

---

# COMP9318: Data Warehousing and Data Mining

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

- 
- What is Cluster Analysis?

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# What is Cluster Analysis?

---

- Cluster: a collection of data objects
  - Similar to one another within the same cluster
  - Dissimilar to the objects in other clusters
- Cluster analysis <https://eduassistpro.github.io/>
  - Grouping a s lusters
- Clustering *belongs to unsuper fication*: no predefined classes
- Typical applications
  - As a **stand-alone tool** to get insight into data distribution
  - As a **preprocessing step** for other algorithms

# General Applications of Clustering

---

- Pattern Recognition
- Spatial Data Analysis
  - create thematic maps in GIS by clustering feature spaces
  - detect spatial data mining <https://eduassistpro.github.io/>
- Image Processing [Add WeChat edu\\_assist\\_pro](#)
- Economic Science (especially market research)
- WWW
  - Document classification
  - Cluster Weblog data to discover groups of similar access patterns

# Examples of Clustering Applications

---

- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- Land use: Identify land use in an earth observatory
- Insurance: Identifying groups of insurance policy holders with a high average claim cost
- City-planning: Identifying groups of houses according to their house type, value, and geographical location
- Earthquake studies: Observed earth quake epicenters should be clustered along continent faults

# What Is Good Clustering?

---

- A good clustering method will produce high quality clusters with
  - high intra-cl
  - low inter-cl
- The quality of a clustering method both the similarity measure used by the method and its implementation.
- The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns.

# Requirements of Clustering in Data Mining

---

- Scalability
- Ability to deal with different types of attributes
- Discovery of clusters with arbitrary shape
- Minimal requirement to determine input
- Able to deal with noise and outliers
- Insensitive to order of input records
- High dimensionality
- Incorporation of user-specified constraints
- Interpretability and usability

# Chapter 8. Cluster Analysis

---

- Preliminaries

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Typical Inputs

Key component for clustering:  
the dissimilarity/similarity  
metric:  $d(i, j)$

- Data matrix
  - N objects, each represented by a m-dimensional vector
- Dissimilarity matrix
  - A square matrix giving distances between all pairs of objects.
  - If similarity functions are used → similarity matrix

Assignment Project Exam Help

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

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1m} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{im} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{nm} \end{bmatrix}$$

$n \times m$

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

$n \times n$

# Comments

---

- The definitions of distance functions are usually very different for interval-scaled, boolean, categorical, ordinal and ratio variables.
- Weights should be appropriate based on application variables  
https://eduassistpro.github.io/
- There is a separate “quality” function that measures the “goodness” of a cluster.
- It is hard to define “similar enough” or “good enough”
  - the answer is typically highly subjective.

# Type of data in clustering analysis

---

- Interval-scaled variables:
- Binary variables: Assignment Project Exam Help
- Nominal, ordinal <https://eduassistpro.github.io/>
- Variables of missing values Add WeChat edu\_assist\_pro

# Interval-valued variables

---

- Standardize data
  - Calculate the *mean absolute deviation*:

**Assignment Project Exam Help**

$$s_f = \frac{1}{n}(|x_{if} - m_f| + |x_{if} - m_f| + \dots + |x_{nf} - m_f|)$$

where

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

**Add WeChat  $m_f$   $\bar{x}_f$   $s_f^2$   $x_{if}$   $x_{nf}$**

- Calculate the standardized measurement (*z-score*)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- Using mean absolute deviation is more **robust** than using standard deviation

# Similarity and Dissimilarity Between Objects

---

- Distances are normally used to measure the similarity or dissimilarity between two data objects
- A popular choice is the *Minkowski distance, or the  $L_p$  norm of difference* <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

- Special cases:
  - if  $p = 1$ ,  $d$  is the **Manhattan distance**
  - if  $p = 2$ ,  $d$  is the **Euclidean distance**
  - if  $p = \infty$ ,  $\|\mathbf{x}_i - \mathbf{x}_j\|_\infty = \max_{k=1}^m |\mathbf{x}_{ik} - \mathbf{x}_{jk}|$

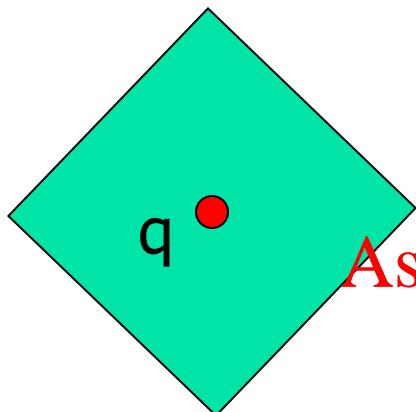
# Similarity and Dissimilarity Between Objects (Cont.)

- Other similarity/distance functions:
  - Mahalanobis distance
  - Jaccard, Dice, cosine similarity, Pearson correlation coefficient  
<https://eduassistpro.github.io/>
- Metric distance
  - Properties
    - $d(i,j) \geq 0$
    - $d(i,i) = 0$
    - $d(i,j) = d(j,i)$
    - $d(i,j) \leq d(i,k) + d(k,j)$
  - Add WeChat [edu\\_assist\\_pro](#)  
to all distance functions

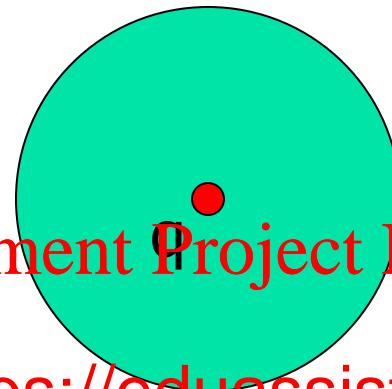
positiveness  
symmetry  
reflexivity

triangular inequality

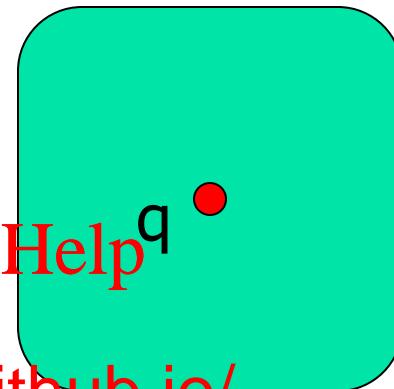
# Areas within a unit distance from q under different $L_p$ distances



$L_1$



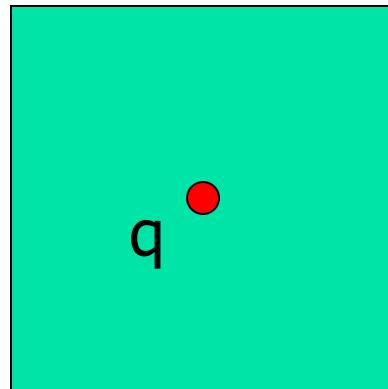
Assignment Project Exam Help



$q$

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

Add WeChat  $edu\_assist\_pro$



$L_\infty$

# Binary Variables

Obj	Vector Representation
i	[0, 1, 0, 1, 0, 0, 1, 0]
j	[0, 0, 0, 0, 1, 0, 1, 1]

- A contingency table for binary data

		Object j		
		1	0	sum
Object i	Assignment	a	b	a+b
	Project	c	d	c+d
		sum	a+c	b+d
		a+b+c+d	Add WeChat	edu_assist_pro

- Simple matching coefficient (invariant, if the binary variable is *symmetric*):  $d(i, j) = \frac{b+c}{a+b+c+d}$
- Jaccard coefficient (noninvariant if the binary variable is *asymmetric*):  $d(i, j) = \frac{b+c}{a+b+c}$

# Dissimilarity between Binary Variables

$$d(i, j) = \frac{b+c}{a+b+c}$$

## Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	A	S	s	P	E	X
Mary	F					P	N
Jim	M					N	N

- gender is a symmetric attribute
- the remaining attributes are asymmetric
- let the values Y and P be set to 1, and the value N be set to 0

$$d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$$

$$d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$$

$$d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$$

# Nominal Variables

---

- A generalization of the binary variable in that it can take more than 2 states, e.g., red, yellow, blue, green
- Method 1: Simple matching
  - $m$ : # of mat <https://eduassistpro.github.io/>
- Add WeChat  $d(i,j) = \frac{1}{M}$
- Method 2: One-hot encoding
  - creating a new binary variable for each of the  $M$  nominal states

# Ordinal Variables

---

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank
- Can be treated like interval-scaled
  - replace  $x_{if}$  by  $\frac{r_{if} - 1}{M_f - 1}$
  - map the range of each variable  $[1, M_f]$  by replacing  $i$ -th object in the  $f$ -th variable

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

- compute the dissimilarity using methods for interval-scaled variables

# Ratio-Scaled Variables

---

- Ratio-scaled variable: a positive measurement on a nonlinear scale, approximately at exponential scale, such as ~~Assignment Project Exam Help~~
- Methods: <https://eduassistpro.github.io/>
  - treat them like interval-scales—*not a good choice!* (why?—the scale is ~~not a good choice!~~)
  - apply logarithmic transformation
$$y_{if} = \log(x_{if})$$
  - treat them as continuous ordinal data treat their rank as interval-scaled

