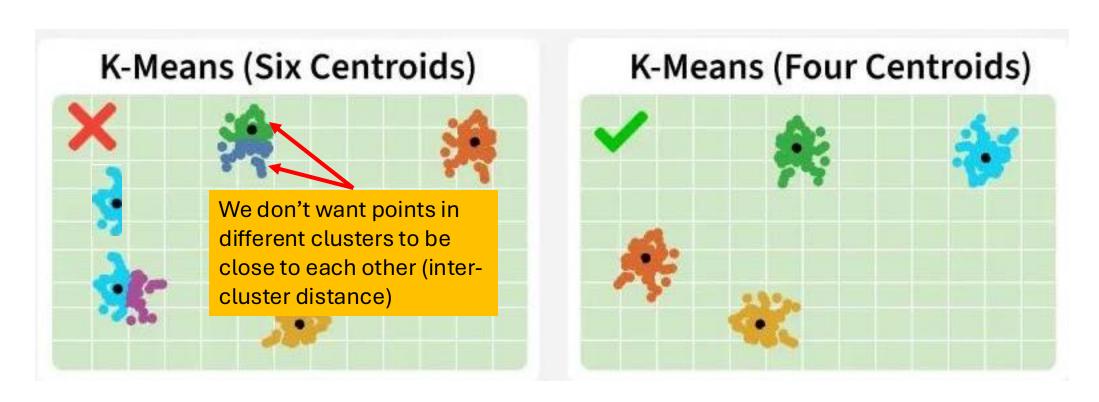
- Data may have multiple levels of clusters (clusters within clusters).
- This can produce two elbow points.
- Either choice may be reasonable—pick the one that best fits the analysis scale.

# What makes a clustering model good/bad?



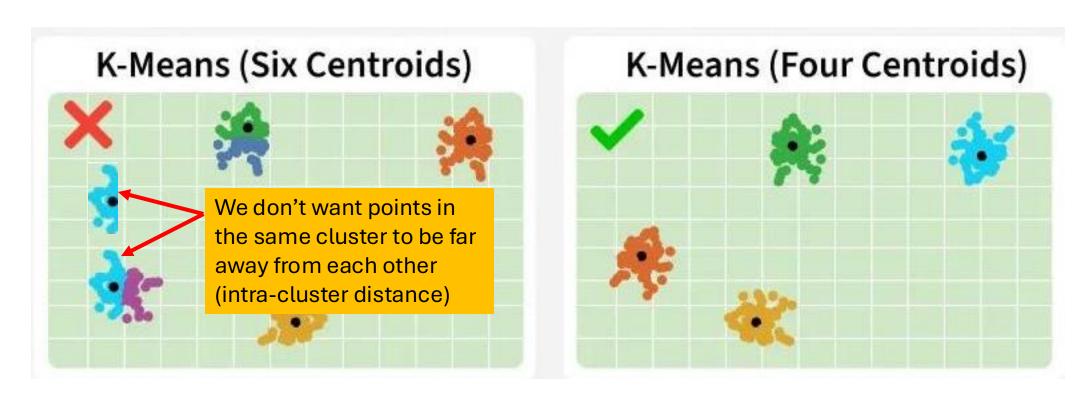
Img source: <a href="https://www.geeksforgeeks.org/machine-learning/what-is-silhouette-score/">https://www.geeksforgeeks.org/machine-learning/what-is-silhouette-score/</a> (modified)

# What makes a clustering model good/bad?



Img source: https://www.geeksforgeeks.org/machine-learning/what-is-silhouette-score/

# What makes a clustering model good/bad?



Img source: https://www.geeksforgeeks.org/machine-learning/what-is-silhouette-score/

• The elbow method started to quantify how good a given "k" is, but the ultimately resorted to drawing a plot, and using human judgement. The Silhouette score helps quantify this choice

- The elbow method started to quantify how good a given "k" is, but the ultimately resorted to drawing a plot, and using human judgement. The Silhouette score helps quantify this choice
- The silhouette score measures how well a point fits within its assigned cluster compared to other clusters.

- The elbow method started to quantify how good a given "k" is, but the ultimately resorted to drawing a plot, and using human judgement. The Silhouette score helps quantify this choice
- The silhouette score measures how well a point fits within its assigned cluster compared to other clusters.
- For each point  $x_i$ :
- Cohesion (within-cluster distance)
  - Compute a(i): the average distance from  $x_i$  to all other points in its own cluster.
  - What does a small a(i) mean?

- The elbow method started to quantify how good a given "k" is, but the ultimately resorted to drawing a plot, and using human judgement. The Silhouette score helps quantify this choice
- The silhouette score measures how well a point fits within its assigned cluster compared to other clusters.
- For each point  $x_i$ :
- Cohesion (within-cluster distance)
  - Compute a(i): the average distance from  $x_i$  to all other points in its own cluster.
  - Small a(i) means  $x_i$  is close to its cluster buddies  $\rightarrow$  good.

- The elbow method started to quantify how good a given "k" is, but the ultimately resorted to drawing a plot, and using human judgement. The Silhouette score helps quantify this choice
- The silhouette score measures how well a point fits within its assigned cluster compared to other clusters.
- For each point  $x_i$ :
- Cohesion (within-cluster distance)
  - Compute a(i): the average distance from  $x_i$  to all other points in its own cluster.
  - Small a(i) means  $x_i$  is close to its cluster mates  $\rightarrow$  good.
- Separation (distance to other clusters)
  - Compute b(i) :the smallest average distance from  $x_i$  to all points in another cluster (the "next-best" cluster).
  - What does a large b(i) means?

- The elbow method started to quantify how good a given "k" is, but the ultimately resorted to drawing a plot, and using human judgement. The Silhouette score helps quantify this choice
- The silhouette score measures how well a point fits within its assigned cluster compared to other clusters.
- For each point  $x_i$ :
- Cohesion (within-cluster distance)
  - Compute a(i): the average distance from  $x_i$  to all other points in its own cluster.
  - Small a(i) means  $x_i$  is close to its cluster mates  $\rightarrow$  good.
- Separation (distance to other clusters)
  - Compute b(i) : the smallest average distance from  $x_i$  to all points in another cluster (the "next-best" cluster).
  - Large b(i) means  $x_i$  is far from other clusters  $\rightarrow$  also good.

### For each point $x_i$ :

- Cohesion (within-cluster distance)
  - Compute a(i): the average distance from  $x_i$  to all other points in its own cluster.
  - Small a(i) means  $x_i$  is close to its cluster mates  $\rightarrow$  good.

### For each point $x_i$ :

- Cohesion (within-cluster distance)
  - Compute a(i): the average distance from  $x_i$  to all other points in its own cluster.
  - Small a(i) means  $x_i$  is close to its cluster mates  $\rightarrow$  good.
- This assumes some model similar to k-means, k-medioid, or mean-link HAC. We want to quantify a disjoint clustering  $S_1, S_2, \ldots, S_k$ . We then quantify the average *inter-cluster distance* for a point  $x_i \in X$  in cluster j as:

$$a(i) = rac{1}{|S_j|-1} \sum_{x \in S_j; x 
eq x_i} \mathbf{d}(x_i, x)$$

#### For each point $x_i$ :

- Separation (distance to other clusters)
  - Compute b(i): the smallest average distance from  $x_i$  to all points in another cluster (the "next-best" cluster).
  - Large b(i) means  $x_i$  is far from other clusters  $\rightarrow$  also good.

### For each point $x_i$ :

- Separation (distance to other clusters)
  - Compute b(i): the smallest average distance from  $x_i$  to all points in another cluster (the "next-best" cluster).
  - Large b(i) means  $x_i$  is far from other clusters  $\rightarrow$  also good.
- We also consider the average replacement cluster score again for a point  $x_i \in X$  in cluster j where

$$b(i) = \min_{j' 
eq j} rac{1}{|S_{j'}|} \sum_{x \in S_{j'}} \mathbf{d}(x_i, x)$$

This is what the score a(i) would be for  $b(i) = \min_{j' \neq j} \frac{1}{|S_{j'}|} \sum_{x \in S_{i'}} \mathbf{d}(x_i, x)$   $\mathbf{x}_i$  if it could not use cluster  $S_j$  and instead had to use the next best x<sub>i</sub> if it could not use cluster S<sub>i</sub> and option (which would be cluster S<sub>i'</sub>)

With these values, we can define the Silhouette score for a point  $x_i \in X$  as

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}.$$

Note that a(i) is not defined if the cluster for  $x_i$  is of size 1. In this case, we define s(i) = 0.

With these values, we can define the Silhouette score for a point  $x_i \in X$  as

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}.$$

Note that a(i) is not defined if the cluster for  $x_i$  is of size 1. In this case, we define s(i) = 0.

How to interpret s(i) > 0

With these values, we can define the Silhouette score for a point  $x_i \in X$  as

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}.$$

Note that a(i) is not defined if the cluster for  $x_i$  is of size 1. In this case, we define s(i) = 0.

Note that  $s(i) \in [-1, 1]$ , and that if s(i) > 0, then its "better off" in its cluster than the replacement one scored by b(i).

The average silhouette score for the entire clustering is reported as

$$\operatorname{sil}(S, X) = \frac{1}{|X|} \sum_{x_i \in X} s(i)$$

Finally, we can choose k as the value which results in the clustering with highest average silhouette score. This is a popular way to fully automate the decision on k, but note that it assumes a specific model of what constitutes a good cluster. For instance, this may not be meaningful for density based clustering or single-link HAC.

 The BIC is a model selection criterion that tries to balance two things:

#### Goodness of fit

- How well the model explains the data (e.g., likelihood).
- A more complex model (higher k) will usually fit better.

### Model simplicity (penalty for complexity)

- Adding parameters always risks overfitting.
- BIC penalizes models with more parameters, especially when the dataset is large.

- We want to fit each cluster  $S_j$  with a generative likelihood model fj. Given a data element x, we can measure  $f_j(x)$  which evaluates how likely a point is to come from the model. It is a positive probability density function, so it is normalized so its integral is 1.
- Let's make this concrete by fitting each cluster  $S_j \subset X$  with an isotropic Gaussian (multi-dimensional normal) distribution. The *isotropic* term means we consider the same variance  $\sigma$  in each direction; for simplicity, assume  $\sigma$  is fixed. To define this we need a center parameter  $s_i \in R^d$ . Then

$$f_j(x) = \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-\frac{\|x - s_j\|^2}{2\sigma^2}\right)$$

 $\bullet$  And the the likelihood for a cluster, assuming the data is iid from  $f_j$  is

$$f_j(S_j) = \prod_{x \in S_j} f_j(x) = \prod_{x \in S_j} \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-\frac{\|x - s_j\|^2}{2\sigma^2}\right)$$

Indeed if we were to maximize this for  $S_j$  over the choice of center  $s_j$ , the mean  $\frac{1}{|S_j|} \sum_{x \in S_j}$  is the optimal choice: the *maximum likelihood estimator*.

We can assume we have defined a likelihood for a set  $f(S_k)$  for the best clustering  $S_k = \{S_1, S_2, ..., S_k\}$ .

It is typically more numerically stable to work with the *negative log-likelihood*:

$$\ell(\mathcal{S}_k) = -\ln(f(\mathcal{S}_k))$$

We can assume we have defined a likelihood for a set  $f(S_k)$  for the best clustering  $S_k = \{S_1, S_2, ..., S_k\}$ .

It is typically more numerically stable to work with the *negative log-likelihood*:

$$\ell(\mathcal{S}_k) = -\ln(f(\mathcal{S}_k))$$

For one cluster S<sub>i</sub>

$$\ell(S_j) = -\ln(f_j(S_j)) = -\sum_{x \in S_j} \ln(f_j) = \sum_{x \in S_j} \left( \frac{\|x - s_j\|^2}{2\sigma^2} - d\ln(\frac{1}{\sqrt{2\pi}\sigma}) \right)$$

Because we negated it, we seek to minimize  $\ell(\cdot)$  when we sought to maximize the likelihood  $f(\cdot)$ .\

- Likelihood challenge: Increasing the number of clusters k always increases the likelihood.
- $\bullet$  As k grows, the model can fit the data more closely (sometimes overfitting).
- $\bullet$  Correspondingly, the negative log-likelihood always decreases with larger k.
- This makes it tricky to use likelihood or negative log-likelihood alone for choosing the "right" k.

- To avoid overfitting, we penalize models with more parameters.
- In k-means clustering in  $\mathbb{R}^d$ :
  - Each cluster center  $s_i$  has d parameters.
  - Total parameters = kd
- Using information-theoretic arguments and a Bayesian perspective, we derive the Bayesian Information Criterion (BIC).

• For a model M with m parameters and n observations, the BIC is:

$$\mathsf{BIC}(M) = -2\ln(f(M)) + m\ln(n)$$

• For a model M with m parameters and n observations, the BIC is:

$$\mathsf{BIC}(M) = -2\ln(f(M)) + m\ln(n)$$

• Since our k-means algorithm has kd parameters, for a best-fit clustering  $S_k$  of size k we have

$$\mathsf{BIC}(\mathcal{S}_k) = -2\ln(f(\mathcal{S}_k)) + kd\ln(|X|)$$

- Now the first term is twice the negative log-likelihood, and so decreases with k increasing. On the other hand, the second term kd ln(|X|) has a fixed quantity dln(|X|) and so increases linearly with k.
- The value k which minimizes BIC(Sk) provides a choice for k.

- The BIC method applies broadly to any model with a well-defined likelihood function.
- It is well-defined for:
  - Mixture of Gaussians
  - k-means
- For Hierarchical Agglomerative Clustering (HAC):
  - Application of BIC is possible but debated.
- For Spectral clustering or DBScan:
  - Applying BIC is difficult and not straightforward.