

BANA 273 Session 8

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro  
Prof. Vibs

The Paul Merage School of Business  
University of California, Irvine

# Agenda

- Assignment 4 due on Canvas soon
- Please work on your projects
- Clustering using k-means algorithm

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

Add WeChat edu\_assist\_pro

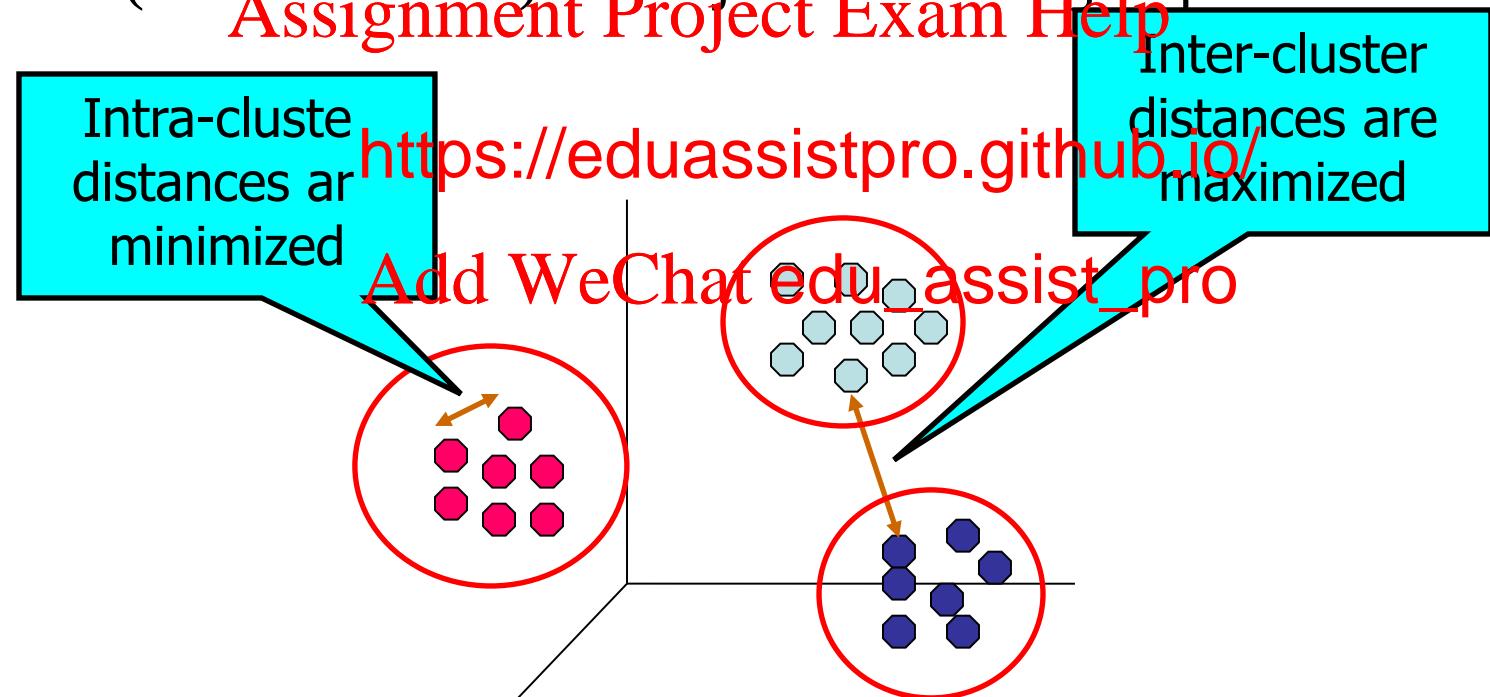
# Clustering Definition

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters
  - Data points similar to one another.
  - Data points in separate clusters similar to one another.
- Similarity Measures:
  - Euclidean Distance if attributes are continuous.
  - Other Problem-specific Measures.

# What is Cluster Analysis?

- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups

Assignment Project Exam Help



□ Euclidean Distance Based Clustering in 3-D space.

# Clustering: Application 1

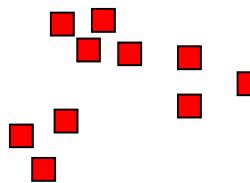
- Market Segmentation:
  - Goal: subdivide a market into distinct subsets of customers where any subset may conceivably be selected as a market target to be reached with
  - Approach: <https://eduassistpro.github.io/>
    - Collect different attributes of customers related to their geographical and lifestyle information
    - Find clusters of similar customers
    - Measure the clustering quality by observing buying patterns of customers in same cluster vs. those from different clusters.

