

Assignment Project Exam Help

Cluster Analysis I

<https://tutorcs.com>

An Unsupervised Learning Technique

WeChat: cstutorcs

Presented by Klaus Ackermann



What is Cluster Analysis

Cluster: collection of data objects

- Grouping similar objects in a cluster
- No pre-defined classes → unsupervised classification
- For analysis of data <https://tutorcs.com>
- Cluster analysis – grouping a set of data objects into cluster
- In what contexts would you want to complete a cluster analysis?

WeChat: estutorcs



Cluster Analysis Applications

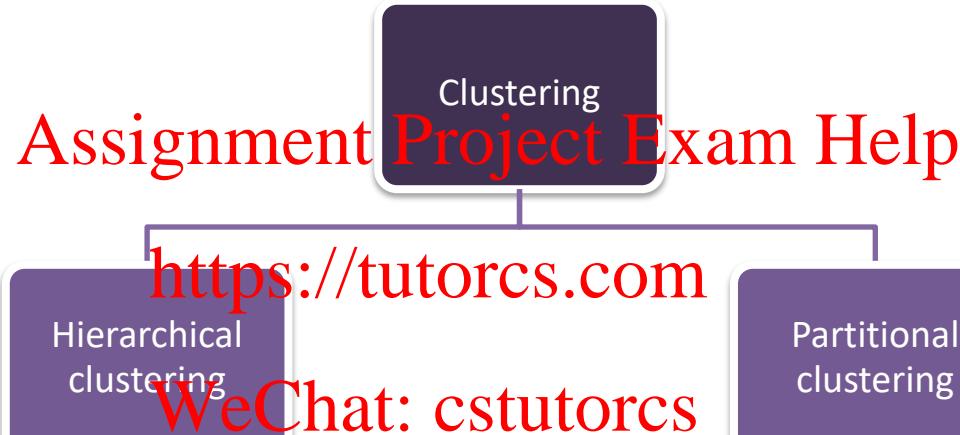
- Market segmentation – targeted marketing
- Market structure analysis – identify groups of similar products according to competitive measures of similarity
- Finance – develop a balanced portfolio
 - eg. Volatility, beta, vega, returns over different periods
- Land use – identify areas of similar land use in Earth observation database
- Medical
 - eg. Types of tumour

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Clustering Methods



Clustering Methods

Hierarchical clustering

Assignment Project Exam Help

<https://tutorcs.com>

Agglomerative

1. Group each observation into cluster of one
2. Merge most similar clusters
3. Merge those a bit further apart
4. Sequence of nested clusters

WeChat: cstutorcs

or

Devisive

1. Observation in one cluster
2. Recursively split into smaller clusters

Clustering Methods

Partitional clustering

Assignment Project Exam Help

<https://tutorcs.com>
K-Means

1. Random assignment of points in the midst of cluster
2. Repeatedly re-assign points cluster
3. Minimise within-cluster distances
4. Maximise between-cluster distances