- 
- A Categorization of Major Clustering Methods

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Major Clustering Approaches

---

- Partitioning algorithms: Construct various partitions and then evaluate them by some criterion
- Hierarchy algorithms: Create a hierarchical decomposition of the set of data <https://eduassistpro.github.io/>
- Graph-based algorithms: Specified
- Density-based: based on connectivity and density functions
- Grid-based: based on a multiple-level granularity structure
- Model-based: A model is hypothesized for each of the clusters and the idea is to find the best fit of that model to each other

---

- Partitioning Methods

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Partitioning Algorithms: Problem Definition

---

- Partitioning method: Construct a “**good**” partition of a database of  $n$  objects into a set of  $k$  clusters
  - Input: a  $n \times m$  data matrix
- Assignment Project Exam Help
- How to measure given partitioning scheme?  
<https://eduassistpro.github.io/>
- Cost of a cluster,  $\text{cost}(C_i) = \sum_{x_j \in C_i} \|x_j - \text{center}(C_i)\|_2^2$ 
  - Note:  $L_2$  distance used
  - Analogy with binning?
  - How to choose the center of a cluster?
    - Centroid (i.e., Avg) of  $x_j \rightarrow$  Minimizes  $\text{cost}(C_i)$
- Cost of  $k$  clusters: sum of  $\text{cost}(C_i)$

# Example (2D)

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Partitioning Algorithms: Basic Concept

---

- It's an optimization problem!
  - Global optimal:
    - NP-hard (for a wide range of cost functions)
    - Requires ex  $n^n \} = \Theta\left(\frac{k^n}{k!}\right)$  partitions
      - Stirling n <https://eduassistpro.github.io/>
  - Heuristic methods:
    - k-means: an instance of the ion-maximization)
    - Many variants

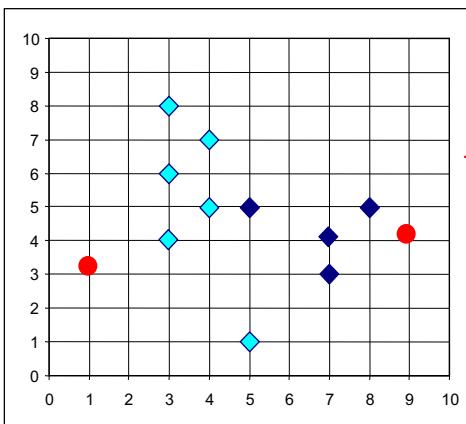
# The *K-Means* Clustering Method

---

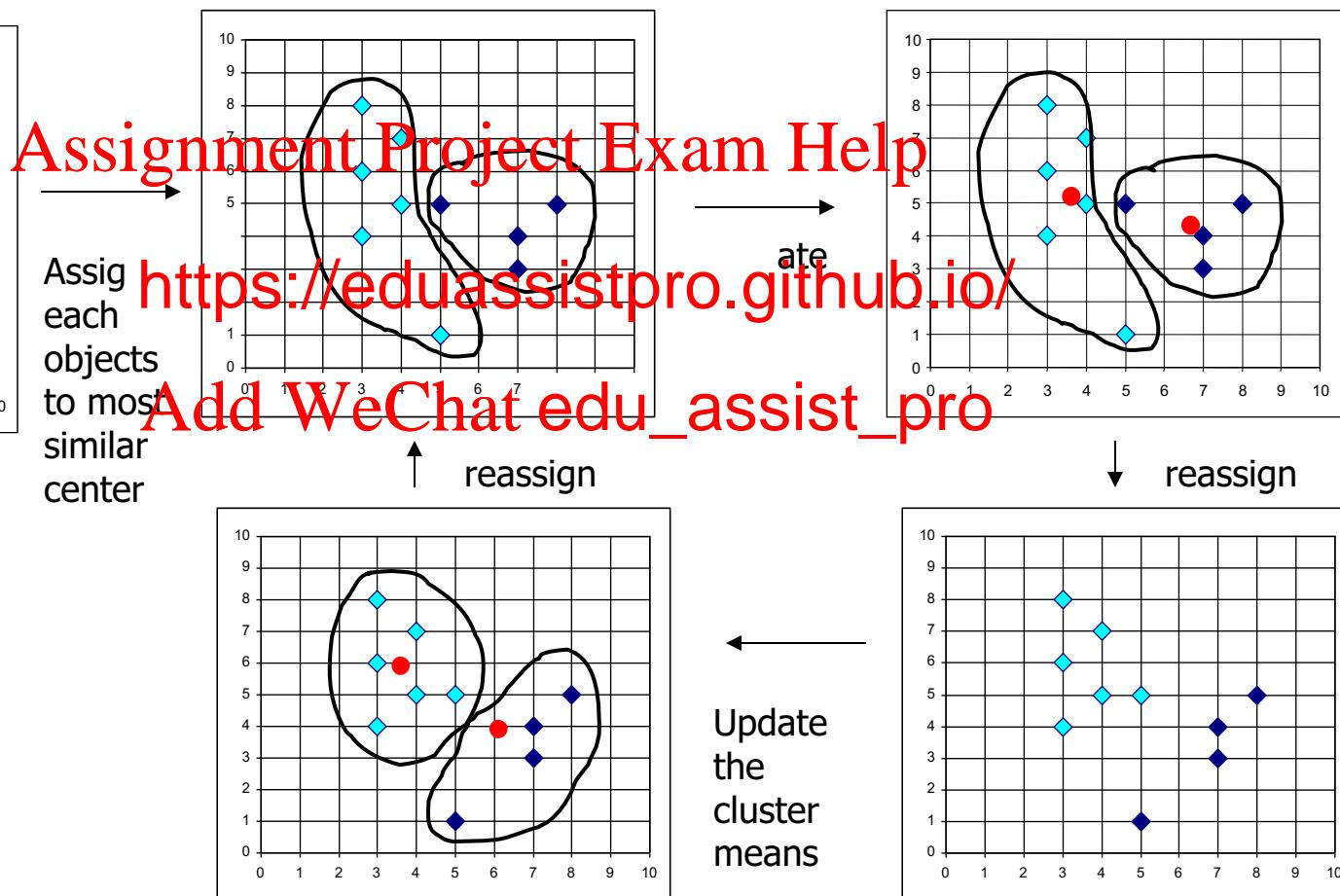
- Lloyds Algorithm:
  1. Initialize  $k$  centers randomly  
[Assignment](#) [Project](#) [Exam](#) [Help](#)
  2. While stop   
     i. Assign <https://eduassistpro.github.io/> ith the nearest et  
         center [Add WeChat edu\\_assist\\_pro](#)  
     ii. Compute the new center for each cluster.
- Stopping condition =?
- What are the final clusters?

# The *K*-Means Clustering Method

## ■ Example



Arbitrarily choose  $K$  object as initial cluster center



# Comments on the *K-Means* Method

- Strength: *Relatively efficient:*  $\mathcal{O}(tkn)$ , where  $n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations. Normally,  $k, t \ll n$ .
  - Comparing: PAM:  $O(k(n-k)2)$ , CLARA:  $O(ks2 + k(n-k))$
- Comment: **Assignment Project Exam Help**
  - Often terminates at optimum may be found using techniques like annealing and genetic algorithms
  - No guarantee on the quality. Use [Add WeChat edu\\_assist\\_pro](https://eduassistpro.github.io/)
- Weakness
  - Applicable only when *mean* is defined, then what about categorical data?
  - Need to specify  $k$ , the *number* of clusters, in advance
  - Unable to handle noisy data and *outliers*
  - Not suitable to discover clusters with *non-convex shapes*

# Variations of the *K-Means* Method

---

- A few variants of the *k-means* which differ in
  - Selection of the initial *k* means
  - Dissimilarity ca  
<https://eduassistpro.github.io/>
  - Strategies to c
- Handling categorical data: *k-modes*
  - Replacing means of clusters with modes
  - Using new dissimilarity measures to deal with categorical objects
  - Using a frequency-based method to update modes of clusters
  - A mixture of categorical and numerical data: *k-prototype* method