# Notion of a Cluster can be Ambiguous

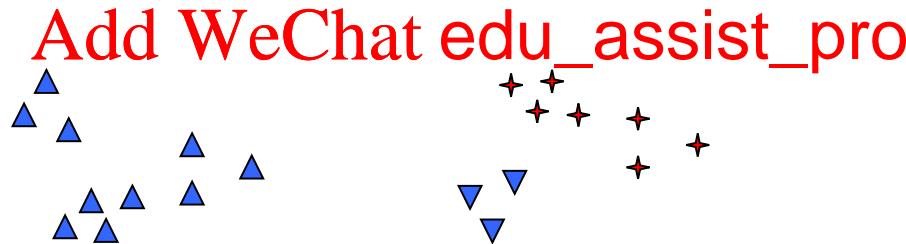


How many clusters?

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



Two Clusters

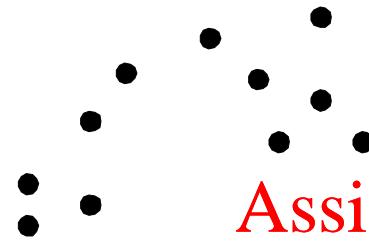


Four Clusters

# Types of Clusterings

- A **clustering** is a set of clusters
- Important distinction between **hierarchical** and **partitional** <https://eduassistpro.github.io/>
- Partitional Clustering
  - A division of data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset
- Hierarchical clustering
  - A set of nested clusters organized as a hierarchical tree

# Partitional Clustering

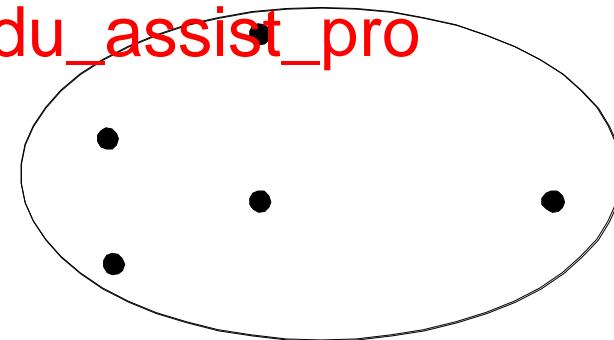


Original Points

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



A Partitional Clustering

# K-Means Clustering

1. Begin by specifying K, the number of clusters
2. Select ~~Assignment Project Exam Help~~ initial centroids
3. Assign each point to the closest centroid  
(Euclidean Distance)  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro
4. Re-compute the centroids of the clusters
5. Repeat steps 3 and 4 until points stop moving between clusters

# Similarity Measure

- Need a distance measure
  - Example of a distance measure:  
Assignment Project Exam Help
    - Manhattan
- <https://eduassistpro.github.io/>
- Add  $D(X, Y) = \sum_{i=1}^n |x_i - y_i|$