WeChat: cstutorcs

Distances
between
points in red;
Names of
points in black

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



4

3

5

6

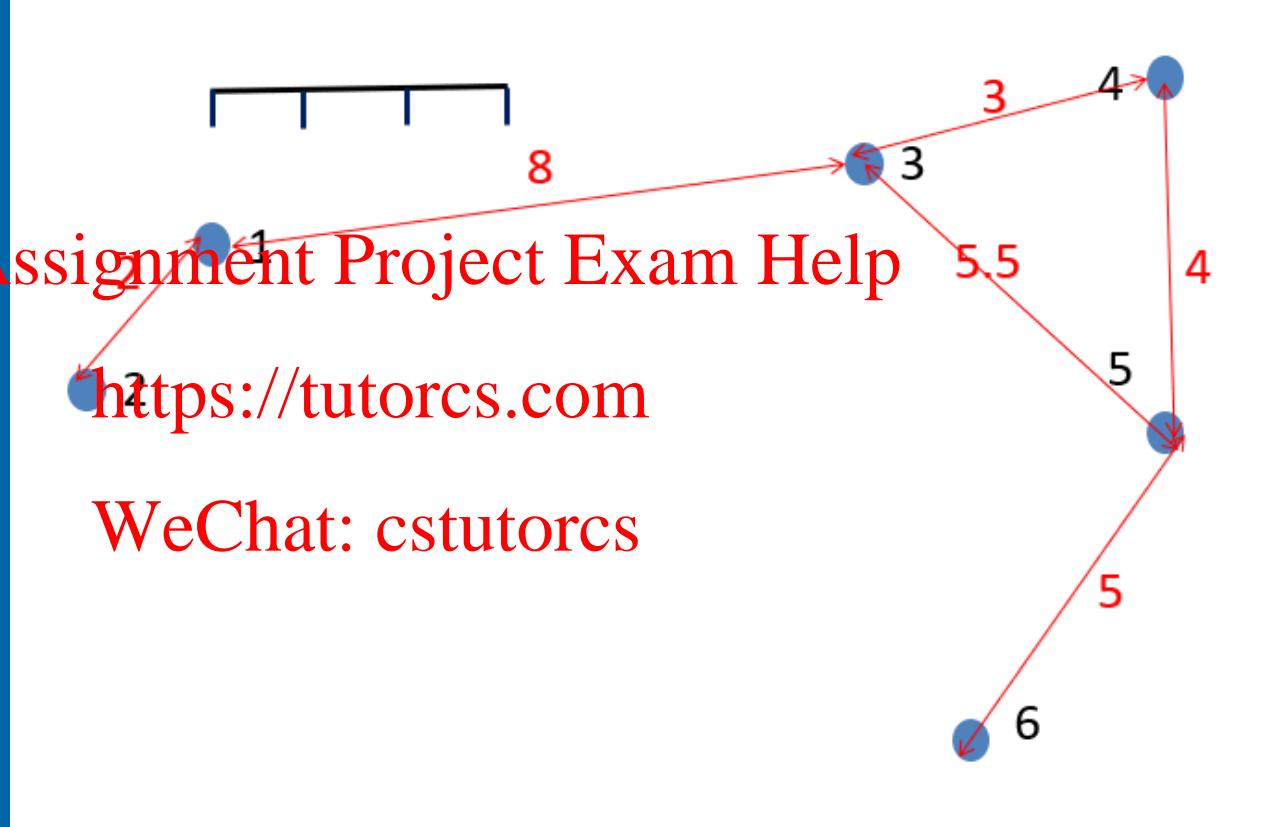


Distances
between
points in red;
Names of
points in black

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs





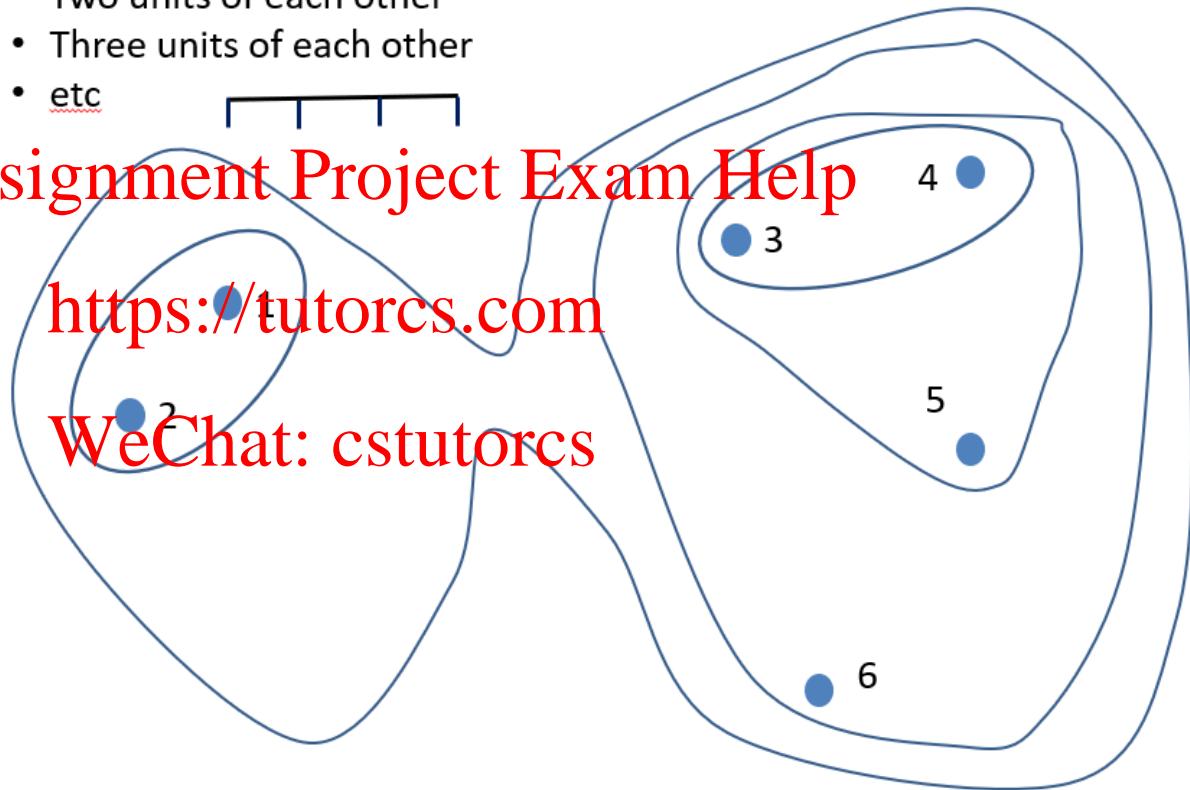
Hierarchical clustering: Suppose in order to be in a cluster, points have to be within:

- Two units of each other
- Three units of each other
- etc

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs





For now, assume distance between clusters is distance between closest points

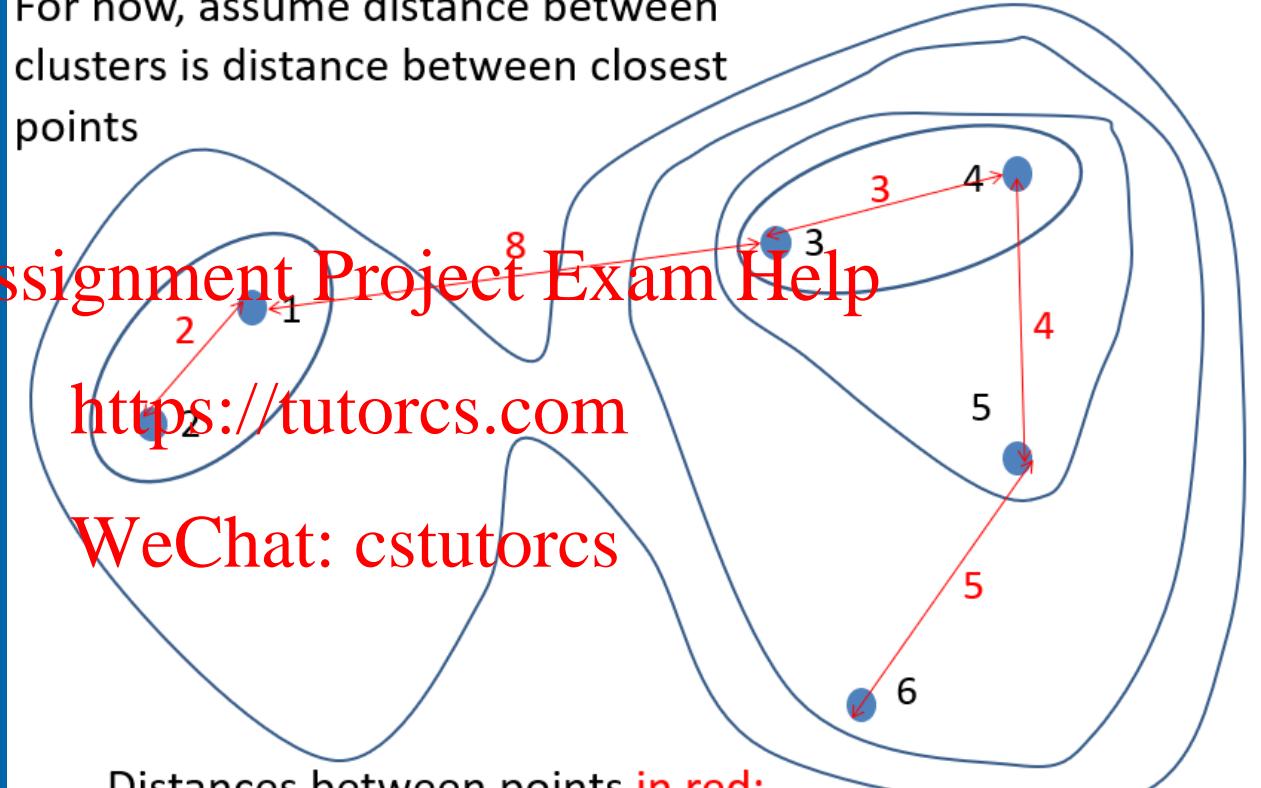
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Distances between points in red;

