

BIG DATA

Analytics & Management

Lecture 7 (03/27, 03/29): Clustering

Decisions, Operations & Information Technologies
Robert H. Smith School of Business
Spring, 2017



Outline

- Introduction
- Similarity Measures
- K-means Algorithm
- MapReduce-based K-means
- Agglomerative Algorithm
- Labelling and Quality of Clusters

Supervised learning vs. unsupervised learning



- **Supervised learning:** discover patterns in the data that relate data attributes with a target (class) attribute.
 - These patterns are then utilized to predict the values of the target attribute in future data instances.
- **Unsupervised learning:** The data have no target attribute.
 - We want to explore the data to find some intrinsic structures in them.

Clustering

- Clustering is a technique for finding **similarity groups** in data, called **clusters**. I.e.,
 - it groups data instances that are similar to (near) each other in one cluster and data instances that are very different (far away) from each other into different clusters.
- Clustering is often called an **unsupervised learning** task as no class values denoting an *a priori* grouping of the data instances are given, which is the case in supervised learning.
- Due to historical reasons, clustering is often considered synonymous with unsupervised learning.
 - In fact, association rule mining is also unsupervised
- This week focuses on clustering.

An illustration

- The data set has three natural groups of data points, i.e., 3 natural clusters.



What is clustering for?

- Let us see some real-life examples
- Example 1: groups people of similar sizes together to make “small”, “medium” and “large” T-Shirts.
 - Tailor-made for each person: too expensive
 - One-size-fits-all: does not fit all.
- Example 2: In marketing, segment customers according to their similarities
 - To do targeted marketing.

What is clustering for?

- Example 3: Given a collection of text documents, we want to organize them according to their content similarities,
 - To produce a topic hierarchy
- In fact, clustering is one of the most utilized data mining techniques.
 - It has a long history, and used in almost every field, e.g., medicine, psychology, botany, sociology, biology, archeology, marketing, insurance, libraries, etc.
 - In recent years, due to the rapid increase of online documents, text clustering becomes important.

Aspects of clustering

- A clustering algorithm
 - Partitional clustering
 - Hierarchical clustering
 - ...
- A distance (similarity, or dissimilarity) function
- Clustering quality
 - Inter-clusters distance \Rightarrow maximized
 - Intra-clusters distance \Rightarrow minimized
- The **quality** of a clustering result depends on the algorithm, the distance function, and the application.

How do we define “similarity”?

- Recall that the goal is to group together “similar” data – but what does this mean?
- No single answer – it depends on what we want to find or emphasize in the data; this is one reason why clustering is an “art”
- The similarity measure is often more important than the clustering algorithm used – don’t overlook this choice!

(Dis)similarity measures

- Instead of talking about similarity measures, we often equivalently refer to dissimilarity measures (I'll give an example of how to convert between them in a few slides...)
- Jagota defines a dissimilarity measure as a function $f(x,y)$ such that $f(x,y) > f(w,z)$ if and only if x is less similar to y than w is to z
- This is always a *pair-wise* measure

Distance functions

- There are numerous distance functions for
 - Different types of data
 - Numeric data
 - Nominal data
 - Different specific applications

Distance functions for numeric attributes



- Most commonly used functions are
 - Euclidean distance and
 - Manhattan distance
- We denote distance with: $dist(\mathbf{x}_i, \mathbf{x}_j)$, where \mathbf{x}_i and \mathbf{x}_j are data points (vectors)
- They are special cases of Minkowski distance. h is positive integer.

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \left((x_{i1} - x_{j1})^h + (x_{i2} - x_{j2})^h + \dots + (x_{ir} - x_{jr})^h \right)^{\frac{1}{h}}$$

Euclidean distance and Manhattan distance

- If $h = 2$, it is the Euclidean distance

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2}$$

- If $h = 1$, it is the Manhattan distance

$$dist(\mathbf{x}_i, \mathbf{x}_j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ir} - x_{jr}|$$

- Weighted Euclidean distance

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \dots + w_r(x_{ir} - x_{jr})^2}$$

Squared distance and Chebychev distance



- **Squared Euclidean distance:** to place progressively greater weight on data points that are further apart.

$$dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2$$

- **Chebychev distance:** one wants to define two data points as "different" if they are different on any one of the attributes.