# Similarity Metric

- Example for a distance measure:
  - Euclidean distance

Assignment Project Exam Help

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

# Example of Euclidean Distance



John:  
Age=35  
Income  
no. of c



Rachel:  
Age=41  
Income=215K  
no. of credit cards=2

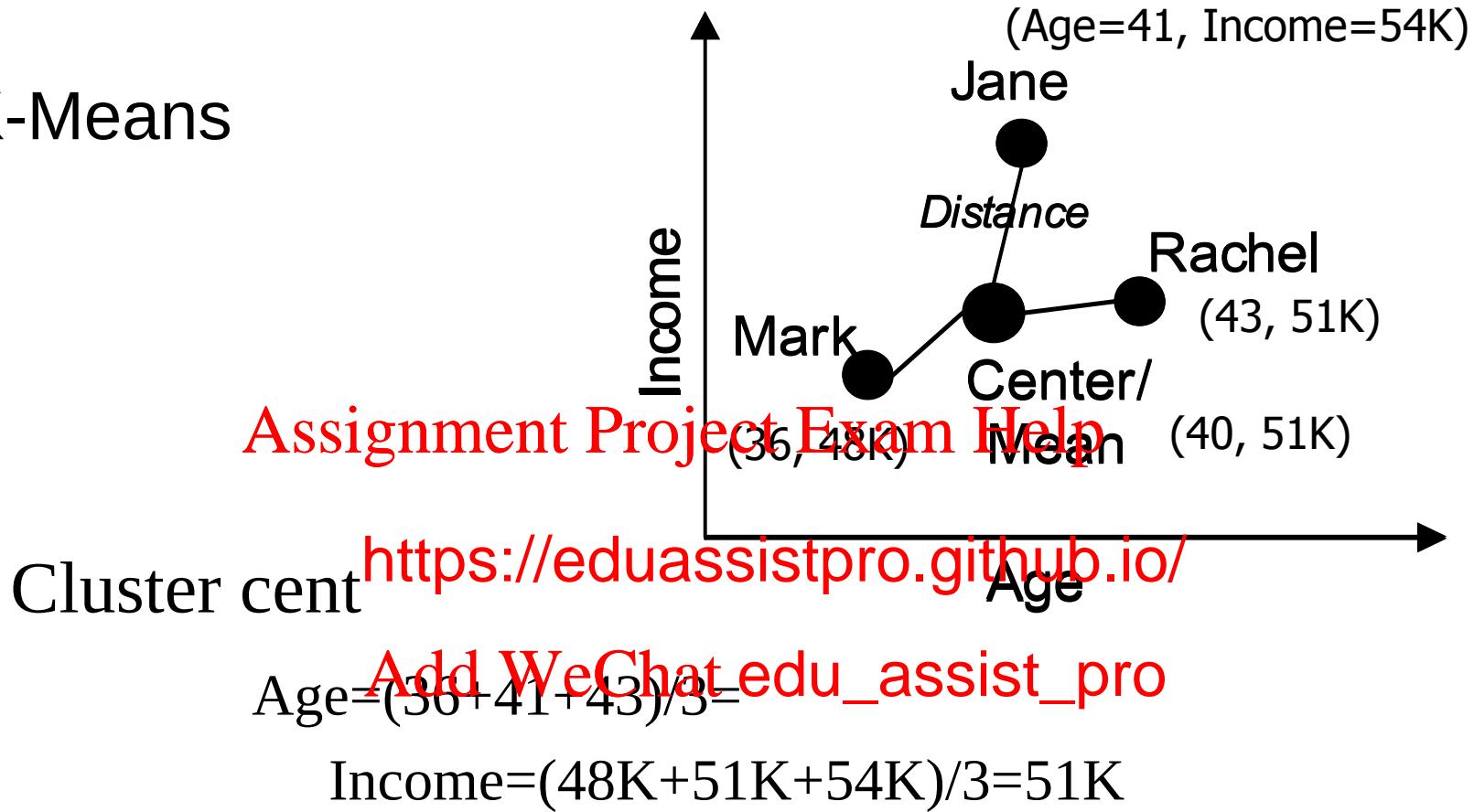
Assignment Project Exam Help  
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