Names of points in black



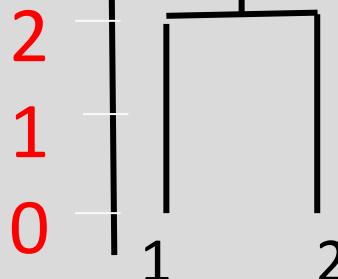


Dendrogram: a diagram that represents a tree

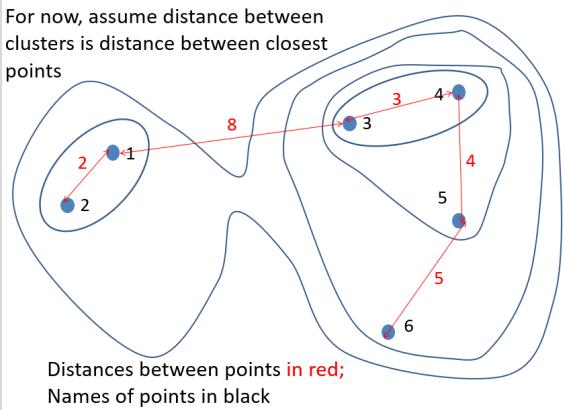
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



For now, assume distance between clusters is distance between closest points





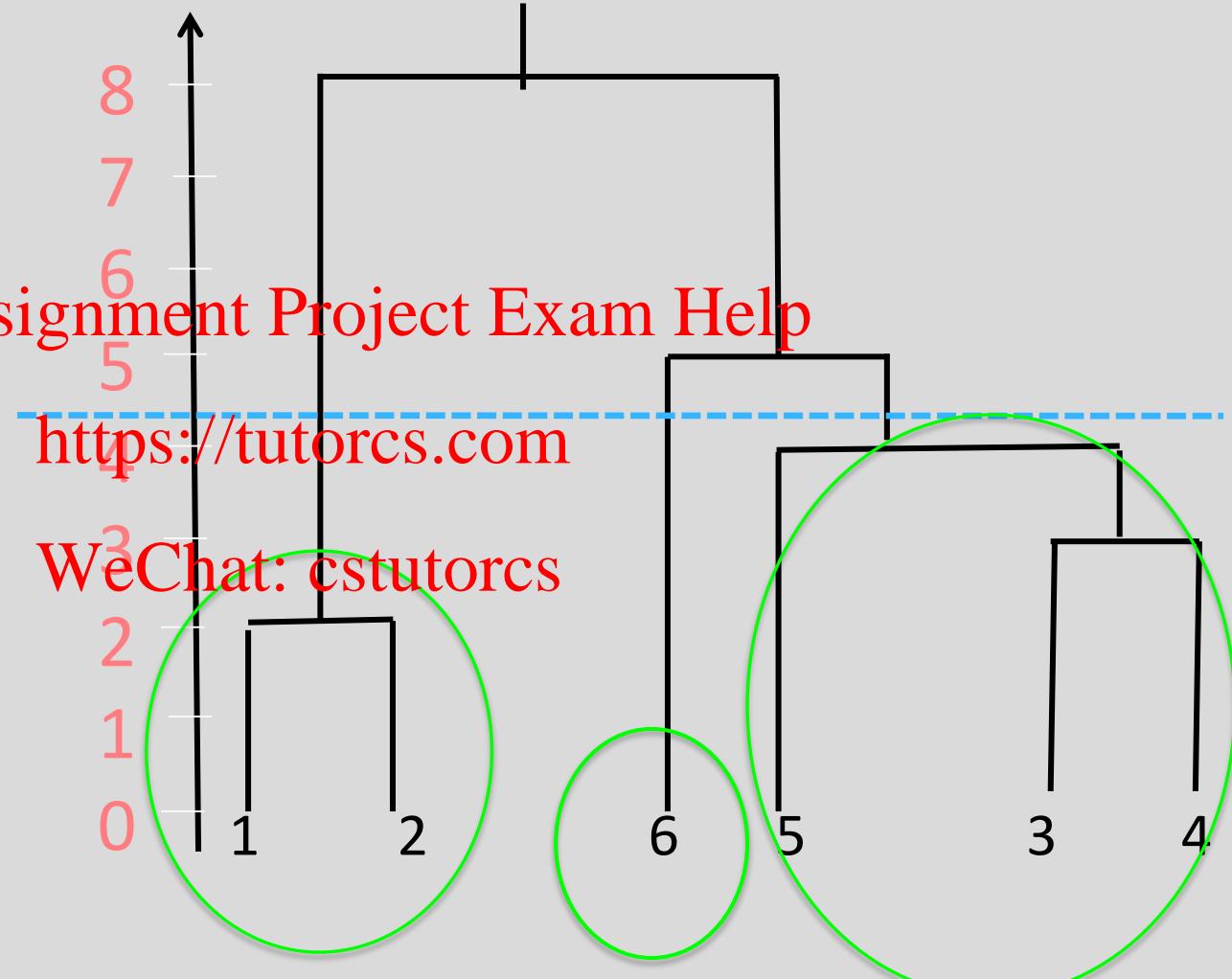
If this distance is just above 4, we create three clusters.