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|, \dots, |x_{ir} - x_{jr}|)$$

Distance functions for binary and nominal attributes



- **Binary attribute:** has two values or states but no ordering relationships, e.g.,
 - Gender: male and female.
- We use a confusion matrix to introduce the distance functions/measures.
- Let the i^{th} and j^{th} data points be \mathbf{x}_i and \mathbf{x}_j (vectors)

Confusion matrix

		Data point j		
		1	0	
Data point i	1	a	b	a+b
	0	c	d	c+d
		a+c		b+d
				a+b+c+d

- a: The number of attributes with the value of 1 for both data points
- b: The number of attributes with which $x_{if} = 1$ and $x_{jf} = 0$, where x_{if} is the value of f^{th} attribute of the data point x_i
- c: The number of attributes for which $x_{if}=0$ and $x_{jf}=1$
- d: The number of attributes with the value of 0 for both data points

Symmetric binary attributes

- A binary attribute is **symmetric** if both of its states (0 and 1) have equal importance, and carry the same weights, e.g., male and female of the attribute “Gender”
- Distance function: **Simple Matching Coefficient**, proportion of mismatches of their values

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \frac{b + c}{a + b + c + d}$$

Symmetric binary attributes: example

X ₁	1	1	1	0	1	0	0
X ₂	0	1	1	0	0	1	0

$$\text{dist}(X_1, X_2) = (2+1)/(2+2+1+2) = 0.429$$

Asymmetric binary attributes

- **Asymmetric:** if one of the states is more important or more valuable than the other.
 - By convention, state 1 represents the more important state.
 - Jaccard coefficient is a popular measure

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \frac{b + c}{a + b + c}$$

- We can have some variations, adding weights

Nominal attributes

- Nominal attributes: with more than two states or values.
 - ❑ the commonly used distance measure is also based on the simple matching method.
 - ❑ Given two data points \mathbf{x}_i and \mathbf{x}_j , let the number of attributes be r , and the number of values that match in \mathbf{x}_i and \mathbf{x}_j be q .

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \frac{r - q}{r}$$

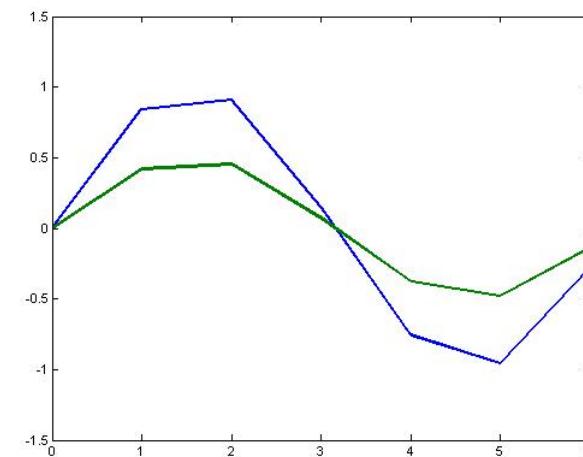
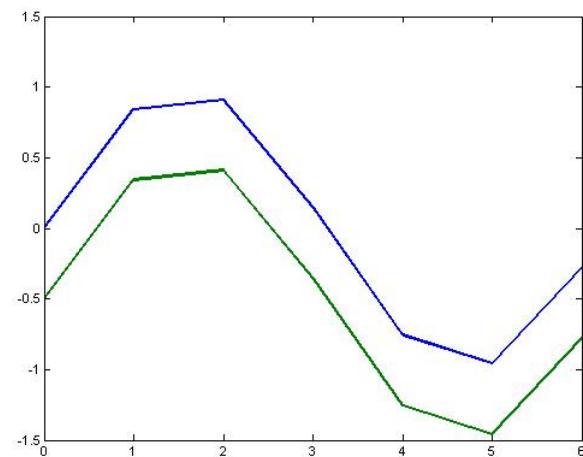
Distance function for text documents



- A text document consists of a sequence of sentences and each sentence consists of a sequence of words.
- To simplify: a document is usually considered a “bag” of words in document clustering.
 - Sequence and position of words are ignored.
- A document is represented with a vector just like a normal data point.
- It is common to use similarity to compare two documents rather than distance.
 - The most commonly used similarity function is the **cosine similarity**.