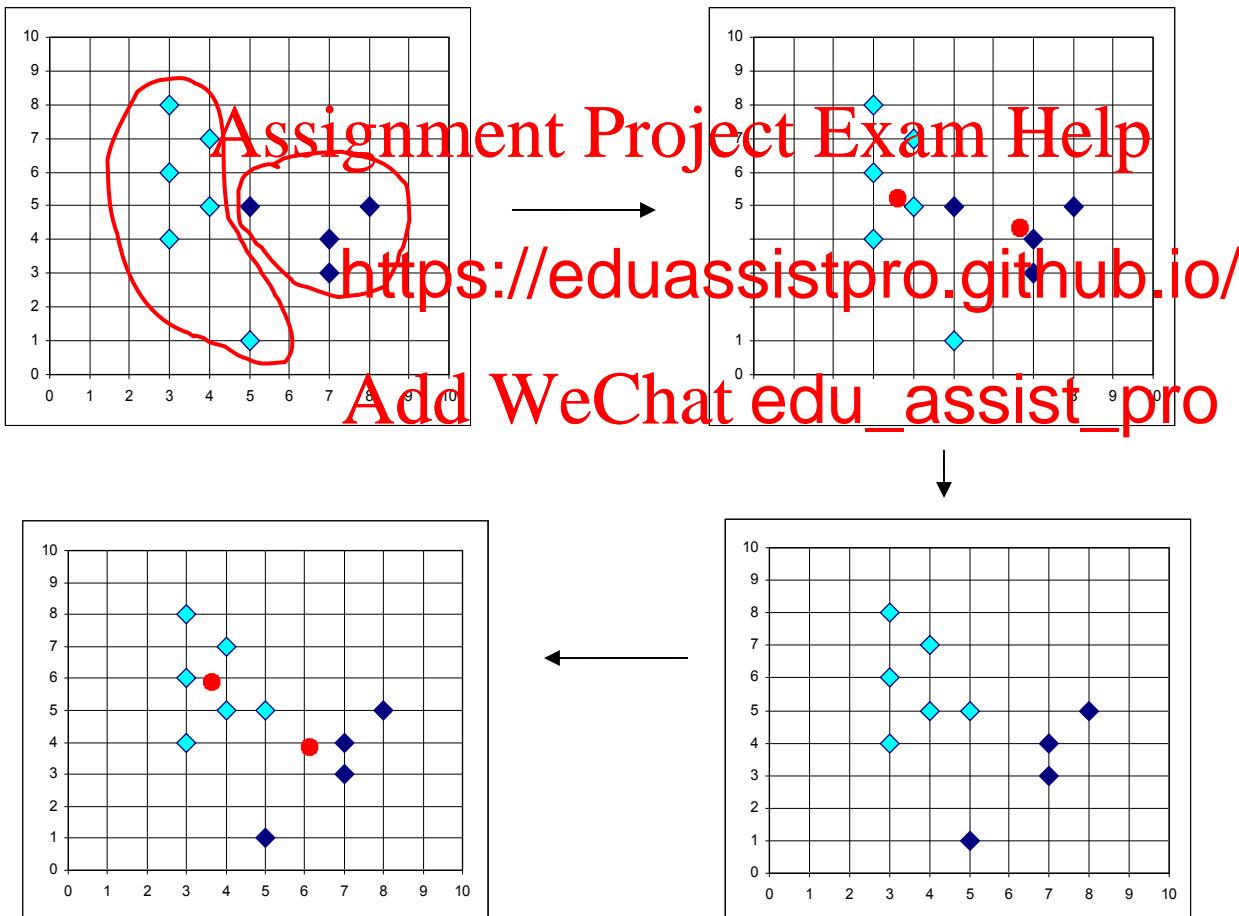
$$D(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Distance (John, Rachel)=sqrt [(35-41)<sup>2</sup>+(95-215)<sup>2</sup>+(3-2)<sup>2</sup>]

# K-Means



# Example: 2-Means



# K-means clustering

## 1. Select inputs



Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# K-means clustering

1. Select inputs
2. Select k cluster centers



# K-means clustering



1. Select inputs
2. Select k cluster centers
3. Assign cases to closest

# K-means clustering



1. Select inputs
2. Select k cluster centers
3. Assign cases to closest

center  
need to define “close”  
to cluster centers

# K-means clustering

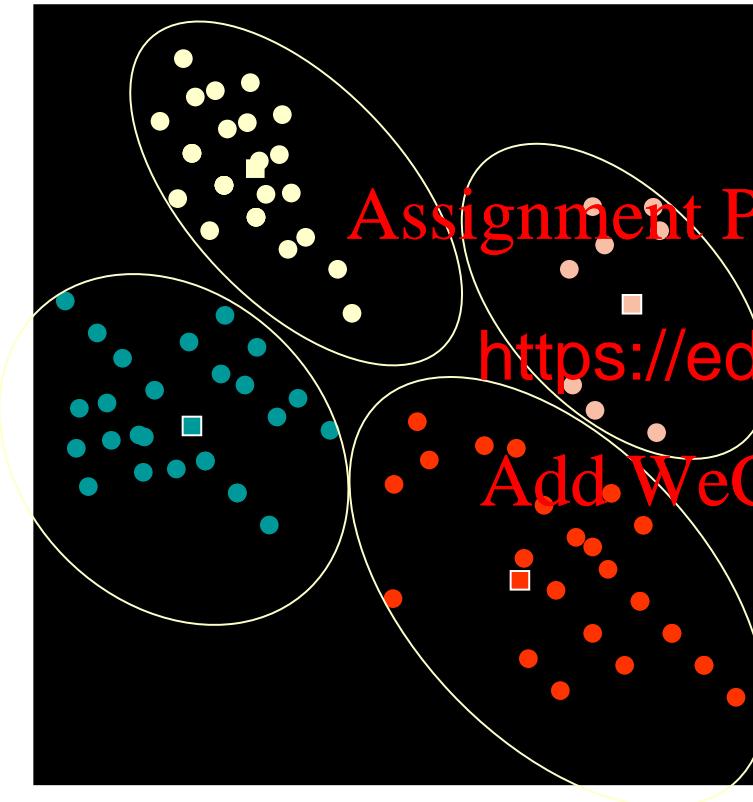


1. Select inputs
2. Select k cluster centers
3. Assign cases to closest

center  
eed to define “close”  
te cluster centers

gn cases

# K-means clustering



1. Select inputs
  2. Select k cluster centers
  3. Assign cases to closest center
- Need to define “close” to cluster centers
- steps 4 and 5 until changes in cluster centers & assigned cases are insignificant