Name the members of the clusters

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

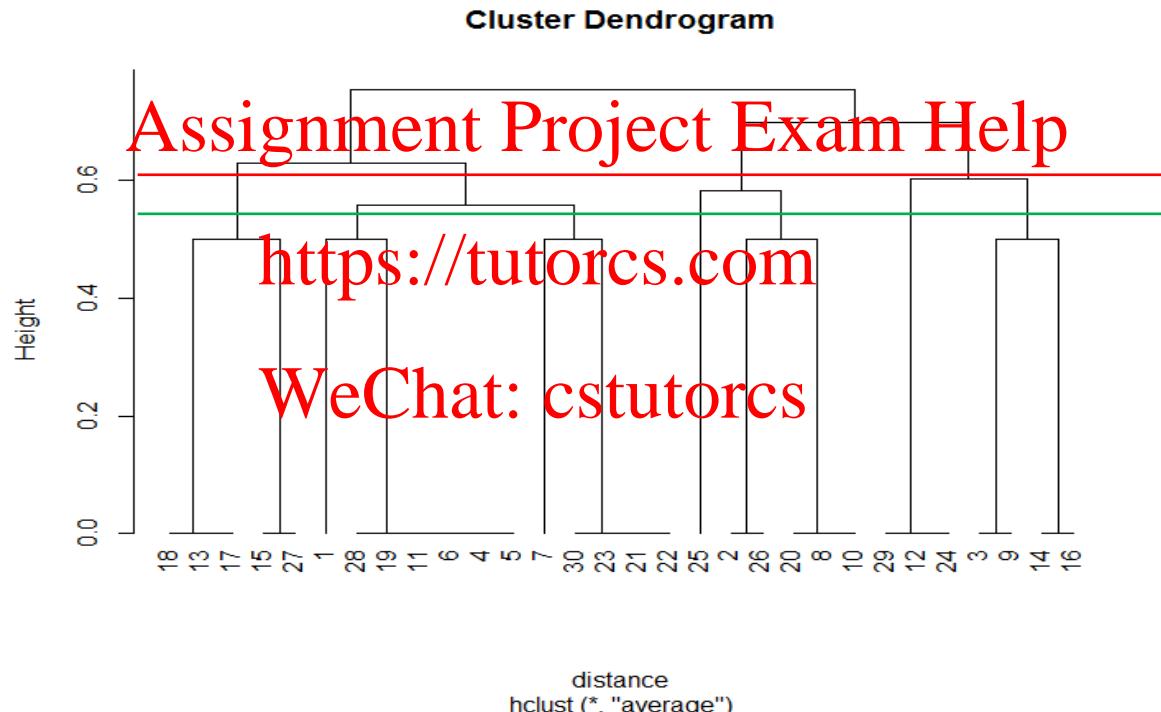


Dendrogram



- Visually, see what the possible numbers of clusters are by introducing a horizontal line
- **Assignment Project Exam Help**
<https://tutorcs.com>
WeChat: cstutorcs
- In the Hierarchical clustering workshop, we use R to select the clusters and then investigate cluster properties

Dendrogram 4 and 7 clusters



Distance

We need to decide **Assignment Project Exam Help**

- How to define the distance between two data points
 - Recall that the data may be of different types
- How to define the distance between two clusters
 - This will depend on the distance between two points, but once that distance is decided there are still many alternative ways of defining distance between clusters