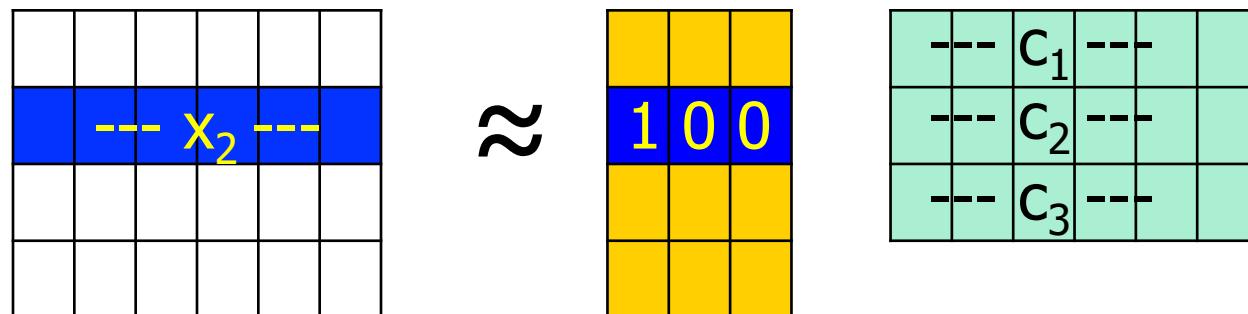
# k-Means++ [Arthur and Vassilvitskii, SODA 2007]

---

- A simple initialization routine that guarantees to find a solution that is  $O(\log k)$  competitive to the optimal  $k$ -means solution.  
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- Algorithm:
  1. Find first  $c_1$  <https://eduassistpro.github.io/>
  2. For each data point  $x_i$ , compute the distance to its nearest center
  3. Randomly sample one point as the new center, with probabilities proportional to  $D^2(x)$
  4. Goto 2 if less than  $k$  centers
  5. Run the normal  $k$ -means with the  $k$  centers

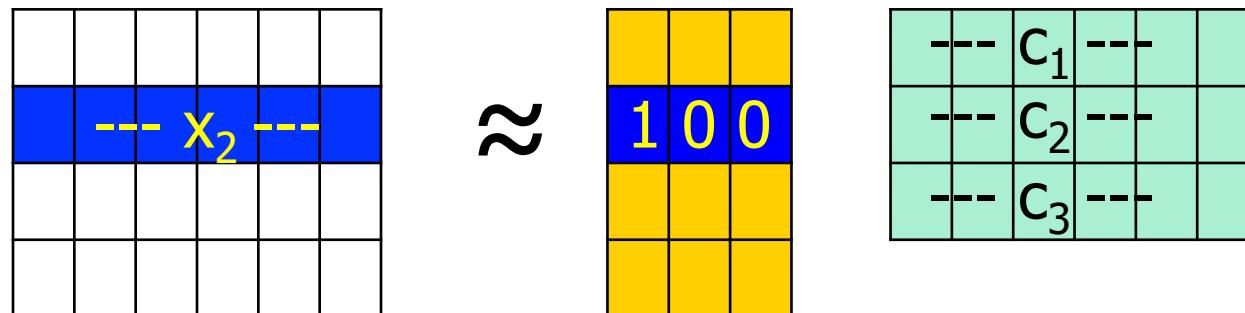
# k-means: Special Matrix Factorization

- $X^{n \times d} \approx U^{n \times k} V^{k \times d}$
- **Loss function:**  $\| X - UV \|_F^2$ 
  - Squared Frobenius norm
- **Constraints:** Assignment Project Exam Help
  - Rows of  $U$  m <https://eduassistpro.github.io/>
- Alternative view
  - $X_{j,*} \approx U_{j,*} V \rightarrow X_{j,*}$  can be expressed as a "special" linear combination of rows in  $V$



# Expectation Maximization Algorithm

- $X^{n \times d} \approx U^{n \times k} V^{k \times d}$
- **Loss function:**  $\| X - UV \|_F^2$
- Finding the best  $U$  and  $V$  simultaneously is hard, but  
**Assignment Project Exam Help**
- Expectation step:
  - Given  $V$ , find <https://eduassistpro.github.io/>
- Maximization step:
  - Given  $U$ , find the best  $V \rightarrow$  **Add WeChat edu\_assist\_pro**
- Iterate until converging at a local minima.

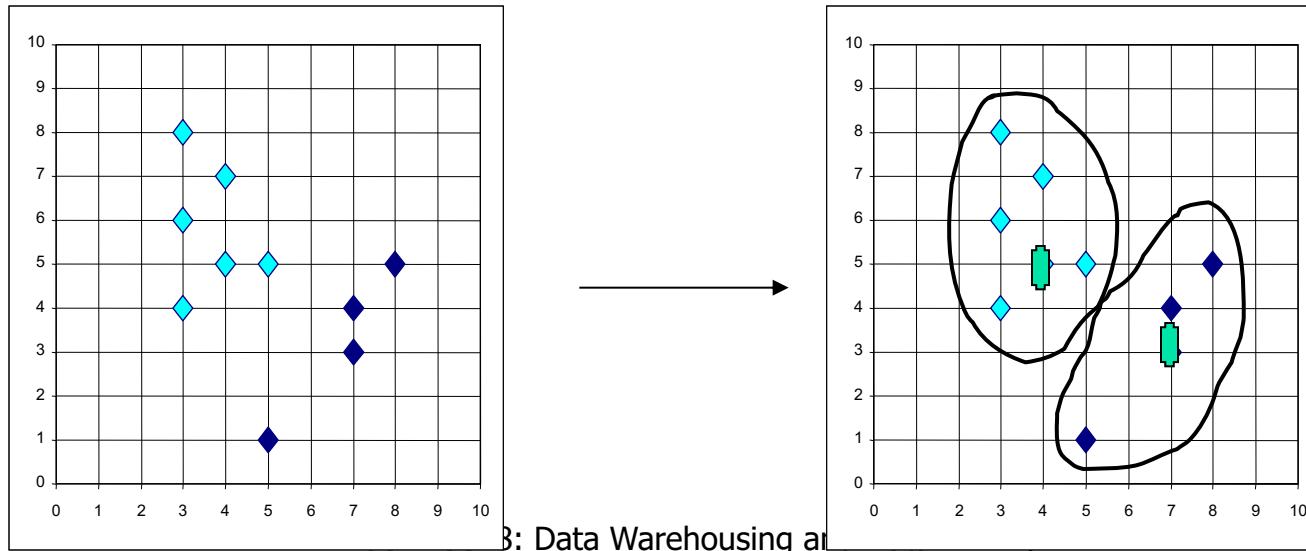


# What is the problem of k-Means Method?

- The k-means algorithm is sensitive to outliers !
  - Since an object with an extremely large value may substantially distort the distribution of the data.
- K-Medoids: Instead of the object in a cluster as a reference, used, which is the **most centrally located** object in a cluster.  
<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)



# K-medoids (PAM)

---

- *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by **one** of the objects in the cluster

**Assignment Project Exam Help**

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

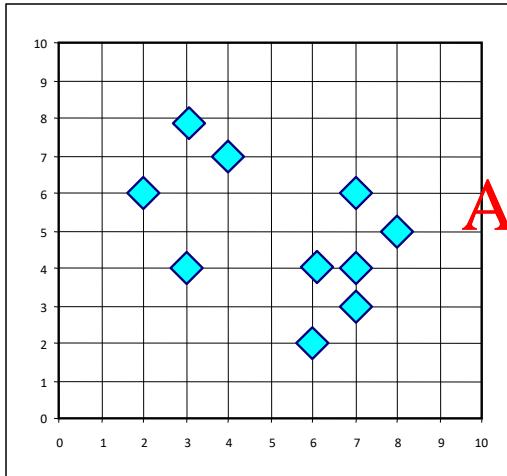
Add WeChat **edu\_assist\_pro**

# The *K-Medoids* Clustering Method

---

- Find *representative* objects, called medoids, in clusters
- *PAM* (Partitioning Around Medoids, 1987)
  - starts from an initial set of medoids and iteratively replaces one of the medoids if it improves the total distance of the resulting
- *PAM* works effectively for small data sets but does not scale well for large data sets
- *CLARA* (Kaufmann & Rousseeuw, 1990)
- *CLARANS* (Ng & Han, 1994): Randomized sampling
- Focusing + spatial data structure (Ester et al., 1995)

# Typical k-medoids algorithm (PAM)



K=2

**Do loop**  
**Until no change**

Arbitrary choose k object as init me

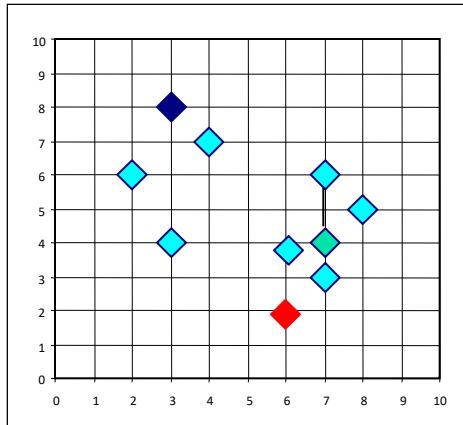
Assignment Project Exam Help

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

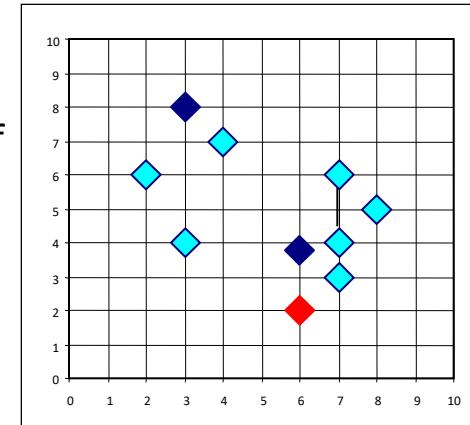
Add WeChat edu\_assist\_pro

Total Cost = 26

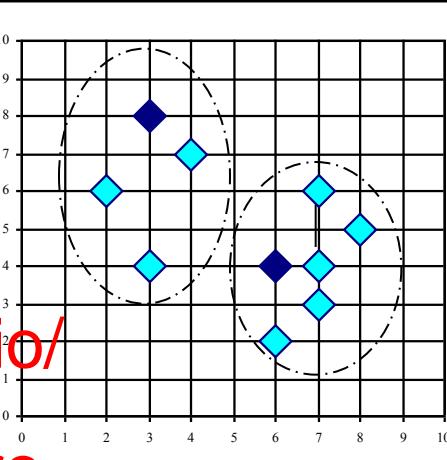
Swapping O and O<sub>a</sub>  
If quality is improved.