Correlation

- We might care more about the overall shape of expression profiles rather than the actual magnitudes
- That is, we might want to consider similarity when they are “up” and “down” together
- When might we want this kind of measure? What experimental issues might make this appropriate?



Pearson Linear Correlation

$$\rho(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\bar{x} = \frac{1}{n} \sum_i^n x_i$$

$$\bar{y} = \frac{1}{n} \sum_i^n y_i$$

- We're shifting the expression profiles down (subtracting the means) and scaling by the standard deviations (i.e., making the data have mean = 0 and std = 1)

Pearson Linear Correlation

- Pearson linear correlation (PLC) is a measure that is invariant to scaling and shifting (vertically) of the expression values
- Always between -1 and $+1$ (perfectly anti-correlated and perfectly correlated)
- This is a similarity measure, but we can easily make it into a dissimilarity measure:

$$d_p = \frac{1 - \rho(\mathbf{x}, \mathbf{y})}{2}$$

Data standardization

- In the Euclidean space, standardization of attributes is recommended so that all attributes can have equal impact on the computation of distances.
- Consider the following pair of data points
 - \mathbf{x}_i : (0.1, 20) and \mathbf{x}_j : (0.9, 720).

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(0.9 - 0.1)^2 + (720 - 20)^2} = 700.000457,$$

- The distance is almost completely dominated by $(720-20) = 700$.
- **Standardize attributes:** to force the attributes to have a common value range

Interval-scaled attributes

- Their values are real numbers following a linear scale.
 - The difference in Age between 10 and 20 is the same as that between 40 and 50.
 - The key idea is that intervals keep the same importance through out the scale
- Two main approaches to standardize interval scaled attributes, **range** and **z-score**. f is an attribute

$$\text{range}(x_{if}) = \frac{x_{if} - \min(f)}{\max(f) - \min(f)},$$

Interval-scaled attributes

- **Z-score:** transforms the attribute values so that they have a mean of zero and a **mean absolute deviation** of 1. The mean absolute deviation of attribute f , denoted by s_f , is computed as follows

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|),$$

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf}),$$

Z-score:
$$z(x_{if}) = \frac{x_{if} - m_f}{s_f}.$$

Ratio-scaled attributes

- Numeric attributes, but unlike interval-scaled attributes, their scales are **exponential**,
- For example, the total amount of microorganisms that evolve in a time t is approximately given by Ae^{Bt} ,
 - where A and B are some positive constants.
- Do log transform: $\log(x_{if})$
 - Then treat it as an interval-scaled attribute

Nominal attributes

- Sometime, we need to transform nominal attributes to numeric attributes.
- Transform nominal attributes to binary attributes.
 - The number of values of a nominal attribute is v .
 - Create v binary attributes to represent them.
 - If a data instance for the nominal attribute takes a particular value, the value of its binary attribute is set to 1, otherwise it is set to 0.
- The resulting binary attributes can be used as numeric attributes, with two values, 0 and 1.

Nominal attributes: an example

- Nominal attribute *fruit*: has three values,
 - Apple, Orange, and Pear
- We create three binary attributes called, Apple, Orange, and Pear in the new data.
- If a particular data instance in the original data has Apple as the value for *fruit*,
 - then in the transformed data, we set the value of the attribute Apple to 1, and
 - the values of attributes Orange and Pear to 0

Ordinal attributes

- Ordinal attribute: an ordinal attribute is like a nominal attribute, but its values have a numerical ordering. E.g.,
 - ❑ Age attribute with values: Young, MiddleAge and Old. They are ordered.
 - ❑ Common approach to standardization: treat it as an interval-scaled attribute.

Mixed attributes

- Our distance functions given are for data with all numeric attributes, or all nominal attributes, etc.
- Practical data has different types:
 - Any subset of the 6 types of attributes,
 - interval-scaled,
 - symmetric binary,
 - asymmetric binary,
 - ratio-scaled,
 - ordinal and
 - nominal

Convert to a single type

- One common way of dealing with mixed attributes is to
 - ❑ Decide the dominant attribute type, and
 - ❑ Convert the other types to this type.
- E.g, if most attributes in a data set are interval-scaled,
 - ❑ we convert ordinal attributes and ratio-scaled attributes to interval-scaled attributes.
 - ❑ It is also appropriate to treat symmetric binary attributes as interval-scaled attributes.