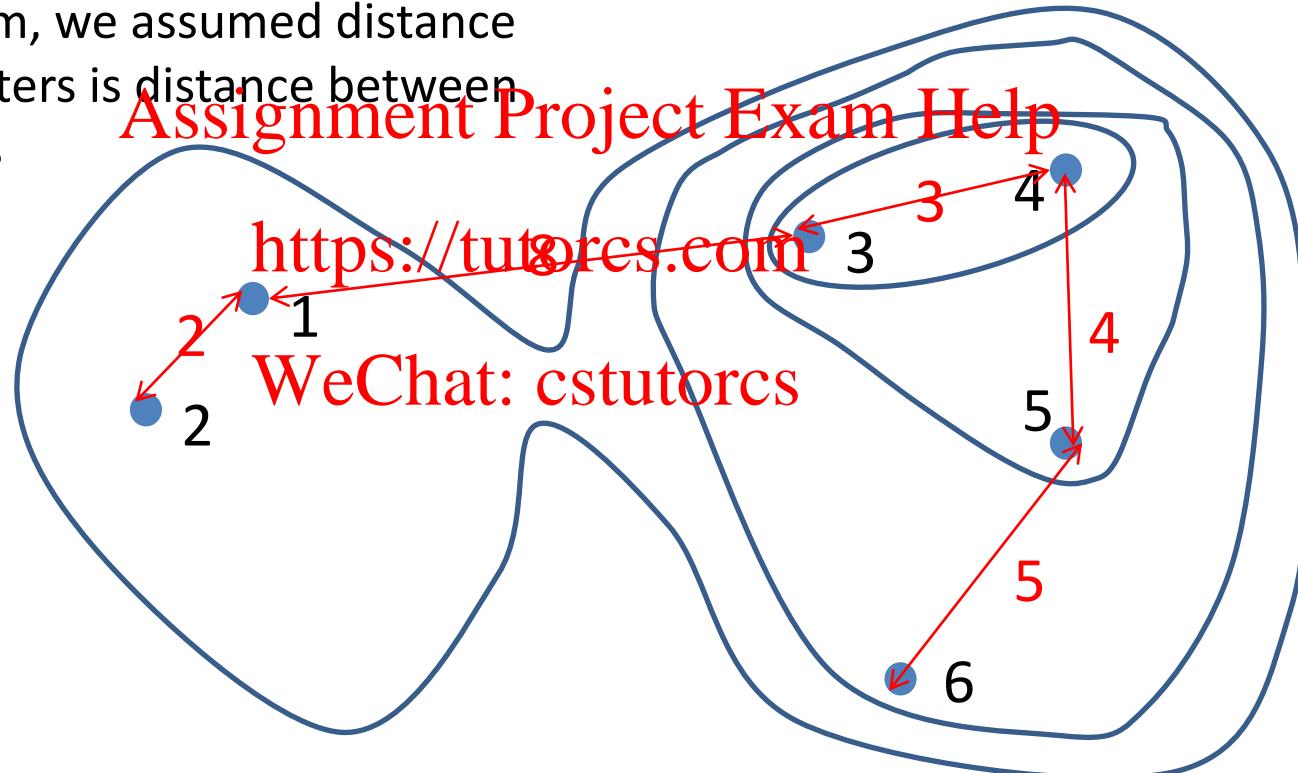
MONASH University

Distance

- In the example, we used
 - Euclidean distance (the usual shortest distance between two points)
 - between-cluster distance was taken to be the distance between the two points, one in each cluster, closest together
- When you have a data frame whose columns represent many different kinds of variable, defining distance between data points takes some thought.

Distance

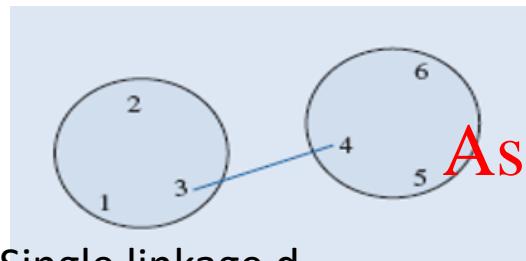
In this diagram, we assumed distance between clusters is distance between closest points



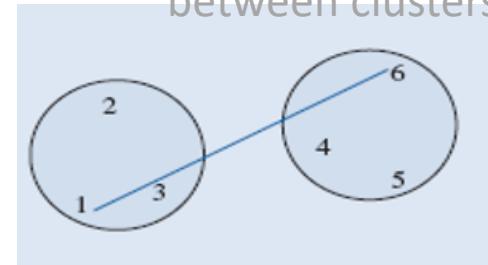
Distance

- Need to decide:
 - What are “distances” between points (=cases=rows)?
Assignment Project Exam Help
- And for the hierarchical method, in addition
 - What are “distances” between clusters?
https://tutores.com
- *Distance* is a dissimilarity measure (larger distance, less similar)
WeChat: cstutorcs
- For categorical variables, see how many attributes match: *Matching* is a similarity measure (more matching, more similar)
- Where a similarity measure s has been used, a corresponding “distance” d is defined to feed into the cluster algorithm:
$$d = \sqrt{1 - s}$$

Distance between Clusters



Single linkage $d_{3,4}$
shortest distance
between clusters

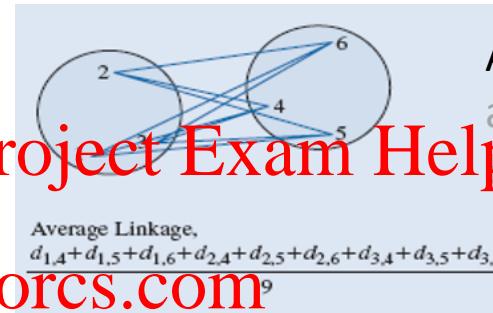


Complete linkage $d_{1,6}$
longest distance
between clusters

Assignment Project Exam Help

<https://tutorcs.com>

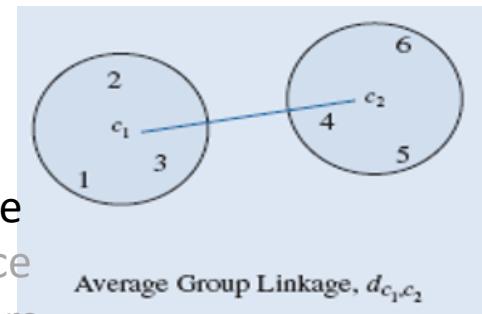
WeChat: cstutorcs



Average Linkage,
 $d_{1,4} + d_{1,5} + d_{1,6} + d_{2,4} + d_{2,5} + d_{2,6} + d_{3,4} + d_{3,5} + d_{3,6}$