Compute total cost of swapping



Total Cost = 20



For each nonmedoid object, O<sub>a</sub>

# PAM (Partitioning Around Medoids) (1987)

---

- PAM (Kaufman and Rousseeuw, 1987), built in Splus
- Use real object to represent the cluster
  - Select  $k$  representative objects arbitrarily
  - For each pair  $i$  and  $h$ , calculate the total cost  $TC_{ih}$
  - For each pair of  $i$  and  $h$ ,
    - If  $TC_{ih} < 0$ ,  $i$  is replaced by  $h$
    - Then assign each non-selected object to the most similar representative object
  - repeat steps 2-3 until there is no change

# What is the problem with PAM?

---

- PAM is more robust than k-means in the presence of noise and outliers because a medoid is less influenced by outliers or ~~Assignment Projects Exam Help~~
- PAM works effi ~~s but does not scale well for~~  
■  $O(k(n-k)^2)$  for each iteration  
~~Add WeChat edu\_assist\_pro~~

where n is # of data,k is # of clusters

→ Sampling based method,

CLARA(Clustering LARge Applications)

# Gaussian Mixture Model for Clustering

---

- $k$ -means can be deemed as a special case of the EM algorithm for GMM
- GMM Assignment Project Exam Help  
<https://eduassistpro.github.io/>
  - allows “soft” clustering  $t$ :
    - model  $\Pr(t|z)$
  - also a good example of:
    - Generative model
    - Latent variable model
  - Use the Expectation-Maximization (EM) algorithm to obtain a local optimal solution

---

- Hierarchical Methods

Assignment Project Exam Help

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

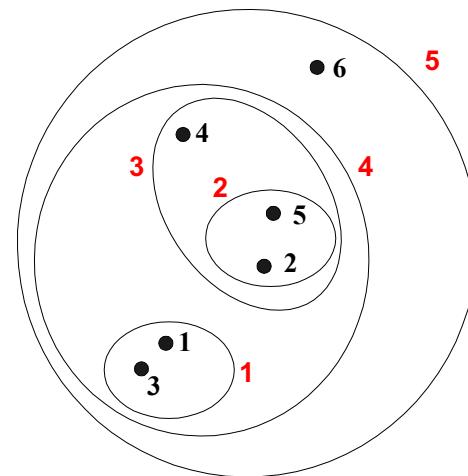
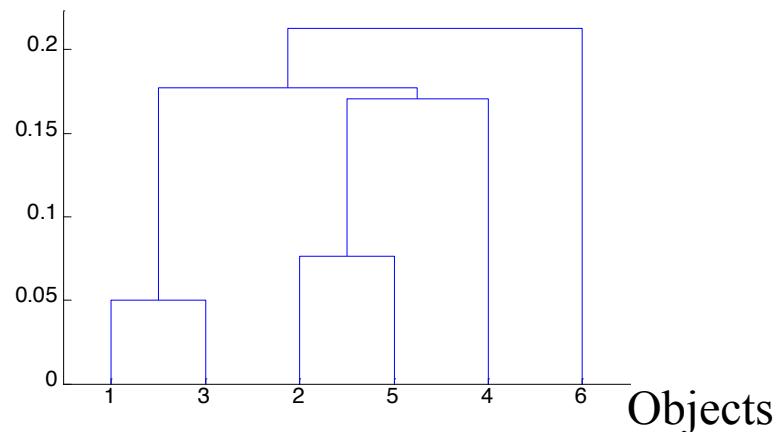
Add WeChat edu\_assist\_pro

# Hierarchical Clustering

- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a **dendrogram**
  - A tree like diagram that records the sequences of merges or splits
  - A clustering of the data objects is obtained by cutting the dendrogram at connected component forms a cluster <https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

Distance



# Strengths of Hierarchical Clustering

---

- Do not have to assume any particular number of clusters
  - Any desired number of clusters can be obtained by 'cutting' the dendrogram at the proper level <https://eduassistpro.github.io/>
- They may correspond to taxonomies
  - Example in biological sciences (e.g., animal kingdom, **phylogeny** reconstruction, ...)

# Hierarchical Clustering

---

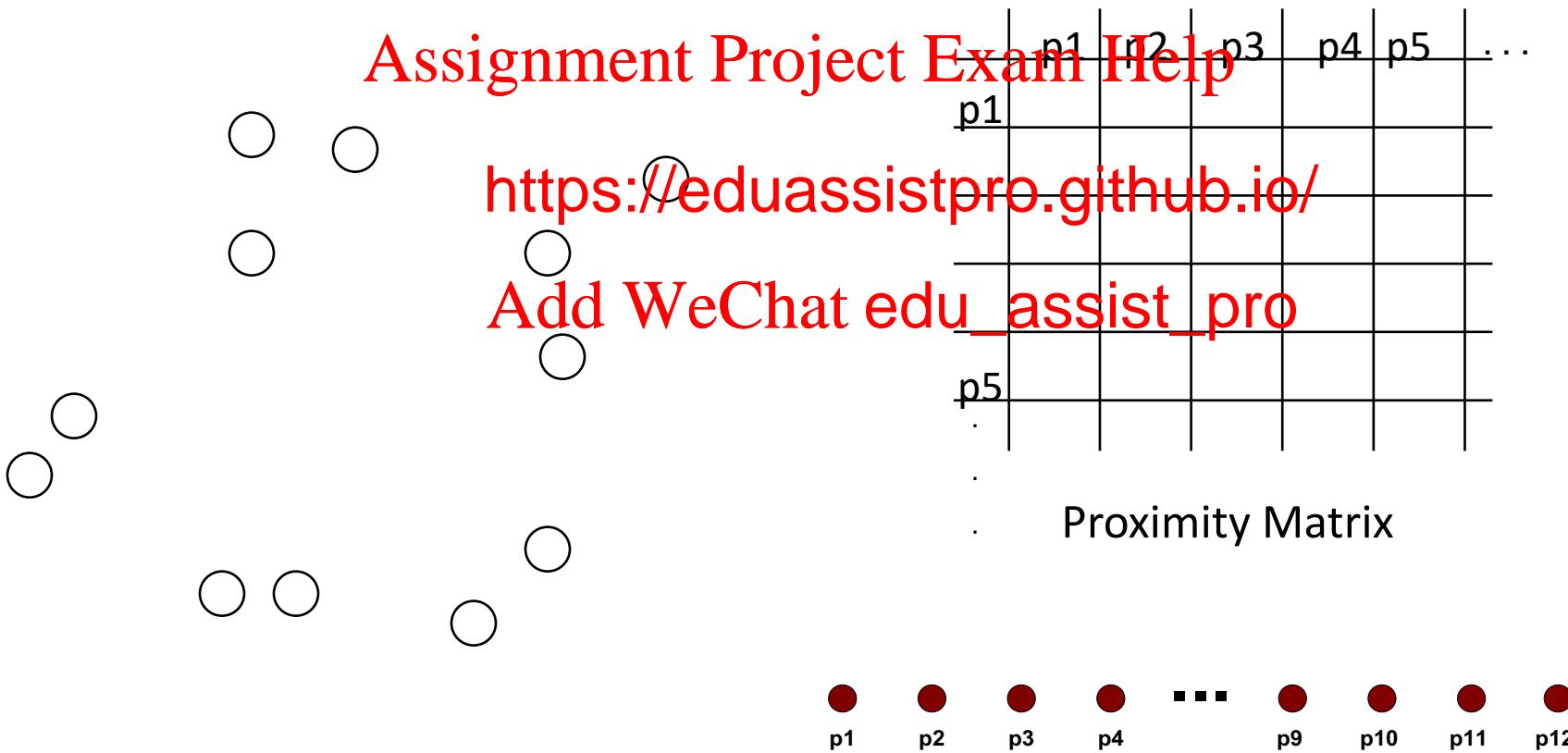
- Two main types of hierarchical clustering
  - Agglomerative:
    - Start with the points as individual clusters
    - At each step, **merge the closest pair of clusters** until only one cluster (or K clusters) left
  - Divisive: <https://eduassistpro.github.io/>
    - Start with one, all-inclusive cluster
    - At each step, split a cluster until it contains a point (or there are k clusters)
- Traditional hierarchical algorithms use a similarity or distance matrix
  - Merge or split one cluster at a time

# Agglomerative Clustering Algorithm

- More popular hierarchical clustering technique
- Basic algorithm is straightforward
  1. Compute the proximity matrix (i.e., matrix of pair-wise distances)
  2. Let each data point be a cluster
  3. Repeat
    - 4. Merge two clusters (Me <https://eduassistpro.github.io/>)
    - 5. Update the proximity
  6. Until only a single cluster remains
- Key operation is the computation of the proximity of two **clusters** ← different from that of two **points**
  - Different approaches to defining the distance between clusters distinguish the different algorithms

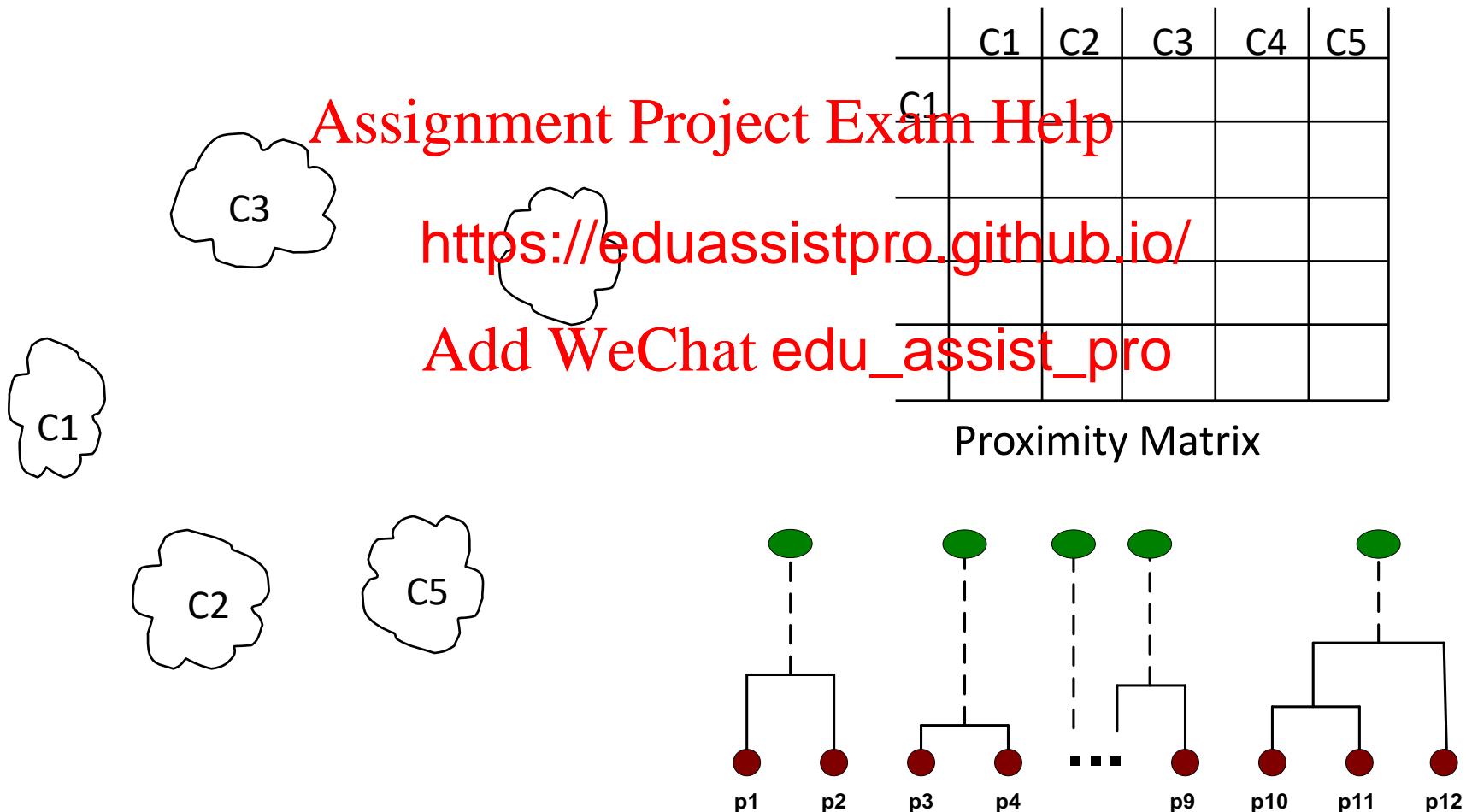
# Starting Situation

- Start with clusters of individual points and a proximity matrix



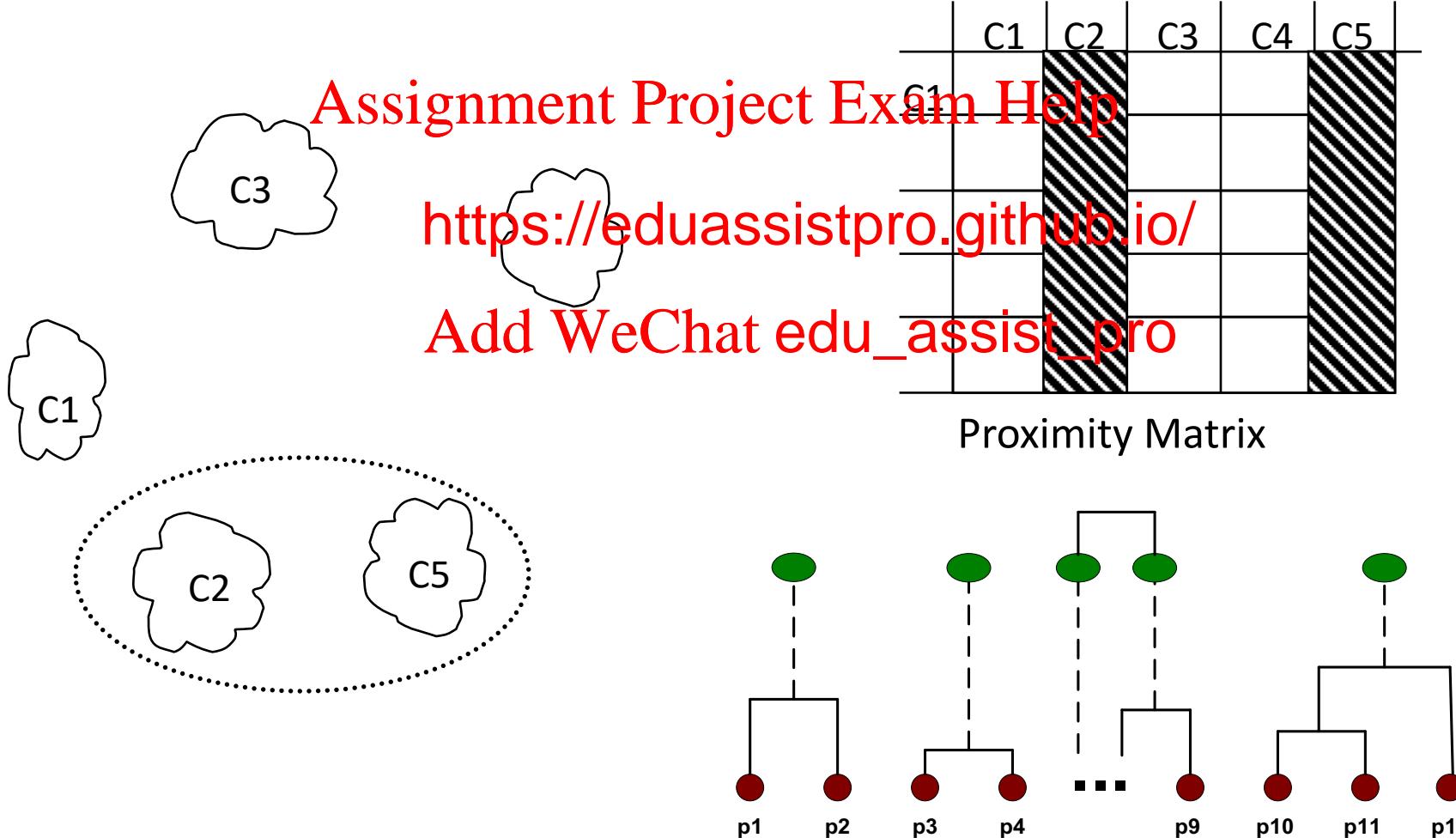
# Intermediate Situation

- After some merging steps, we have some clusters



# Intermediate Situation

- We want to merge the two closest clusters (C2 and C5) and update the proximity matrix.