# “k” in k-means clustering

- Generally,  $k$  is set in advance
- If not known, <https://eduassistpro.github.io/> Assignment Project Exam Help  
values of  $k$  to expect from the clusters. The sum of distances (in the clusters) for clusters one r of clusters one larger  $k$ 's

# Cluster Validity

- Compute ratio
  - = [sum of squared distances for a given  $k$ ] / [sum of squared distances for all records ( $k = 1$ )]
    - If the ratio is not very effective
    - If it is small we have well-separated clusters
- Weka reports sum of squared errors (Intra cluster distance)

# Example

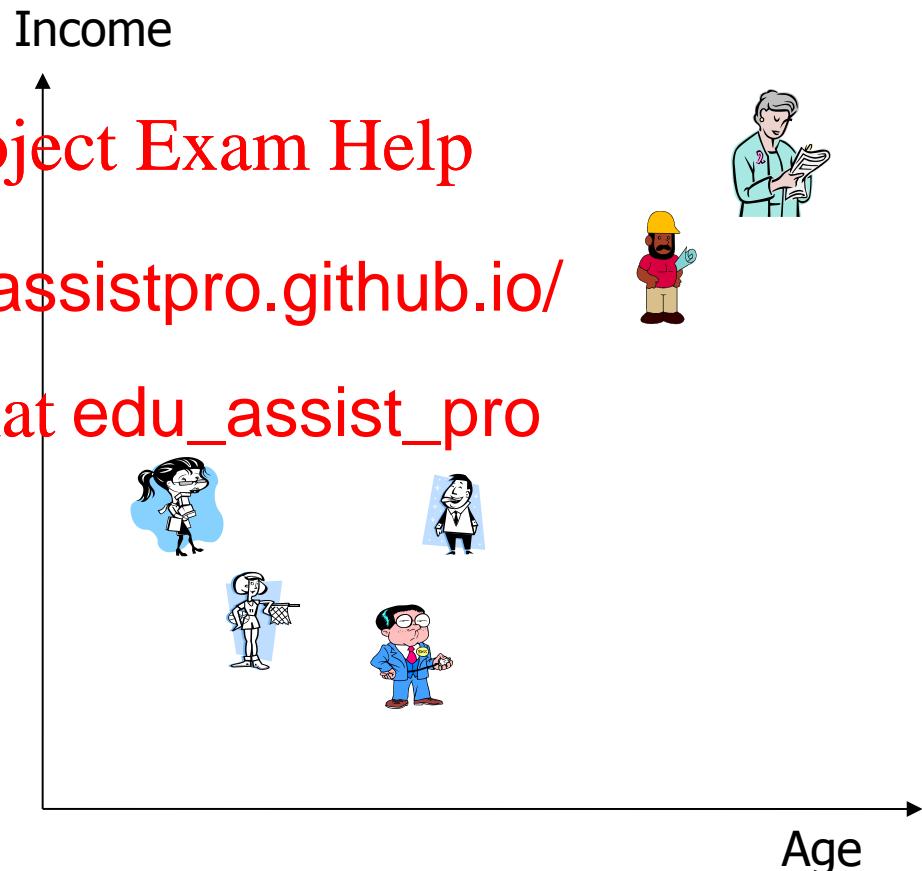
Note: Both Age and Income are normalized.

Customer	Age	Income (K)
John	0.55	0.175
Rachel		
Hannah	1	1
Tom	0.93	0.85
Nellie	0.39	0.2
David	0.58	0.25

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



# K-Means Algorithm: Example

## Step 1:

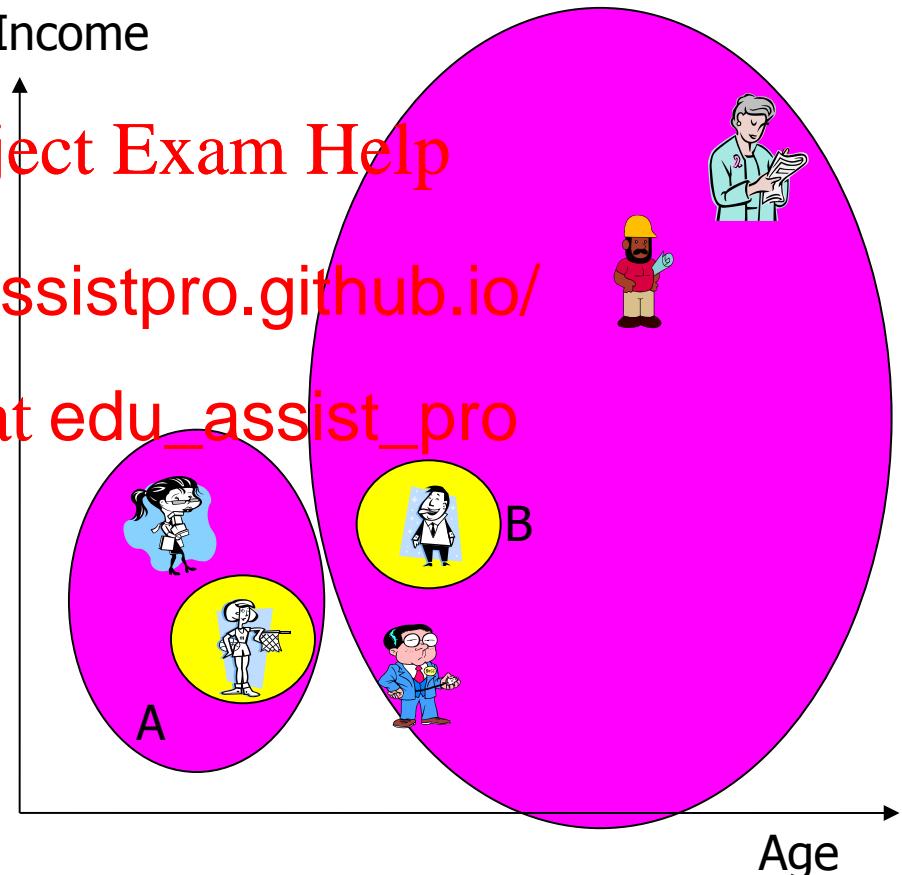
Nellie and David are selected as cluster centers A and B respectively

Customer	Distance from David	Distance from Nellie
John	0.08	
Rachel	0.24	
Hannah	0.86	1.01
Tom	0.69	0.85
Nellie		
David		

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



# K-Means Algorithm: Example

## Calculate cluster center:

Cluster A center:

- Age 0.37, Income=0.23

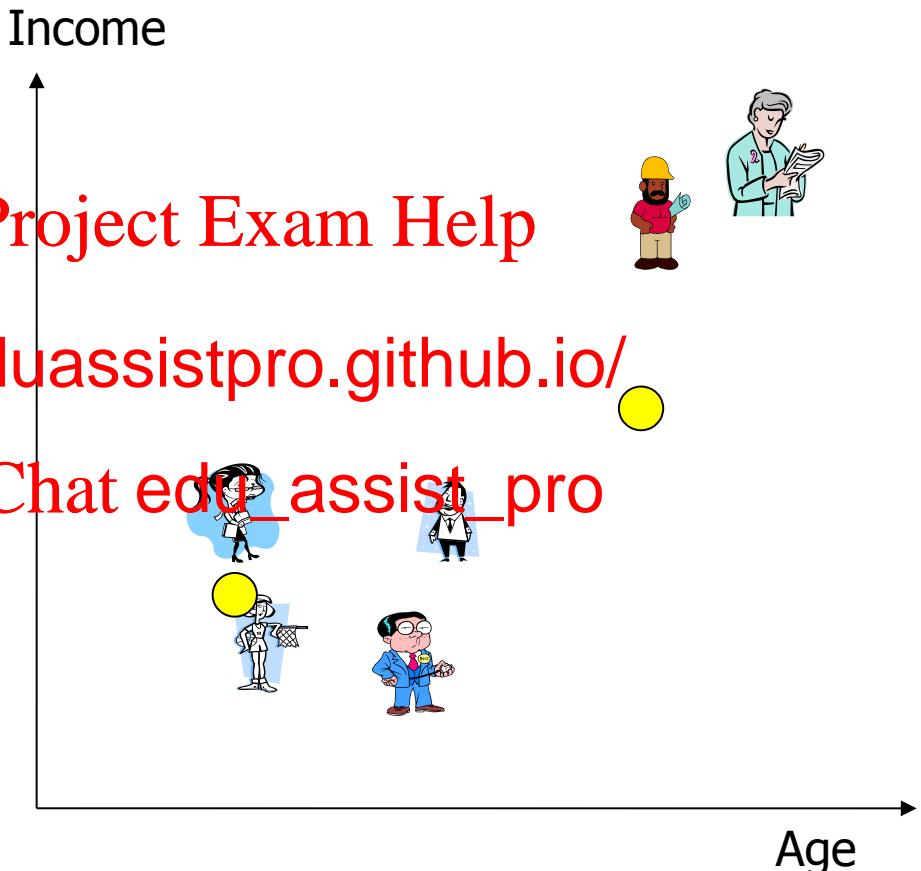
Cluster B center:

- Age 0.77, Inco

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

Assign customers to clusters

based on new cluster centers



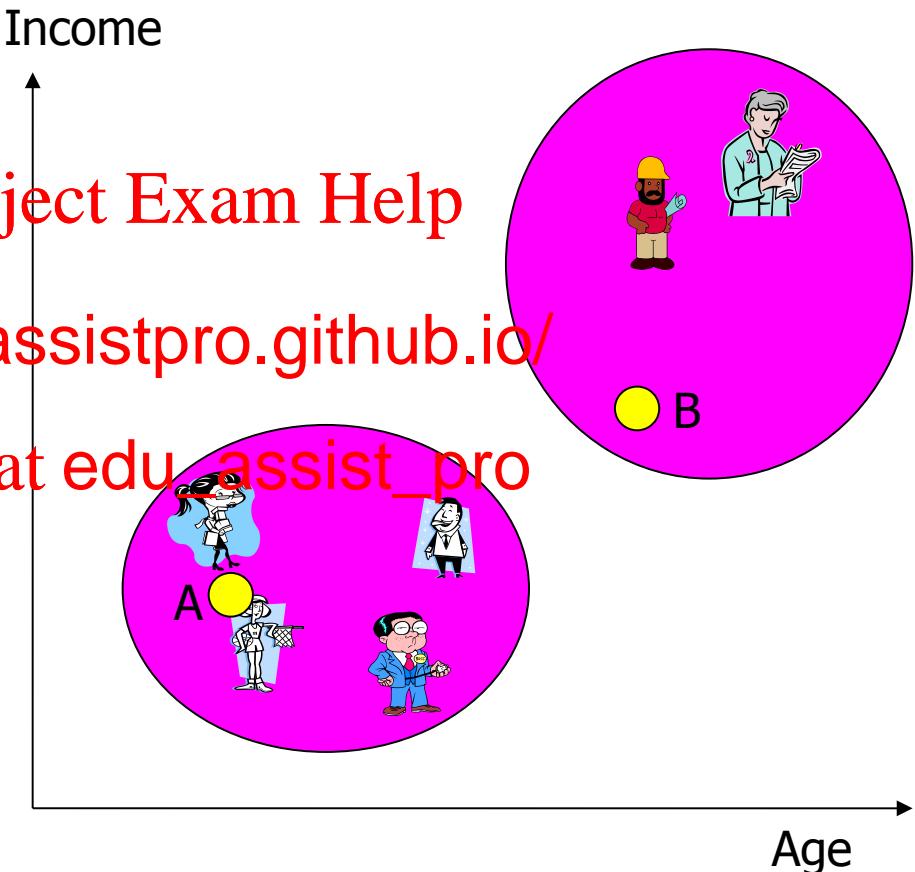
# K-Means Algorithm: Example

Customer	Distance A	Distance B
John	0.19	
Rachel	0.04	
Hannah	0.99	0.49
Tom	0.84	0.32
Nellie	0.04	0.53
David	0.21	0.37

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



## K-Means Algorithm: Example

### Calculate cluster center:

Cluster A center:

- Age 0.47, Income=0.22

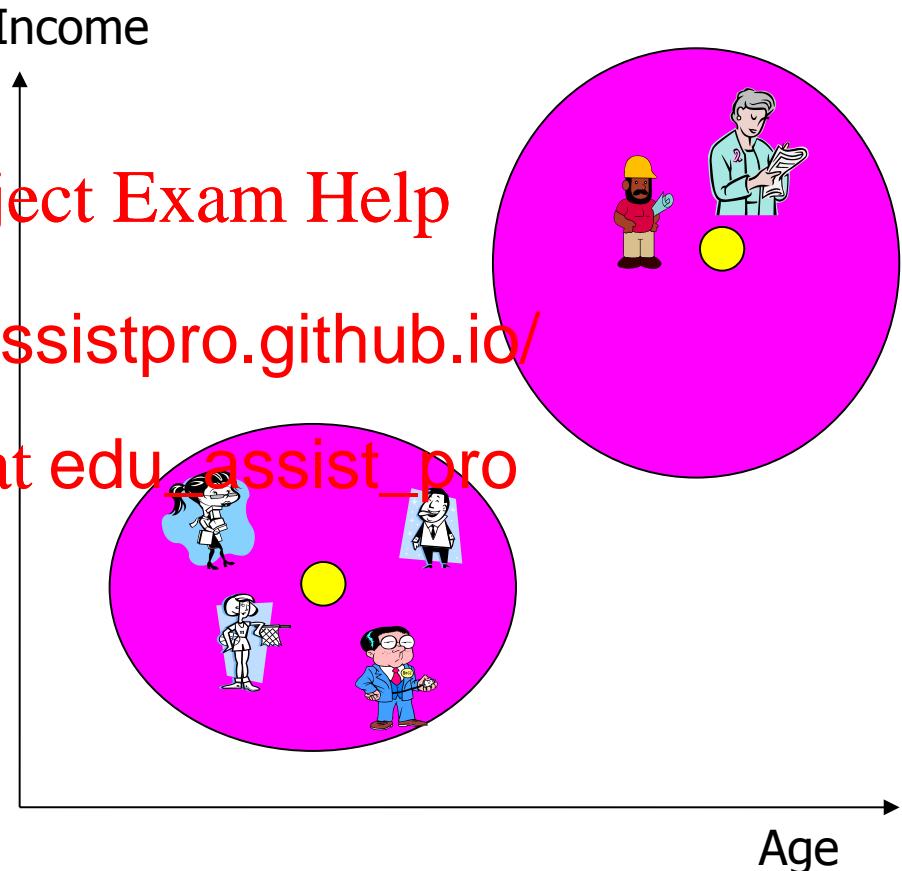
Cluster B center:

- Age 0.97, Inc
- Clusters do not change

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



# Scale and Weigh Data

- Scaling makes sure that the distance is not biased by units (1K, 1M, etc.)
- Weighting can add the information that one variable is more (or less) important.  
<https://eduassistpro.github.io/>
- After scaling to different units, use weights to introduce bias based on the business context.
  - (eg. Two households with the same income are more similar than two households with the same number of pets.)
- Common way to scale:
  - Range:  $(\text{value}-\min)/(\max-\min)$ ; [0,1]
    - E.g.  $\{11,8,4,6,10,1\} \rightarrow \{1, 0.7, 0.3, 0.5, 0.9, 0\}$

# What is a “Good” cluster?

- A. Inter-cluster distance is maximized and intra-cluster distance is minimized
- B. Inter-cluster distance is minimized and intra-cluster distance is maximized
- C. Inter-cluster distance is minimized and intra-cluster distance is minimized
- D. Inter-cluster distance is minimized and intra-cluster distance is minimized
- E. None of the Above

# Clustering in Weka

## Utility Example

Assignment Project Exam Help  
East West Airli

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

<http://facweb.cs.depaul.edu/110/lectures/lect584/WEKA/k-means.html>

# Clustering Exercise

Start with individuals 1 and 4 as initial centroids

Assignment Project Exam Help  
Student A B C D

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

3	3.0	
4	5.0	
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

# Strengths and Weaknesses of the K-Means

- Strength
  - Relatively efficient
  - Simple implementation

Assignment Project Exam Help

- Weakness
  - Need to specify the number of clusters, in advance
  - Unable to handle noisy data well
  - Euclidian Distance does not work well for nominal variables.

# Applications of Clustering

- **Marketing:** Customer segmentation (discovery of distinct groups of customers) for target marketing. Create product differentiation: different offers for different segments (It's not always positive.)  
*Assignment Project Exam Help*
- **Car insurance:** <https://eduassistpro.github.io/> ps with high average claim cost
- **Property:** Identify houses in the same area with similar characteristics  
*Add WeChat edu\_assist\_pro*
- **Image recognition**
- Creating **document collections**, or grouping web pages

# Review of Assignments

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Next Session

- Review of Assignment 4
- Review ~~Assignment Project Exam Help~~
- Other Data <https://eduassistpro.github.io/>
  - Text Minin
  - Collaborative Filtering [Add WeChat edu\\_assist\\_pro](#)