Average linkage
average distance
between all linkage

$$d = \sqrt{1 - s}$$



Centroid linkage
average distance
between clusters

Average Group Linkage, d_{c_1, c_2}

Distance Matrix to Dendrogram

	A	B	C	D	E
A	0	2	8	10	13
B	2	0	10	15	14
C	8	10	0	5	5
D	10	15	1	0	4
E	13	14	5	4	0

Assignment Project Exam Help

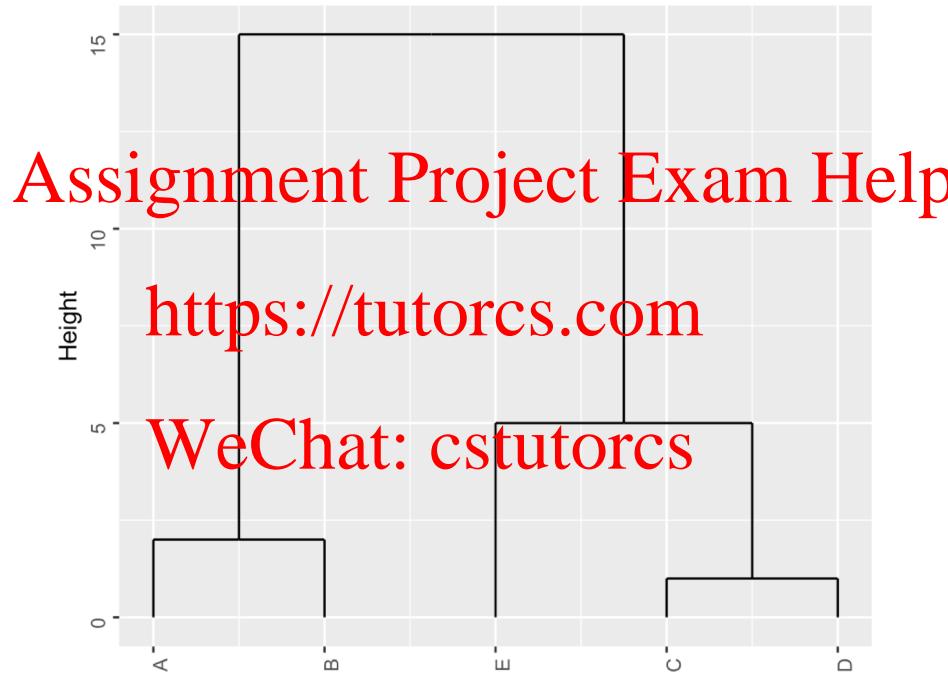
<https://tutorcs.com>

WeChat: cstutorcs

Given these distances,
create a dendrogram for

1. complete linkage
2. single linkage

Complete Linkage



Complete Linkage – Breaking Down Steps

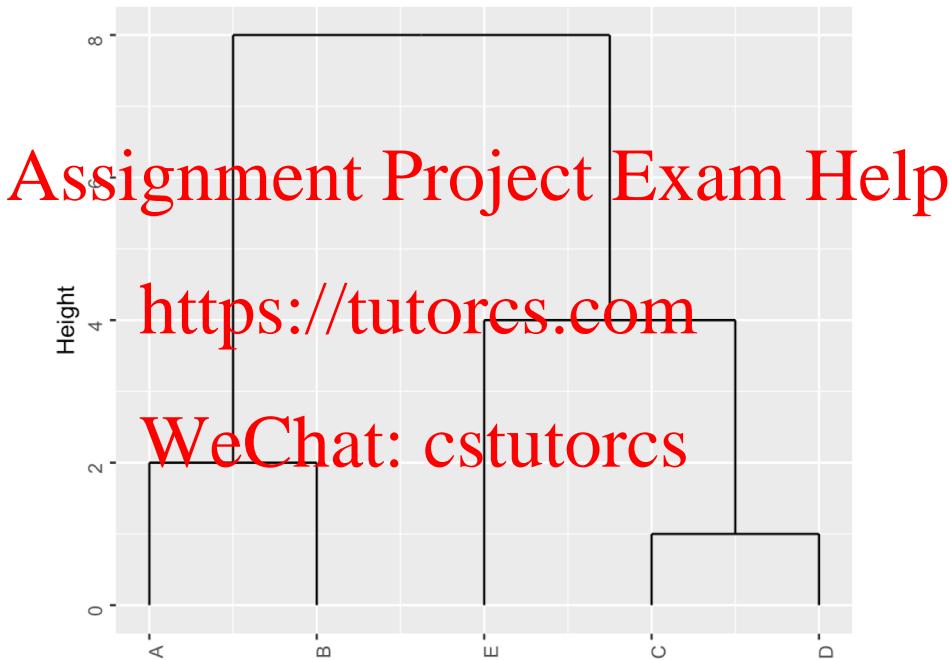
- Let's start with the smallest distance:
 1. C-D (dist 1): This is first cluster, recompute the similarities by considering C and D as one observation = 1
 2. A-B (dist 2): is the second cluster, recompute the similarities = 2
- For the next step, look for the larger distance between D-E or C-E. Which is longer? (D-E = 4), (C-D = 5). For complete linkage, look for maximum distance between cluster
 3. A new cluster C-D-E = 5
- Now look for maximum distance between A-B and C-D-E
 4. A-B to C-D-E = 15

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutors

Single Linkage



Single Linkage – Breaking Down Steps

- Let's start with the smallest distance:
 - C-D (dist 1): This is first cluster, recompute the similarities by considering C and D as one observation = 1
 - A-B (dist 2): is the second cluster, recompute the similarities = 2
- For the next step, look for the shorter distance between D-E or C-E. Which is shorter? (D-E = 4), (C-D = 5). For single linkage, look for minimum distance between cluster
 - A new cluster C-D-E = 4
- Now look for minimum distance between A-B and C-D-E
 - A-B to C-D-E = 8

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Assignment Project Exam Help

- A centroid is the average observation of a cluster
<https://tutorcs.com>
- That is, the average value for each variable across the points in the cluster
WeChat: cstutorcs
- The centroid doesn't have to be an actual observation

Hierarchical Cluster: Pros & Cons

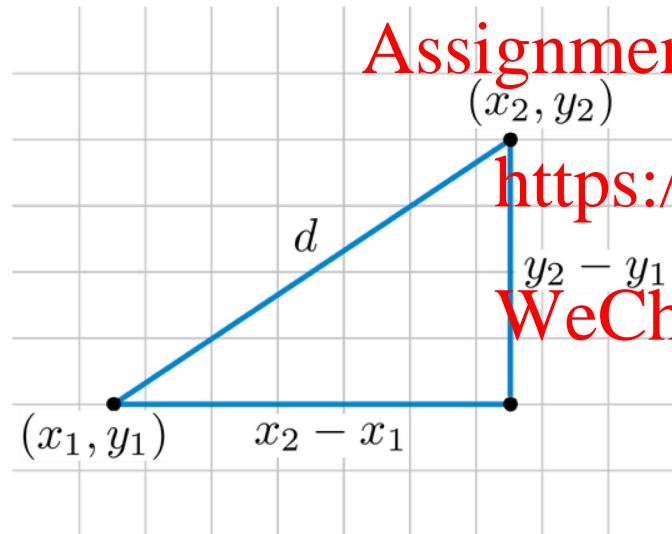
- Hierarchical clustering methods differ in the method of representing similarity between clusters, each with advantages and disadvantages