# After Merging

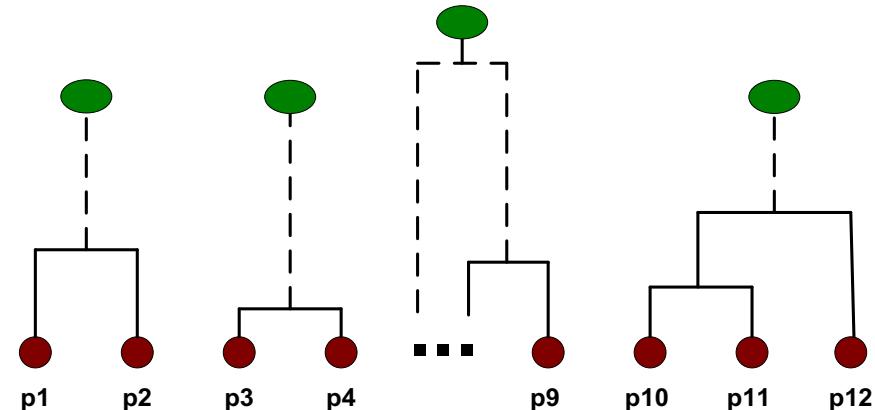
- The question is "How do we update the proximity matrix?"

Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

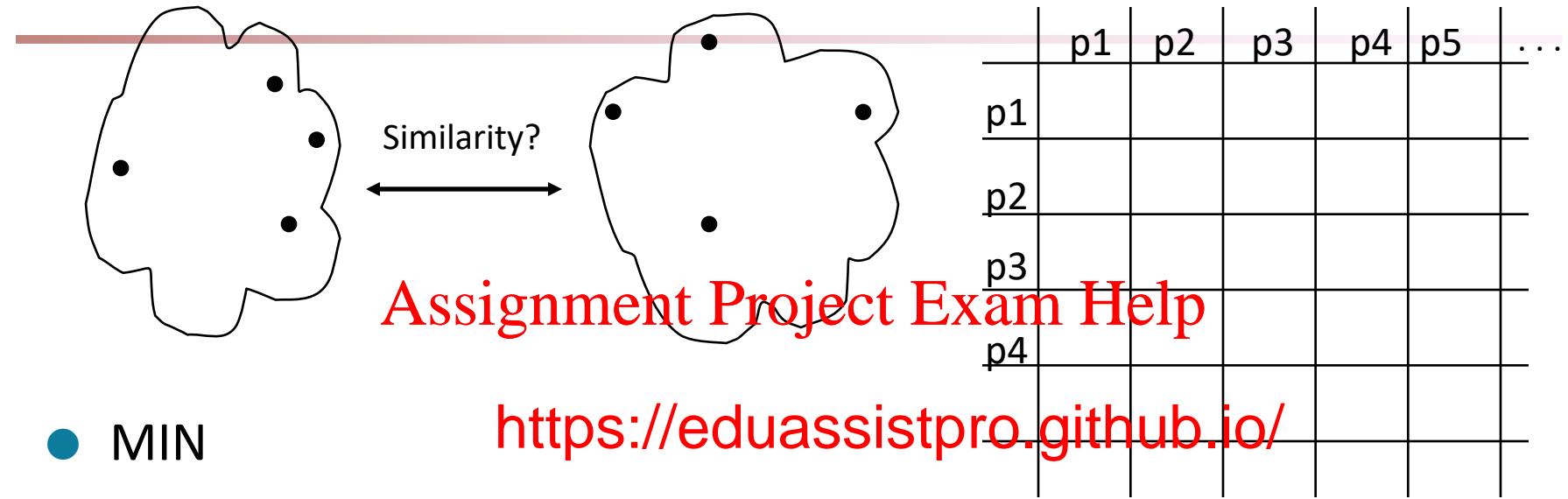
	C1	C2	U	C3	C4
C1	?	?	?	?	?
U	?	?	?	?	?
C3	?	?	?	?	?
C4	?	?	?	?	?

C1

C2 U C5



# How to Define Inter-Cluster Distance



- MIN
- MAX
- Centroid-based
- **Group Average**
- Other methods driven by an objective function
  - Ward's Method uses squared error

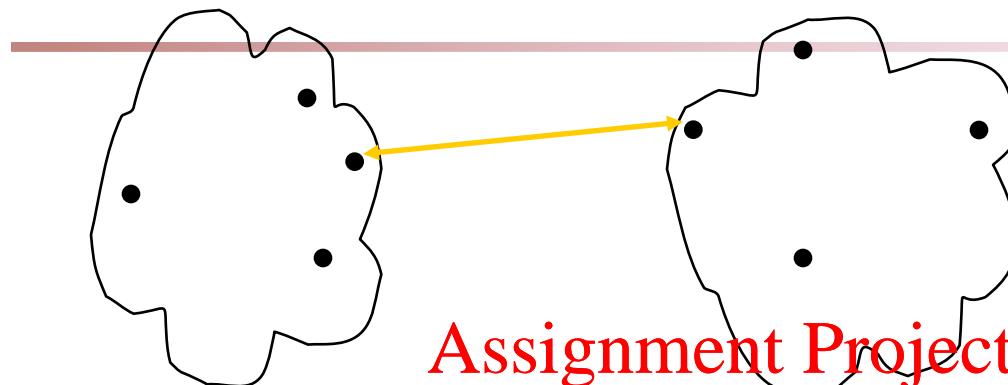
Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

· Proximity Matrix

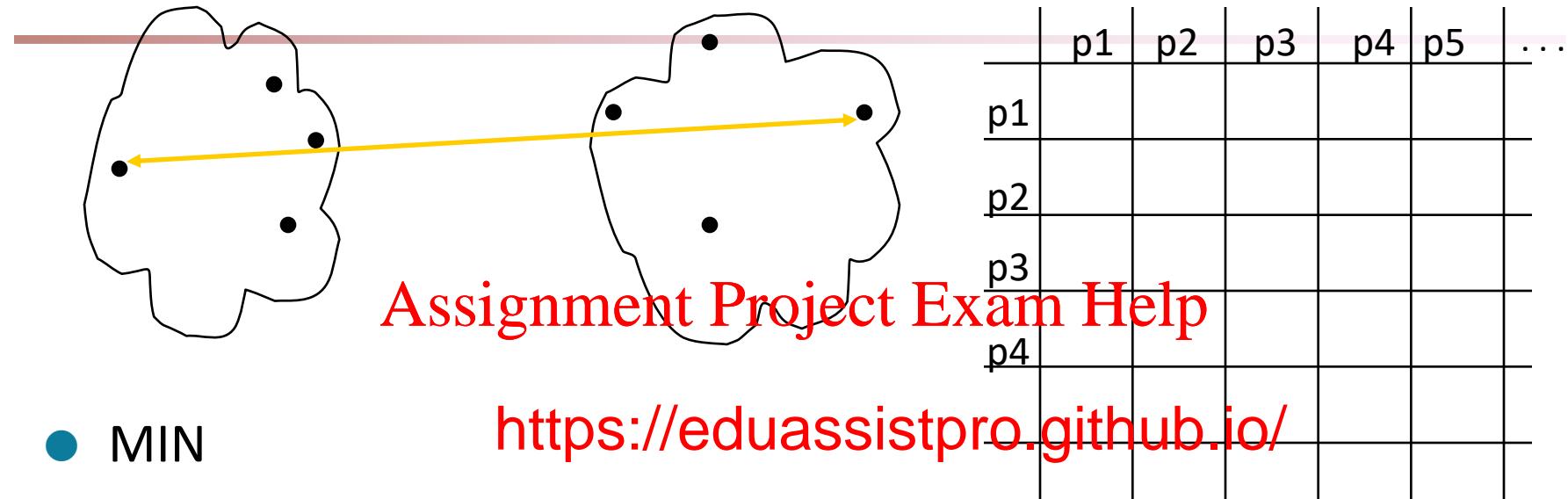
# How to Define Inter-Cluster Similarity



# Assignment Project Exam Help

- MIN
  - MAX
  - Centroid-based
  - Group Average
  - Other methods driven by an objective function
    - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



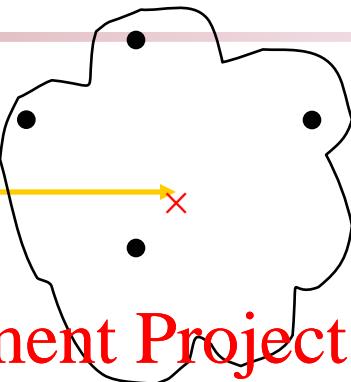
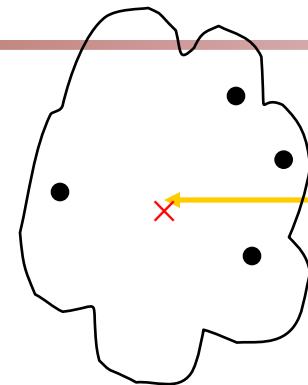
- MIN
- MAX
- Centroid-based
- Group Average
- Other methods driven by an objective function
  - Ward's Method uses squared error

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

Add WeChat edu\_assist\_pro

· Proximity Matrix

# How to Define Inter-Cluster Similarity



	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
...						

Assignment Project Exam Help

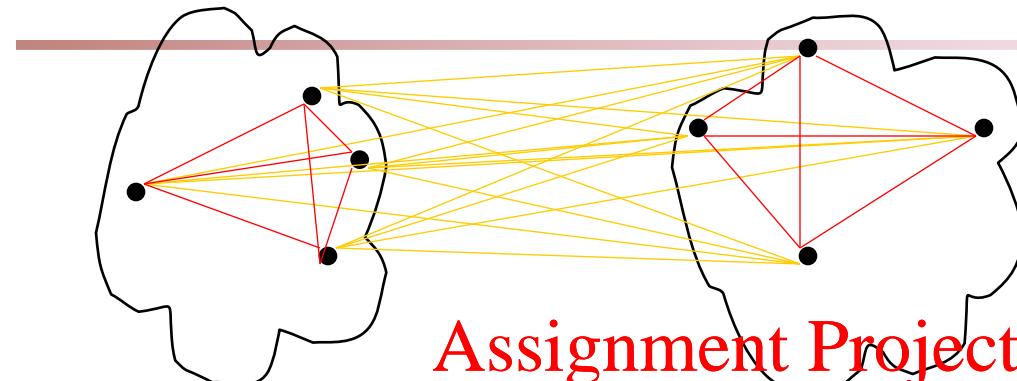
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

- MIN
- MAX
- Centroid-based
- Group Average
- Other methods driven by an objective function
  - Ward's Method uses squared error

· Proximity Matrix

# How to Define Inter-Cluster Similarity



Assignment Project Exam Help

- MIN
- MAX
- Centroid-based
- **Group Average**
- Other methods driven by an objective function
  - Ward's Method uses squared error

Add WeChat edu\_assist\_pro

Note: not simple avg distance between the clusters

Proximity Matrix

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
...						

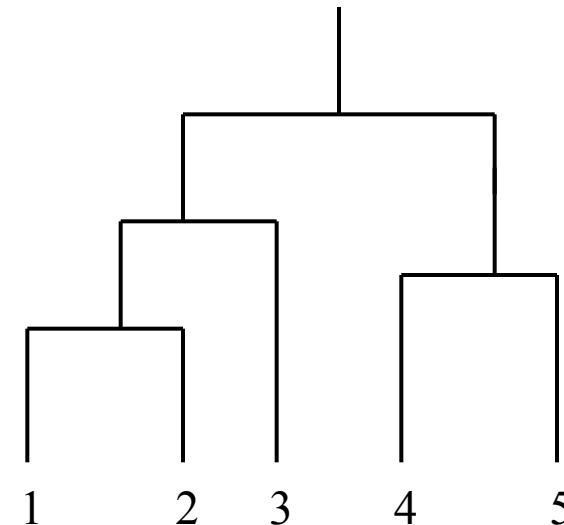
# Cluster Similarity: MIN or Single Link/LINK

- Similarity of two clusters is based on the two most similar (closest) points in the different clusters
  - i.e.,  $\text{sim}(C_i, C_j) = \min(\text{dissim}(p_x, p_y))$  //  $p_x \in C_i, p_y \in C_j$
  - Determined by the minimum distance between points in the proximity graph

[Assignment](#) [Project](#) [Exam](#) [Help](#)

Add WeChat **edu\_assist\_pro**

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00



$$\text{sim}(C_i, C_j) = \max(\text{sim}(p_x, p_y))$$

# Single-Link Example

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	

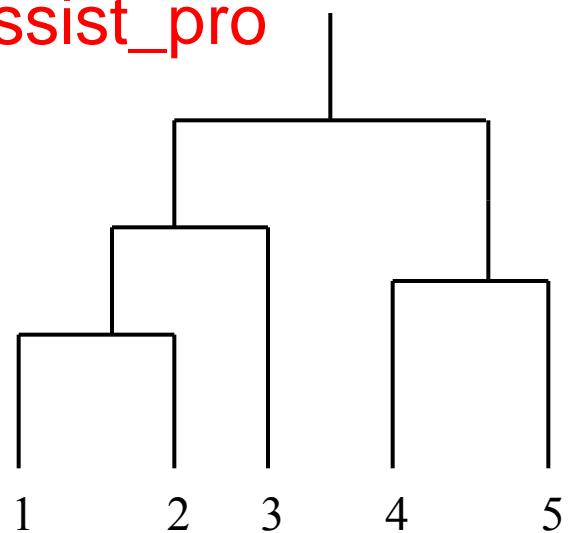
Assignment Project Exam Help  
<https://eduassistpro.github.io/>

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2		1.00	0.70	0.60	0.50
P3			1.00	0.40	0.30
P4				1.00	0.80
P5					1.00

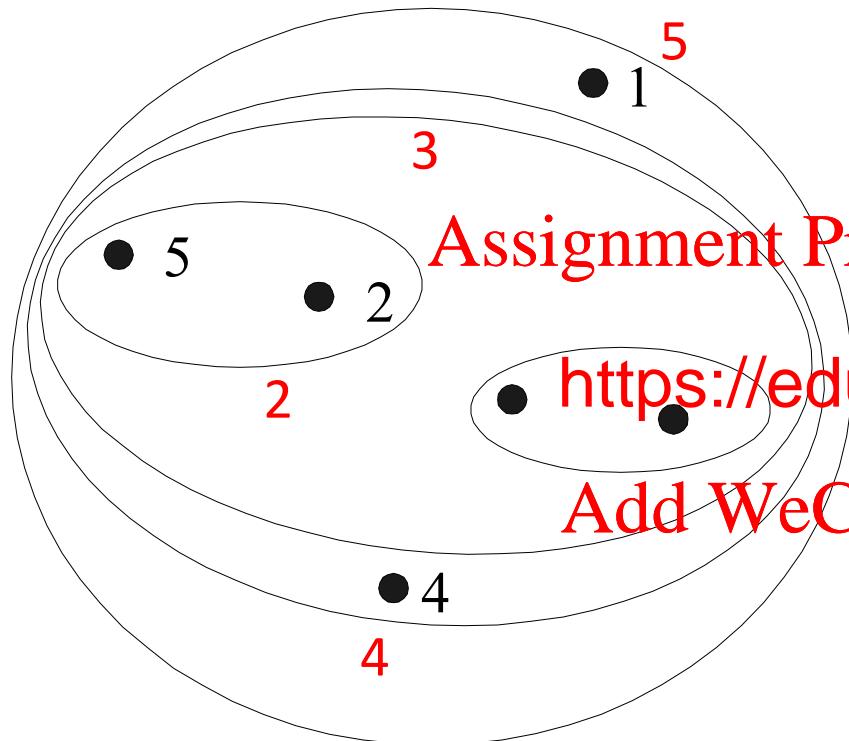
Add WeChat edu\_assist\_pro

	12	P3	P4	P5
12	1.00	0.70	0.65	0.50
P3		1.00	0.40	0.30
P4			1.00	0.80
P5				1.00

	12	P3	
12	1.00	0.70	0.65
P3		1.00	0.40
45			1.00



# Hierarchical Clustering: MIN



Nested Clusters



Dendrogram

# Strength of MIN

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

Original Points

Two Clusters

- Can handle non-elliptical shapes

# Limitations of MIN

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

Original Points

Two Clusters

- Sensitive to noise and outliers

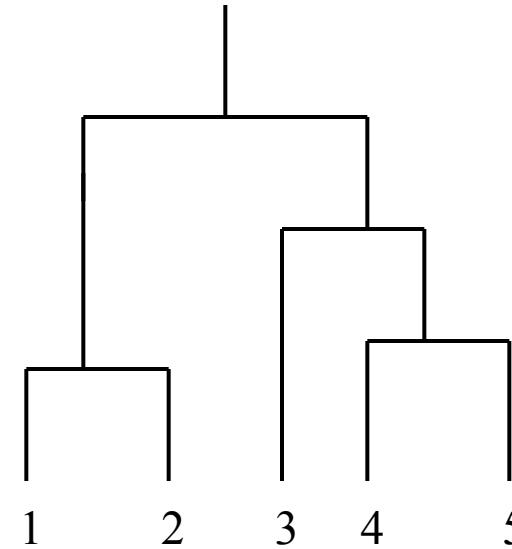
# Cluster Similarity: MAX or Complete Link (CLINK)

- Similarity of two clusters is based on the two least similar (most distant) points in the different clusters
  - i.e.,  $\text{sim}(C_i, C_j) = \max(\text{dissim}(p_x, p_y))$  //  $p_x \in C_i, p_y \in C_j$
  - Determined by  $\max(\text{dissim}(p_x, p_y))$  o clusters

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

Add WeChat edu\_assist\_pro

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00



$$\text{sim}(C_i, C_j) = \min(\text{sim}(p_x, p_y))$$

# Complete-Link Example

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	

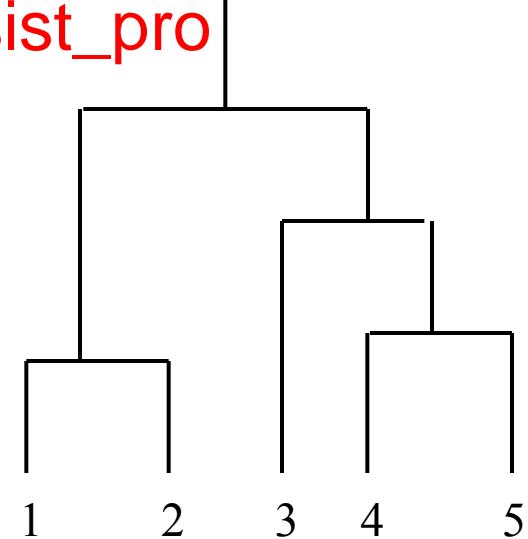
Assignment Project Exam Help  
<https://eduassistpro.github.io/>

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2		1.00	0.70	0.60	0.50
P3			1.00	0.40	0.30
P4				1.00	0.80
P5					1.00

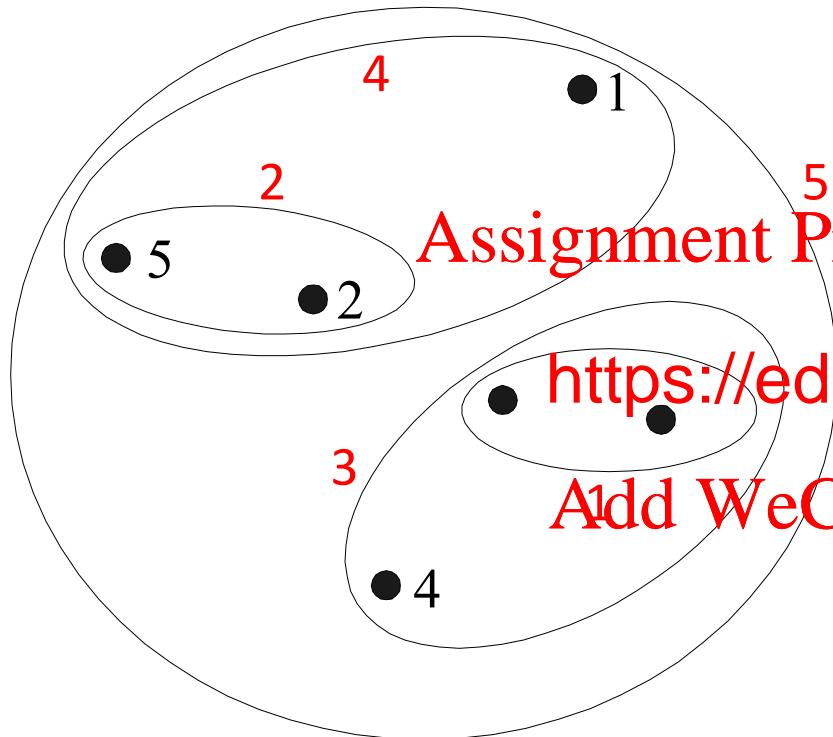
Add WeChat edu\_assist\_pro

	12	P3	P4	P5
12	1.00	0.10	0.60	0.20
P3		1.00	0.40	0.30
P4			1.00	0.80
P5				1.00

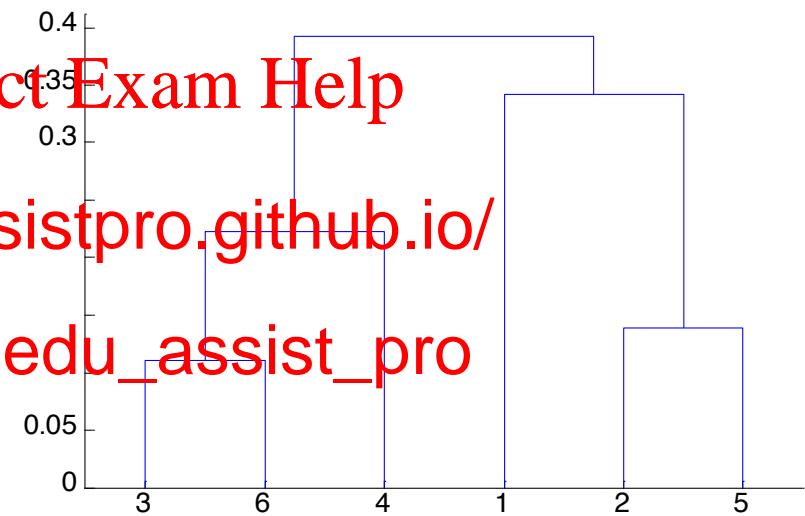
	12	P3	P5
12	1.00	0.10	0.20
P3		1.00	0.30
45			1.00



# Hierarchical Clustering: MAX



Nested Clusters



Dendrogram

# Strength of MAX

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

Original Points

Two Clusters

- Less susceptible to noise and outliers

# Limitations of MAX

---

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

Original Points

Two Clusters

- Tends to break large clusters
- Biased towards globular clusters

# Cluster Similarity: Group Average

- GAAC (Group Average Agglomerative Clustering)
- Similarity of two clusters is the average of pair-wise similarity between points in the two clusters.

$$\text{similarity(Cluster}_i, \text{Cluster}_j) = \frac{\sum_{\substack{p_i, p_j \in \text{Cluster}_i \cup \text{Cluster}_j \\ p_j \neq p_i}} \text{similarity}(p_i, p_j)}{*(|\text{Cluster}_i| + |\text{Cluster}_j| - 1)}$$

Assignment Project Exam Help  
<https://eduassistpro.github.io/>

- Why not using si his method guarantees that ~~Add Websat~~ [Add Websat](https://eduassistpro.github.io/) ~~edu\_assist\_pro~~

(3, 4.5)

(1, 1)

(5, 1)

$$\text{sim}(C_i, C_j) = \text{avg}(\text{sim}(p_i, p_j))$$

# Group Average Example

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	

Assignment Project Exam Help  
<https://eduassistpro.github.io/>

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2		1.00	0.70	0.60	0.50
P3			1.00	0.40	0.30
P4				1.00	0.80
P5					1.00

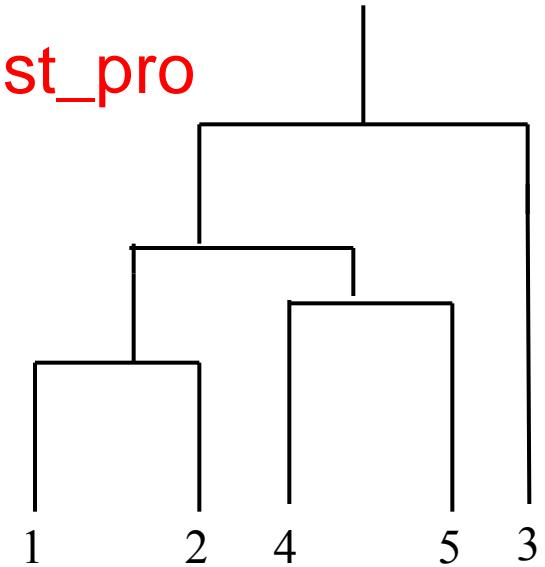
Add WeChat edu\_assist\_pro

	12	P3	P4	P5
12	1.00	0.567	0.717	0.533
P3		1.00	0.40	0.30
P4			1.00	0.80
P5				1.00

	12	P3	
12	1.0	0.567	0.608
P3		1.00	0.5
45			1.00

$$\text{Sim}(12,3)=2*(0.1+0.7+0.9)/6 = 0.5666666$$

$$\text{Sim}(12,45)=2*(0.9+0.65+0.2+0.6+0.5+0.8)/12 = 0.608$$



# Hierarchical Clustering: Centroid-based and Group Average

---

- Compromise between Single and Complete Link
  - Assignment Project Exam Help**
- Strengths <https://eduassistpro.github.io/>
  - Less susceptible to noise outliers
- Limitations
  - Biased towards globular clusters

# More on Hierarchical Clustering Methods

---

- Major weakness of agglomerative clustering methods
  - do not scale well: time complexity of at least  $O(n^2)$ , where  $n$  is the number of total objects
  - can never uniquely determine the clusters
- Integration of hierarchical clustering
  - BIRCH (1996): uses CF-tree to quickly adjusts the quality of sub-clusters
  - CURE (1998): selects well-scattered points from the cluster and then shrinks them towards the center of the cluster by a specified fraction
  - CHAMELEON (1999): hierarchical clustering using dynamic modeling

# Spectral Clustering

---

- See additional slides.

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro