

INTRO to DATA SCIENCE

CLUSTER ANALYSIS

LAST TIME:

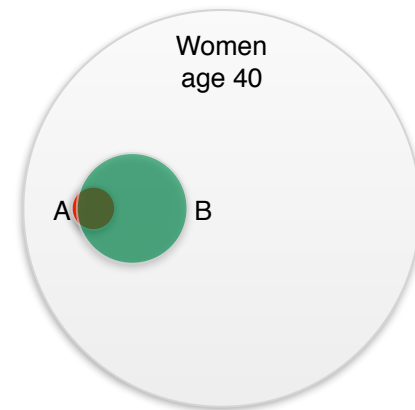
I. INTRO TO PROBABILITY

II. NAÏVE BAYESIAN CLASSIFICATION

EXERCISES:

III. NAÏVE BAYES CLASSIFICATION IN PYTHON

QUESTIONS?



INTRO TO DATA SCIENCE

QUESTIONS?

WHAT WAS THE MOST INTERESTING THING YOU LEARNT?

WHAT WAS THE HARDEST TO GRASP?

- I. CLUSTER ANALYSIS**
- II. THE K-MEANS ALGORITHM**
- III. CHOOSING K**
- IV. EXAMPLE**

I. CLUSTER ANALYSIS

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	???	???
<i>unsupervised</i>	???	???

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

Q: What does categorical mean in this context?

Q: What is a cluster?

Q: What is a cluster?

*A: A group of **similar** data points.*

Q: What is a cluster?

*A: A group of **similar** data points.*

The concept of similarity is central to the definition of a cluster, and therefore to cluster analysis.

Q: What is a cluster?

*A: A group of **similar** data points.*

The concept of similarity is central to the definition of a cluster, and therefore to cluster analysis.

In general, greater similarity between points leads to better clustering.

Q: What is the purpose of cluster analysis?

Q: What is the purpose of cluster analysis?

A: To enhance our understanding of a dataset by dividing the data into groups.

Q: What is the purpose of cluster analysis?

A: To enhance our understanding of a dataset by dividing the data into groups.

Clustering provides a layer of abstraction from individual data points.

Q: What is the purpose of cluster analysis?

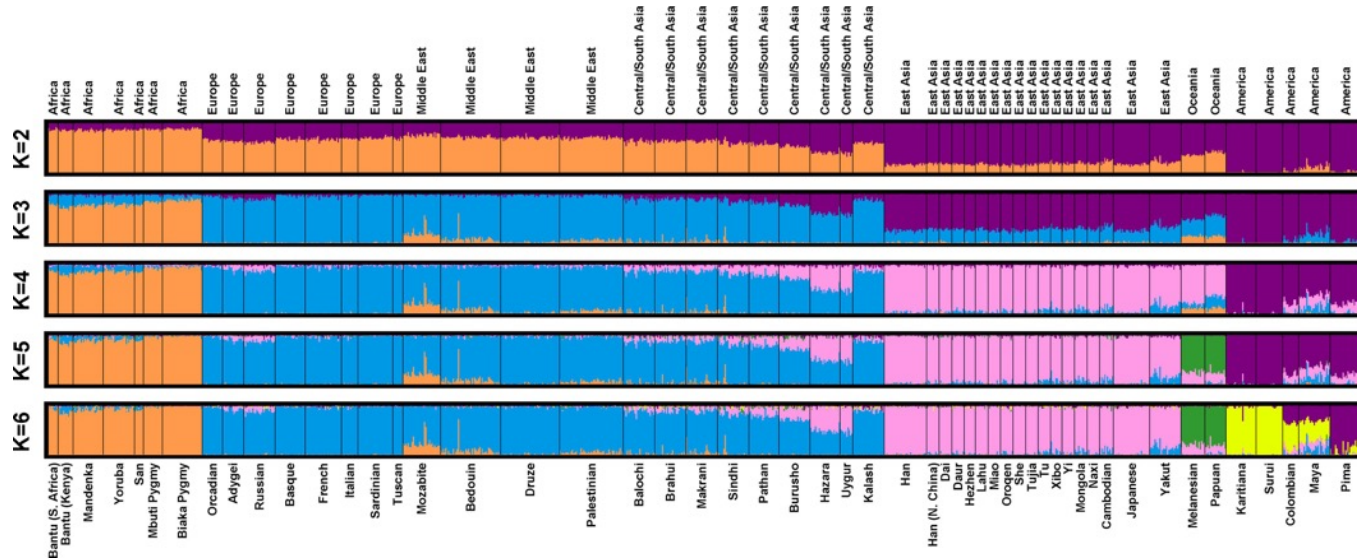
A: To enhance our understanding of a dataset by dividing the data into groups.