Assignment Project Exam Help https://tutorcs.com WeChat: cstutortutor			
Single linkage	Complete linkage	Average linkage	Centroid linkage
<ul style="list-style-type: none">Most versatile algorithmPoorly delineated cluster structures within data can produce unacceptable snakelike “chains” for clusters	<ul style="list-style-type: none">Eliminates the chaining problemOnly considers outermost observations in a cluster → impacted by outliers	<ul style="list-style-type: none">Based on average similarity of all individuals in a cluster and tends to generate clusters with small within-cluster variationless affected by outliers	<ul style="list-style-type: none">Measures distance between cluster centroidsLess affected by outliers

Numerical Distance

- For numerical variables, dissimilarity is measured by:
Assignment Project Exam Help
 - Euclidean distance
 - Manhattan distance
 - Many others
- https://tutorcs.com**
- WeChat: cstutorcs**

Numerical Distance



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$d_E = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Manhattan distance

$$d_M = |x_2 - x_1| + |y_2 - y_1|$$

Euclidean Distance

Euclidean distance: Most common method to measure **distance** between observations, when observations include continuous variables.

Let observations $u = (Age_u, Income_u, Kids_u)$ and

$$v = (Age_v, Income_v, Kids_v)$$

each comprise measurements of 3 variables.

The Euclidean distance between observations u and v is

$$d_{u,v} = \sqrt{(Age_u - Age_v)^2 + (Income_u - Income_v)^2 + (Kids_u - Kids_v)^2}$$

Manhattan Distance

Manhattan distance: Distance between observations that include continuous variables.

Assignment Project Exam Help

- Let observations u and v each comprise measurements of the variables *Age*, *Income*, and *Kids*.
- The Manhattan distance between observations u and v is

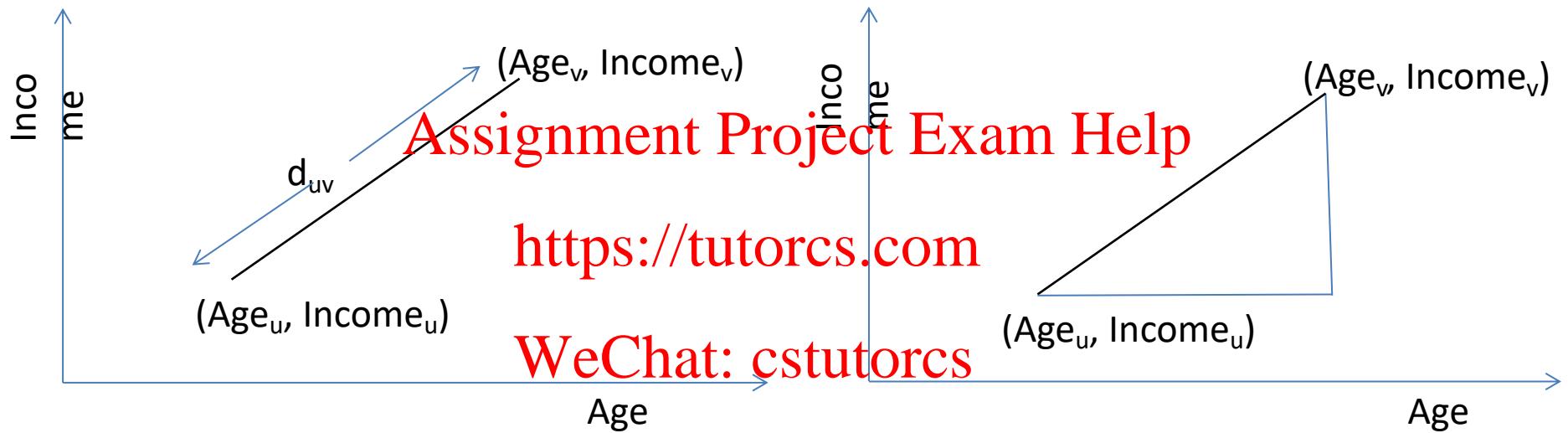
WeChat: cstutorcs

$$d_{u,v} =$$

$$|Age_u - Age_v| + |Income_u - Income_v| + \dots + |Kids_u - Kids_v|$$

Also need to standardise variables for Manhattan distance
Manhattan distance is more robust to outliers

Examples with Two Variables



Euclidean distance becomes smaller as observations become more similar with respect to their variable values. Same is true for Manhattan distance

Example: Know Thy Customer

- KTC is a Financial Advising Company
- Wishes to perform market segmentation of its clients. Thus each client is a case or “point”
 - *Assignment Project Exam Help*
 - *https://tutorcs.com*
 - *WeChat: cstutorcs*
 - *What type is each variable?*
 - *How should similarity between cases be measured in relation to values of*
 - *Numerical variables?*
 - *Categorical variables?*

Example: Know Thy Customer

Corresponding to each customer, KTC has a vector of measurements on seven customer variables: Age, Female, Income, Married, Children, Car Loan, Mortgage

Assignment Project Exam Help

Age	Female	Income	Married	Children	Car Loan	Mortgage
...
61	1	57881	1	2	0	1
27	0	55539	1	0	0	1

<https://tutorcs.com>

WeChat: cstutorcs

- 61-year-old female with an annual income of \$57,881, married with two children, no car loan and but a mortgage.
- 27-year-old male with annual income of \$55,539, married, no children, no car loan but mortgage

Start by considering only the binary variables: Female, Married, Car Loan, Mortgage

Measuring (Dis)Similarity Between *Observations*

- For binary variables, we discuss similarity measures by:

- Matching coefficient
- Jaccard's coefficient

Assignment Project Exam Help

<https://tutorcs.com>

- Many other measures are available

WeChat: cstutorcs

Similarity for Categorical Variables

Matching coefficient: Simplest overlap measure of **similarity** between two observations

Assignment Project Exam Help

- Find the proportion of variables with matching values
- Used when clustering observations solely on the basis of categorical variables encoded as 0–1
- Matching coefficient

WeChat: cstutorcs

$$s_M = \frac{\text{number of variables with matching value for observations } u \text{ and } v}{\text{total number of variables}}$$

Similarity for Categorical Variables

- Weakness of matching coefficient - If two observations both have a 0 entry for a categorical variable this is counted as a sign of similarity between the two observations
- However, matching 0 entries do not necessarily imply similarity
 - If 1 means “Owns a Toyota Land Cruiser” and 0 means “Does not own a Toyota Land Cruiser” then two observations sharing a zero does not indicate significant similarity

<https://tutorcs.com>

WeChat: cstutorcs

Similarity for Categorical Variables

Jaccard's coefficient: A similarity measure that does not count matching zero entries

Assignment Project Exam Help

$$S_J = \frac{\text{number of variables with matching } \textit{nonzero} \text{ value for observations } u \text{ and } v}{(\text{total number of variables}) - (\text{number of variables with matching zero values for observations } u \text{ and } v)}$$

Similarity for Categorical Variables

- For the following example, which coefficient should we use?

Assignment Project Exam Help

- Matching?
- Jaccard?

<https://tutorcs.com>

WeChat: cstutorcs

0 0 0 0 0 1 1 1 1 1

0 0 0 0 1 1 0 0 0 0

Similarity for Categorical Variables

	Female	Married	Car loan	Mortgage
U	1	1	0	1
Assignment	Project	Exam	Help	
V	0	1	0	1
W	https://tutorcs.com			1

Similarity between	WeChat: Jacob	Jacob:tutorcs	Matching
U and V			
V and W			
W and U			

Similarity for Categorical Variables

- Binary variables in KTC data:
 - Gender
 - Married or not
 - Car loan or not
 - Home mortgage or not
 - Given these variables, would you use Matching or Jaccard?
- Assignment Project Exam Help
<https://tutorcs.com>
WeChat: cstutorcs

Example: Know Thy Customer

Again, corresponding to each customer, KIT has a vector of measurements on seven customer variables: Age, Female, Income, Married, Children, Car Loan, Mortgage

<https://tutorcs.com>

Age	Female	Income	Married	Children	Car Loan	Mortgage
...
61	1	57881	1	2	0	1
27	0	55539	1	0	0	1

WeChat: cstutorcs

Euclidean Distance

$$d_{1,2} = \sqrt{(61-27)^2 + (57881-55539)^2 + (2-0)^2}$$

Assignment Project Exam Help

$$\begin{aligned}
 &= \sqrt{1156 + 5486124 + 4} \\
 &= \sqrt{5486124} \\
 &= 2342
 \end{aligned}$$

ID	Age	Income	Children
...
1	61	57881	2
2	27	55539	0

WeChat: cstutorcs

What's the problem here?

Manhattan Distance

$$d_{1,2}^{(M)} = |61 - 27| + |57881 - 55539| + |2 - 0|$$

Assignment Project Exam Help

$$= 34 + 2342 + 2$$

<https://tutorcs.com>

ID	Age	Income	Children

1	61	57881	2
2	27	55539	0

WeChat: cstutorcs

The differences of scale are less, but the variables age and number of children are still swamped by Income

Scaling Distances

- Euclidean and Manhattan distance are highly influenced by the scale on which variables are measured.
 - Therefore, it is common to standardize the units of each variable of each observation u ;
 - Example: Income_u , the value of variable Income in observation u , is replaced with its z-score.
- The conversion to z-scores also makes it easier to identify outlier measurements, which can distort the Euclidean distance between observations.
- Recall: What are z-scores?

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: astutiores

Z-Scores

Consider values of a variable x x_1, x_2, \dots, x_n

Assignment Project Exam Help

Suppose the mean and standard deviation of this set of value is given by \bar{x} and s

Then the z - score corresponding to x_i is given by

$$z_i = \frac{x_i - \bar{x}}{s}$$

If certain variables are known to be more important, deliberately choose unequal weighting

Variables of Mixed Types

- A database may contain all types of variables
Assignment Project Exam Help
- A weighted formula is used to combine the effects of numerical and categorical variables
<https://tutorcs.com>
 - [See the *gower* metric in the command *daisy* in the package *cluster*]
- We focus on cluster analysis using just one type of variable
WeChat: cstutorcs

Limitations of Hierarchical Clustering

- Need to calculate and store an $n \times n$ distance matrix; infeasible for large datasets
 - Low stability: change a few records and a very different result may be obtained
 - Influenced by choice of metric, especially when using “average”
 - Sensitive to outliers
 - Only one pass through the data: misallocations can never be fixed
- Assignment Project Exam Help
<https://tutorcs.com>
WeChat: cstutorcs

Interpreting, Profiling, & Validating Clusters

- The clusters centroid, a mean profile of the cluster on each clustering variable, is particularly useful in the interpretation stage
 - Interpretation involves examining the distinguishing characteristics of each cluster's profile and identifying substantial differences between clusters
 - If there is not substantial difference between clusters, then other cluster solutions should be examined
 - The cluster centroid should also be assessed for correspondence with the researcher's prior expectations based on theory or practical experience

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Interpreting, Profiling, & Validating Clusters

- Validation is essential in cluster analysis since the clusters are descriptive of structure and require additional support for their relevance:
 - Cross-validation empirically validates a cluster solution by creating two sub-samples (randomly splitting the sample) and then comparing the two cluster solutions for consistency with respect to number of clusters and the cluster profiles.
 - Validation is also achieved by examining differences on variables not included in the cluster analysis but for which there is a theoretical and relevant reason to expect variation across the clusters

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs