

Lecture 12:
Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Lecture 12: Clustering

Assignment Project Exam Help

<https://eduassistpro.github.io>

Sarah Erfani and Karin Verspoor, CIS

Add WeChat Search edu_assist_pr



THE UNIVERSITY OF

MELBOURNE

Example clusters for the weather dataset

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Outlook	Temperature	Humidity	Wind	Play
rainy	cool	n		no
rainy	cool	n		no
				yes
				yes
				yes
				o

<https://eduassistpro.github.io>

Add: WeChat edu_assist_pr

A possible clustering of the weather dataset

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Outlook	Temperature	Humidity	Windy	Cluster
sunny	hot	high	FALSE	0
sunny	hot	high	TRUE	0
overcast	hot	high	FALSE	0

<https://eduassistpro.github.io>

sunny	mild	high		
sunny	cool	normal		
rainy	mild	normal		
sunny	mild	normal		
overcast	mild	high	TRUE	1
overcast	hot	normal	FALSE	0
rainy	mild	high	TRUE	1

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Outlook	Temperature	Humidity	Windy	Cluster	Play
sunny	hot	high	FALSE	0	no
sunny	hot	high	TRUE	0	no
overcast	hot	high	FALSE	0	yes
		high	FAL	1	yes
		n			yes
		n			no
		n			yes
sunny	mild	high	FALSE	0	
sunny	cool	normal	FALSE	1	
rainy	mild	normal	FALSE	1	
sunny	mild	normal	TRUE	1	
overcast	mild	high	TRUE	1	yes
overcast	hot	normal	FALSE	0	yes
rainy	mild	high	TRUE	1	no

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

■
■
■
■
<https://eduassistpro.github.io>

- Applications in pattern recognition, s
diagnosis . . .

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

- Exclusive vs. overlapping clustering
 - Can an item be in more than one cluster?

- Deterministic vs. probabilistic clustering (Hard vs. soft clustering)

- Partial vs. complete

- In some cases, we only want to cluster

- Heterogeneous vs. homogeneous

- Clusters of widely different sizes, sh

- Incremental vs. batch clustering

- Is the whole set of items clustered in one go?

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

<i>Instance</i> <i>Cluster</i>		<i>Cluster</i>				
		<i>Instance</i>	1	2	3	4
5	2					
6	2					
7	4					
⋮	⋮					
⋮	⋮					

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

-
-
- Able to deal with noise and outliers
- Insensitive to order of input records

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Unsupervised.

- Measures the goodness of a clustering structure without respect to external information. Includes measures of cluster cohesion (compactness, tightness), and measures of cluster separation

<https://eduassistpro.github.io>

- Measures the extent to which the clustering algorithm matches some instance. *Entropy* can measure the extent to which a clustering matches externally supplied class labels.

Relative.

- Compares different clusterings or clusters (using an unsupervised or supervised measure for the purpose of comparison).

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example
Description

Evaluation

Methods

Similarity

k -means

Hierarchical

Most common measure is Sum of Squared Error (SSE)
or *Scatter*

- For each point, the error is the distance to the nearest cluster



<https://eduassistpro.github.io>

- Can show that the m_i that minimizes the SSE (the mean) of the cluster

- Given two clusters, we can choose the

- One easy way to reduce SSE is to increase k , the number of clusters

- However, a good clustering with smaller k can have a lower SSE than a poor clustering with higher k

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

A key component of any clustering algorithm is a measurement of the distance between any points.

- Data points in Euclidean space

- Euclidean distance



b	1	0	1
c	1	1	0

For two bit strings, the number of positions where the symbols are different

- Documents

- Cosine similarity
- Jaccard measure

- Other measures

- Correlation
- Graph-based measures

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

Given k , the k -means algorithm is implemented in four steps.



1. Choose k initial centroids

2. Assign each data point to the nearest centroid

3. Compute the new centroid of each cluster

4. Repeat steps 2 and 3 until the centroids don't change

5. Exclusive, deterministic, partitioni

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

- Initial centroids are often chosen randomly.
- Clusters produced vary from one run to another.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

- Most of the convergence happens i
- Often the stopping condition is chan
'Until relatively few points change cl
(this way the stopping criterion will no
or dimensionality)

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Strengths:

- relatively efficient:
 - $O(ndki)$, where n is no. instances, d is no. attributes, k is no. clusters and i is no. iterations; normally $k, i \ll n$

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

- tends to converge to local minimum; s
try multiple iterations with different s
- need to specify k in advance
- not able to handle non-convex clusters, or clusters of differing densities or sizes
- “mean” ill-defined for nominal or categorical attributes
- may not work well when the data contains outliers

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Bottom-up (= agglomerative) clustering

- Start with single instance clusters
- At each step, join the two closest clusters (in terms of margin)

<https://eduassistpro.github.io>

- Start with one universal cluster
- Find two partitioning clusters
- Proceed recursively on each sub-set
- Can be very fast

In contrast to *k*-means clustering, hierarchical clustering only requires a measure of similarity between *groups* of data points (no seeds, no *k* value).

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

4

new cluster and the original clusters

5

until Only one cluster remains

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

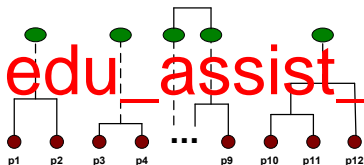
Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

	C1	C2	C3	C4	C5
C1					
C2					
C3					



Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

<https://eduassistpro.github.io>

Updating the proximity matrix:

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

- Single Link: *Minimum* distance between clusters. (most similar members)
- Complete Link: *Maximum* distance between any two points in the two clusters. (most dissimilar members)
- Group Average: *Average* distance between all points (pairwise).

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

	1	2	3	4	5
5	0.20	0.50	0.30	0.80	1.00

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

What are the two closest points?

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

	1	2	3	4	5
1	1.00	0.90	0.10	0.65	0.20
2	0.90	1.00	0.70	0.60	0.50
3	0.10	0.70	1.00	0.40	0.30
4	0.65	0.60	0.40	1.00	0.80
5	0.20	0.50	0.30	0.80	1.00

<https://eduassistpro.github.io>

Update (single link):

Add WeChat edu_assist_pro

	1	2	3	4	5	6
6	—	—	0.70	0.65	0.50	1.00

Update (complete link):

	1	2	3	4	5	6
6	—	—	0.10	0.60	0.20	1.00

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

	1	2	3	4	5
1	1.00	0.90	0.10	0.65	0.20
2	0.90	1.00	0.70	0.60	0.50
3	0.10	0.70	1.00	0.40	0.30
4	0.65	0.60	0.40	1.00	0.80

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Single link

Complete link

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

Assignment Project Exam Help

Clustering is in the eyes of the beholder

■
<https://eduassistpro.github.io>

- Algorithms for Clustering Data (198
<http://homepages.inf.ed.ac.uk/jain/>
Jain_Dubes.pdf

Add WeChat edu_assist_pro

Lecture 12: Clustering

COMP90049
Knowledge
Technologies

Clustering

An example

Description

Evaluation

Methods

Similarity

k-means

Hierarchical

- What basic contrasts are there in different clustering methods?
- How does k-means operate, and what are its strengths and weaknesses?
-
-

<https://eduassistpro.github.io>

Resources:

Tan, Steinbach, Kumar (2006) Introduction to Cluster Analysis

<http://www-users.cs.umn.edu/~kumar/dmbook/ch8.pdf>

Jain, Dubes (1988) Algorithms for Clustering Data. http://homepages.inf.ed.ac.uk/rbf/BOOKS/JAIN/Clustering_Jain_Dubes.pdf