Clustering provides a layer of abstraction from individual data points.

The goal is to extract and enhance the natural structure of the data

Clustering can be useful in a wide variety of domains, including genetics, consumer internet and business.

*Clustering can be useful in a wide variety of domains, including **genetics**, consumer internet and business.*



*Clustering can be useful in a wide variety of domains, including genetics, **consumer internet** and business.*

Snowpiercer

2014

Rent or Buy with 1-Click®

★★★★★ 1,220

7/10

Watch Now While It's In Theaters! A post-apocalyptic ice age forces humanity's last survivors aboard a globe-spanning super train. One man (Chris Evans) will risk everything to lead a revolt for control of the engine and the future of the world.

Starring: Chris Evans, Kang-ho Song
 Runtime: 2 hours 7 minutes
 Available to watch on [supported devices](#).

Rent HD \$4.99

Rent SD \$3.99

Buy HD \$14.99

Buy SD \$14.99

[Redeem a gift card or promotion code](#)

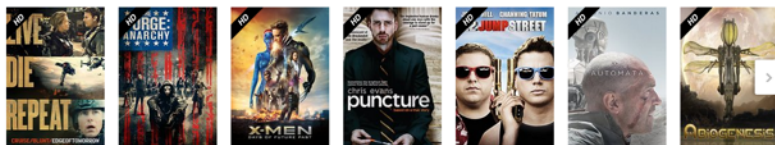
Add to Watchlist

Share

[Send us Feedback](#)

By placing your order, you agree to our [Terms of Use](#). Sold by Amazon Digital Services, Inc. Additional taxes may apply.

Customers Who Watched This Item Also Watched



Google

Search bar with placeholder "What's up?" and search button.

Top navigation links: +Joshia, [Grid Icon], [Bell Icon], [Share Icon], [Profile Picture]

News

U.S. edition | Modern

Top Stories

- San Francisco Giants
- Tim Cook
- International Space Station
- Ledron James
- Thomas Marino
- Facebook
- Halloween
- Taylor Swift
- Burkina Faso
- Arlington, Virginia
- World
- U.S.
- Business
- Technology
- Entertainment
- Sports
- Science
- Health
- Spotlight

FBI Most-Wanted Fugitive Eric Frein Captured Alive

ABC News · 21 minutes ago

Accused cop killer Eric Frein, one of the FBI's 10 Most Wanted Fugitives, was captured after a 48-day manhunt, police said tonight.

Trending on Google+ Eric Frein captured, will face death penalty, Pennsylvania officials say Los Angeles Times

In Depth: Accused Pa. cop killer Eric Frein captured USA TODAY

Related Pennsylvania State Police • Police • Pennsylvania »

ABC News - SF Bay Area

4 hours ago • Google+

CAPTURED! Pennsylvania State Police have captured a man accused of ambushing two troopers, leaving one dead and another seriously injured. Eric Frein was on the run for 48 days. READ MORE ABOUT THE A...

GMA CBS CBS CBS CBS CBS CBS CBS

Attacks on Islamic Site may help Assad admits Chuck Hagel

The Australian · 2 hours ago

PENTAGON Chief Chuck Hagel has acknowledged that US-led air strikes against the Islamic State group could help President Bashar al-Assad.

Trending on Google+ Foreign jihadists flocking to Iraq and Syria on 'viral' edited satellite - UN The Guardian

Opinion: Hagel Worried About Syria Plan Daily Beast

In Depth: Few Arrests of Americans Who Fought In Syria or Iraq, as Feds Focus

*Clustering can be useful in a wide variety of domains, including genetics, consumer internet and **business**.*



There are many kinds of classification procedures. For our class, we will be focusing on K-means clustering, which is one of the most popular clustering algorithms.

K-means is an iterative method that partitions a data set into k clusters.

II. K-MEANS CLUSTERING

Q: How does the algorithm work?

- 1) *choose k initial centroids (note that k is an input)*
- 2) *for each point:*
 - *find distance to each centroid*
 - *assign point to nearest centroid*
- 3) *recalculate centroid positions*
- 4) *repeat steps 2-3 until stopping criteria met*

Q: How do you choose the initial centroid positions?

Q: How do you choose the initial centroid positions?

A: There are several options:

Q: How do you choose the initial centroid positions?

A: There are several options:

- randomly (but may yield divergent behavior)*

Q: How do you choose the initial centroid positions?

A: There are several options:

- randomly (but may yield divergent behavior)*
- perform alternative clustering task, use resulting centroids as initial k -means centroids*

Q: How do you choose the initial centroid positions?

A: There are several options:

- randomly (but may yield divergent behavior)*
- perform alternative clustering task, use resulting centroids as initial k-means centroids*
- start with global centroid, choose point at max distance, repeat (but might select outlier)*

Q: How do you determine which centroid a given point is most similar to?

Q: How do you determine which centroid a given point is most similar to?

The similarity criterion is determined by the measure we choose.

Q: How do you determine which centroid a given point is most similar to?

The similarity criterion is determined by the measure we choose.

*In the case of k -means clustering, the similarity metric is the **Euclidian distance**:*

Q: How do you determine which centroid a given point is most similar to?

The similarity criterion is determined by the measure we choose.

In the case of k-means clustering, the similarity metric is the
Euclidian distance:

$$d(x_1, x_2) = \sqrt{\sum_{i=1}^N (x_{1i} - x_{2i})^2}$$

Q: How do we re-compute the positions of the centers at each iteration of the algorithm?

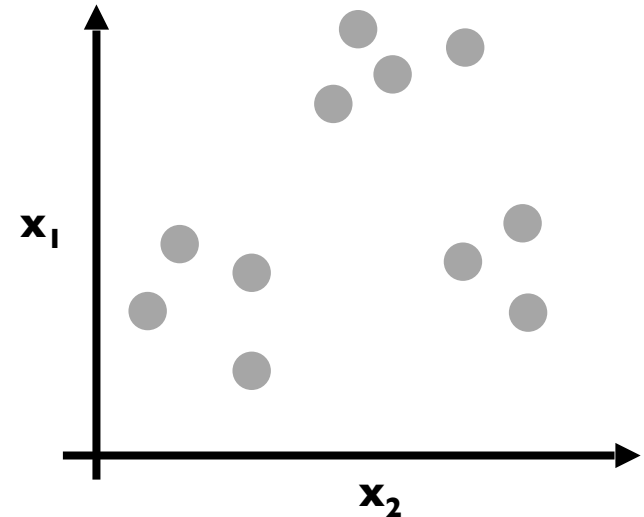
A: By calculating the centroid (i.e., the geometric center)

We iterate until some stopping criteria are met; in general, suitable convergence is achieved in a small number of steps.

We iterate until some stopping criteria are met; in general, suitable convergence is achieved in a small number of steps.

Stopping criteria can be based on the centroids (eg, if positions change by no more than ε) or on the points (eg, if no more than $x\%$ change clusters between iterations).

- 1) *choose k initial centroids (note that k is an input)*
- 2) *for each point:*
 - *find distance to each centroid*
 - *assign point to nearest centroid*
- 3) *recalculate centroid positions*
- 4) *repeat steps 2-3 until stopping criteria met*



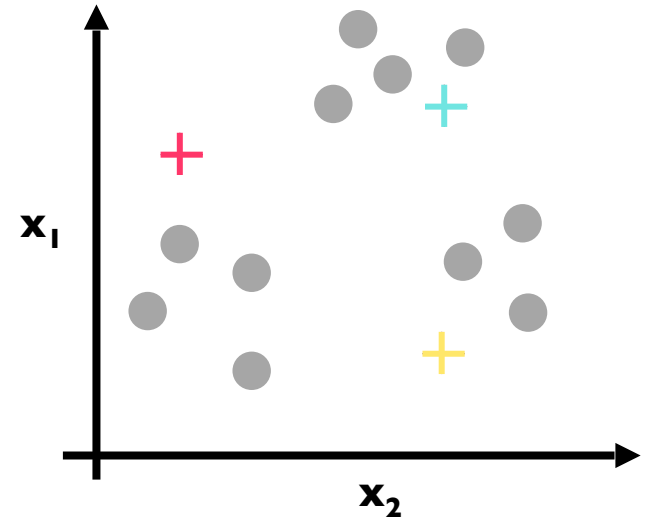
1) choose k initial centroids (note that k is an input)

2) for each point:

- find distance to each centroid
- assign point to nearest centroid

3) recalculate centroid positions

4) repeat steps 2-3 until stopping criteria met



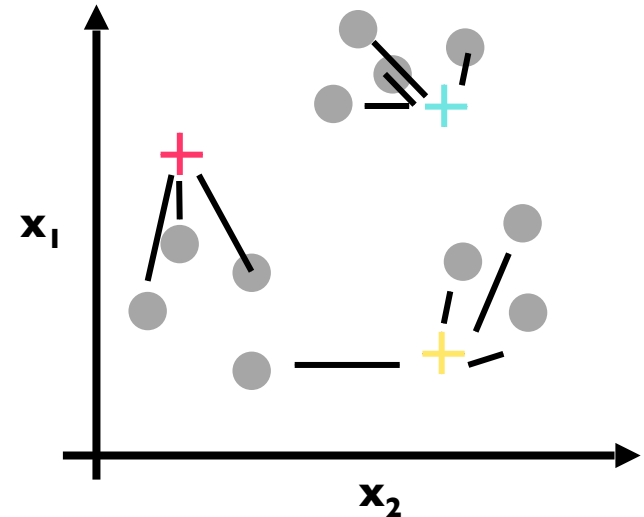
1) *choose k initial centroids (note that k is an input)*

2) ***for each point:***

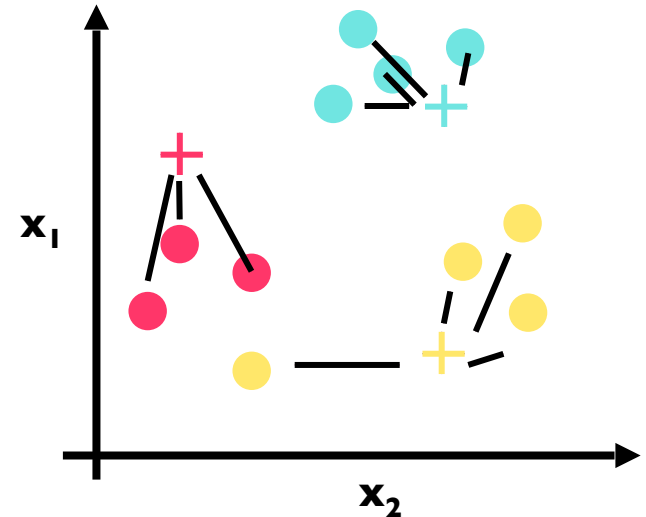
- ***find distance to each centroid***
- *assign point to nearest centroid*

3) *recalculate centroid positions*

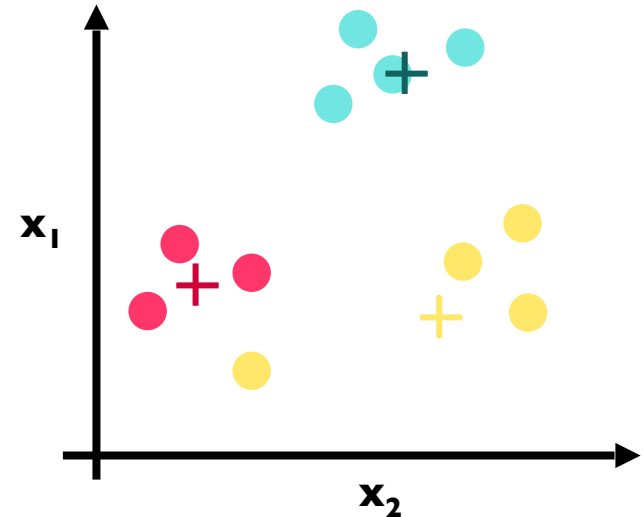
4) *repeat steps 2-3 until stopping criteria met*



- 1) *choose k initial centroids (note that k is an input)*
- 2) ***for each point:***
 - *find distance to each centroid*
 - ***assign point to nearest centroid***
- 3) *recalculate centroid positions*
- 4) *repeat steps 2-3 until stopping criteria met*



- 1) *choose k initial centroids (note that k is an input)*
- 2) *for each point:*
 - *find distance to each centroid*
 - *assign point to nearest centroid*
- 3) ***recalculate centroid positions***
- 4) *repeat steps 2-3 until stopping criteria met*



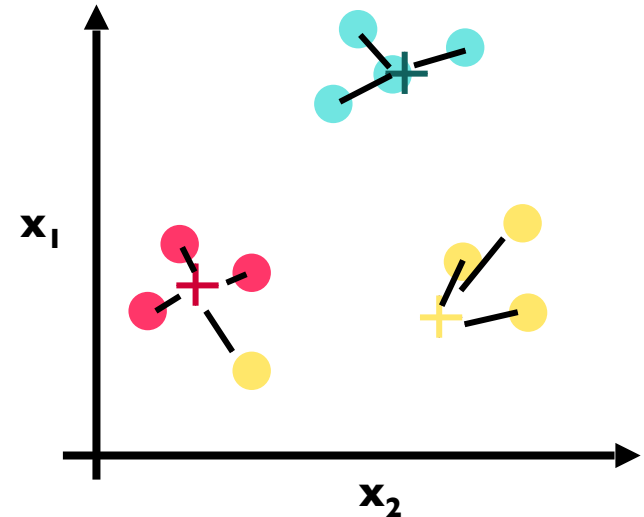
1) *choose k initial centroids (note that k is an input)*

2) *for each point:*

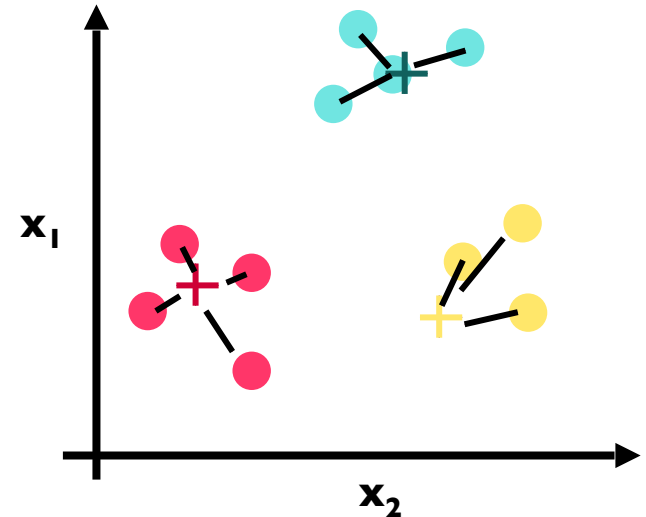
- find distance to each centroid*
- assign point to nearest centroid*

3) *recalculate centroid positions*

4) ***repeat steps 2-3 until stopping criteria met***



- 1) *choose k initial centroids (note that k is an input)*
- 2) *for each point:*
 - *find distance to each centroid*
 - *assign point to nearest centroid*
- 3) *recalculate centroid positions*
- 4) ***repeat steps 2-3 until stopping criteria met***



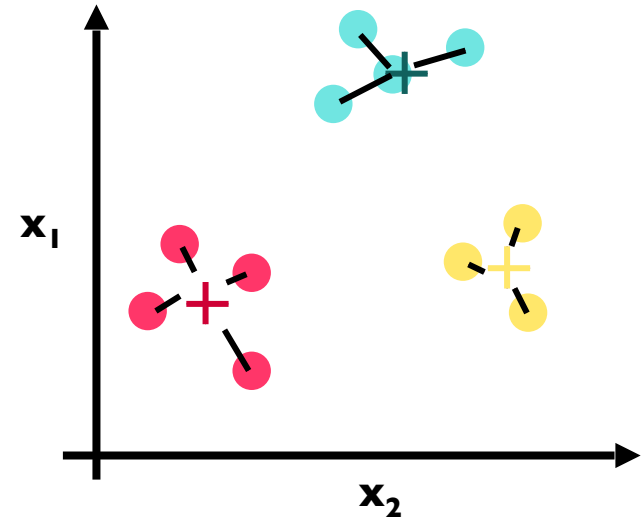
1) *choose k initial centroids (note that k is an input)*

2) *for each point:*

- find distance to each centroid*
- assign point to nearest centroid*

3) *recalculate centroid positions*

4) ***repeat steps 2-3 until stopping criteria met***



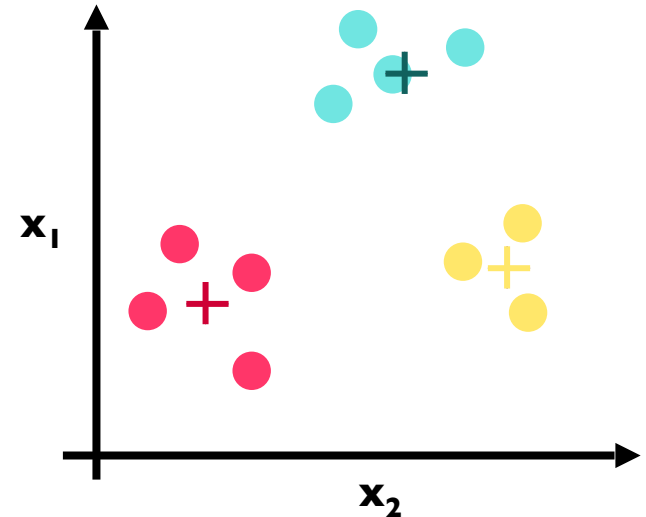
1) *choose k initial centroids (note that k is an input)*

2) *for each point:*

- find distance to each centroid*
- assign point to nearest centroid*

3) *recalculate centroid positions*

4) ***repeat steps 2-3 until stopping criteria met***



III. CLUSTER VALIDATION

In general, k -means will converge to a solution and return a partition of k clusters, even if no natural clusters exist in the data.

In general, k -means will converge to a solution and return a partition of k clusters, even if no natural clusters exist in the data.

*We will look at two validation metrics useful for partitional clustering, **cohesion and separation**.*

Cohesion *measures clustering effectiveness within a cluster.*

$$\hat{C}(C_i) = \sum_{x \in C_i} d(x, c_i)$$

Cohesion *measures clustering effectiveness within a cluster.*

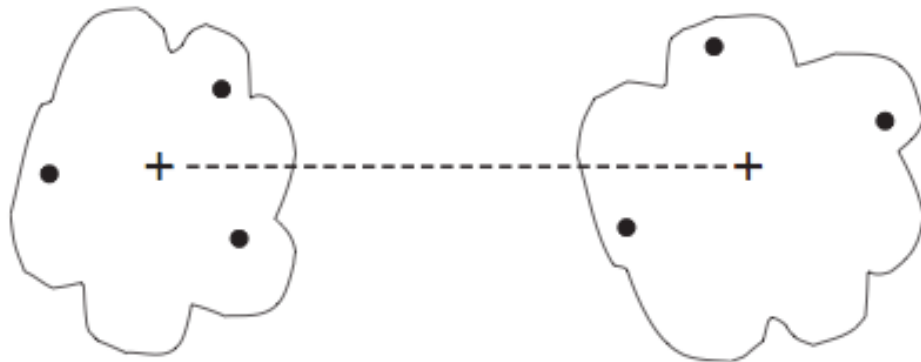
$$\hat{C}(C_i) = \sum_{x \in C_i} d(x, c_i)$$

Separation *measures clustering effectiveness between clusters.*

$$\hat{S}(C_i, C_j) = d(c_i, c_j)$$



(a) Cohesion.



(b) Separation.

Figure 8.28. Prototype-based view of cluster cohesion and separation.

*One useful measure that combines the ideas of cohesion and separation is the **silhouette coefficient**. For point x_i , this is given by:*

$$SC_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

such that:

a_i = average in-cluster distance to x_i

b_{ij} = average between-cluster distance to x_i

$b_i = \min_j(b_{ij})$

The silhouette coefficient can take values between -1 and 1.

In general, we want separation to be high and cohesion to be low. This corresponds to a value of SC close to +1.

A negative silhouette coefficient means the cluster radius is larger than the space between clusters, and thus clusters overlap.

The silhouette coefficient for the cluster C_i is given by the average silhouette coefficient across all points in C_i :

$$SC(C_i) = \frac{1}{m_i} \sum_{x \in C_i} SC_i$$

The silhouette coefficient for the cluster C_i is given by the average silhouette coefficient across all points in C_i :

$$SC(C_i) = \frac{1}{m_i} \sum_{x \in C_i} SC_i$$

The overall silhouette coefficient is given by the average silhouette coefficient across all clusters:

$$SC_{total} = \frac{1}{k} \sum_1^k SC(C_i)$$

The silhouette coefficient for the cluster C_i is given by the average silhouette coefficient across all points in C_i :

$$SC(C_i) = \frac{1}{m_i} \sum_{x \in C_i} SC_i$$

The overall silhouette coefficient is given by the average silhouette coefficient across all clusters:

$$SC_{total} = \frac{1}{k} \sum_1^k SC(C_i)$$

NOTE

This gives a summary measure of the overall clustering quality.

One useful application of cluster validation is to determine the best number of clusters for your dataset.

One useful application of cluster validation is to determine the best number of clusters for your dataset.

Q: How would you do this?

One useful application of cluster validation is to determine the best number of clusters for your dataset.

Q: How would you do this?

A: By computing the SSE or SC for different values of k .

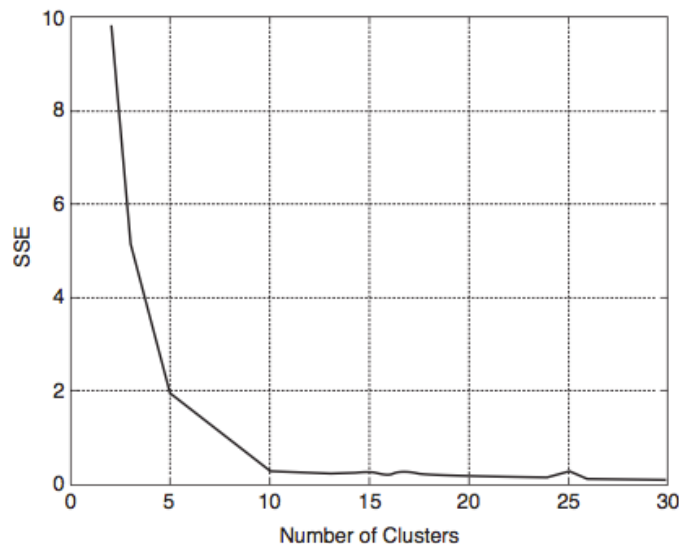


Figure 8.32. SSE versus number of clusters for the data of Figure 8.29.

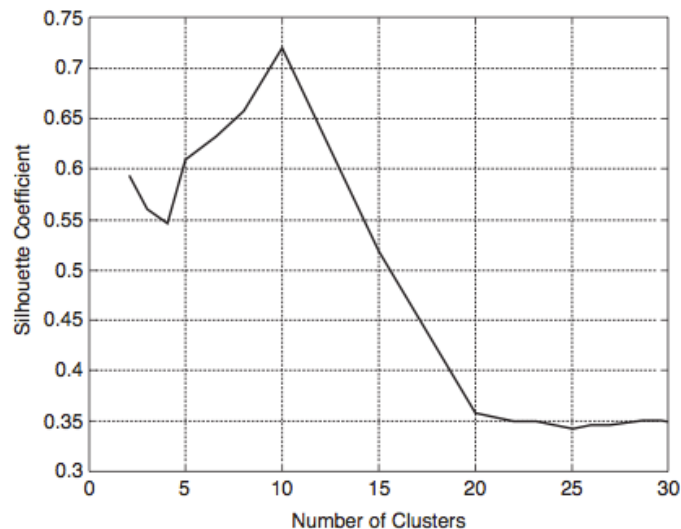


Figure 8.33. Average silhouette coefficient versus number of clusters for the data of Figure 8.29.

Ultimately, cluster validation and clustering in general are suggestive techniques that rely on human interpretation to be meaningful.

Strengths:

K-means is a popular algorithm because of its computational efficiency and simple and intuitive nature.

Strengths:

K-means is a popular algorithm because of its computational efficiency and simple and intuitive nature.

Weaknesses:

However, K-means is highly scale dependent, and is not suitable for data with widely varying shapes and densities.

EX: K-MEANS CLUSTERING