Convert to a single type

- It does not make much sense to convert a **nominal attribute** or an **asymmetric binary attribute** to an interval-scaled attribute,
- Alternatively, a nominal attribute can be converted to a set of (symmetric) binary attributes, which are then treated as numeric attributes.



UNIVERSITY OF
MARYLAND

ROBERT H. SMITH

SCHOOL OF BUSINESS

K-Means Clustering

K-means clustering

- K-means is a **partitional clustering** algorithm
- Let the set of data points (or instances) D be $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, where $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ir})$ is a **vector** in a real-valued space $X \subseteq R^r$, and r is the number of attributes (dimensions) in the data.
- The k -means algorithm partitions the given data into k clusters.
 - Each cluster has a cluster **center**, called **centroid**.
 - k is specified by the user

K-means algorithm

- Given k , the k -means algorithm works as follows:
 - 1) Randomly choose k data points (**seeds**) to be the initial **centroids** (cluster centers)
 - 2) Assign each data point to the closest **centroid**
 - 3) Re-compute the **centroids** using the current cluster memberships.
 - 4) If a convergence criterion is not met, go to 2).

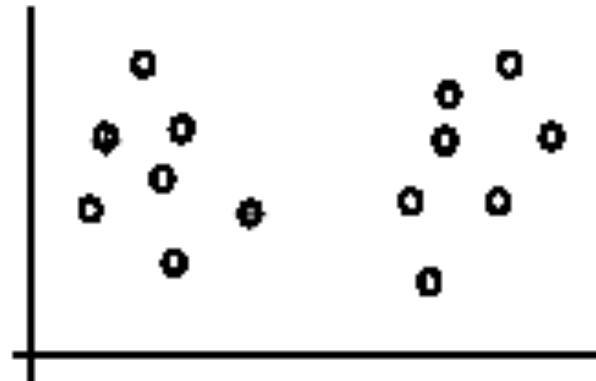
Stopping/convergence criterion

1. no (or minimum) re-assignments of data points to different clusters,
2. no (or minimum) change of centroids, or
3. minimum decrease in the **sum of squared error (SSE)**,

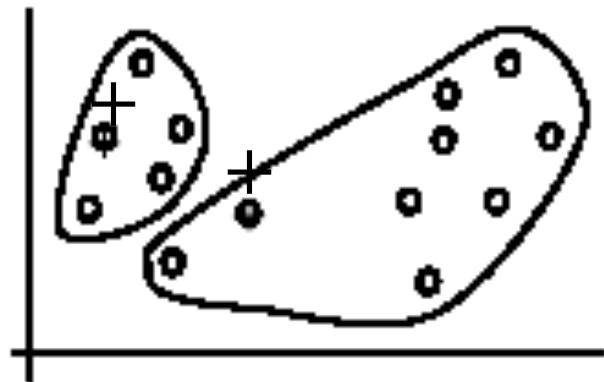
$$SSE = \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} dist(\mathbf{x}, \mathbf{m}_j)^2$$

- C_j is the j^{th} cluster, \mathbf{m}_j is the centroid of cluster C_j (the mean vector of all the data points in C_j), and $dist(\mathbf{x}, \mathbf{m}_j)$ is the distance between data point \mathbf{x} and centroid \mathbf{m}_j .

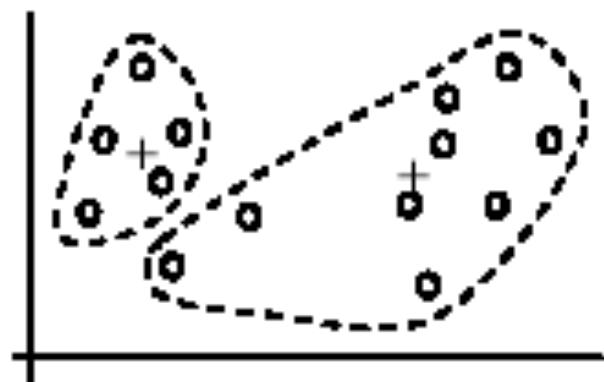
An example



(A). Random selection of k centers

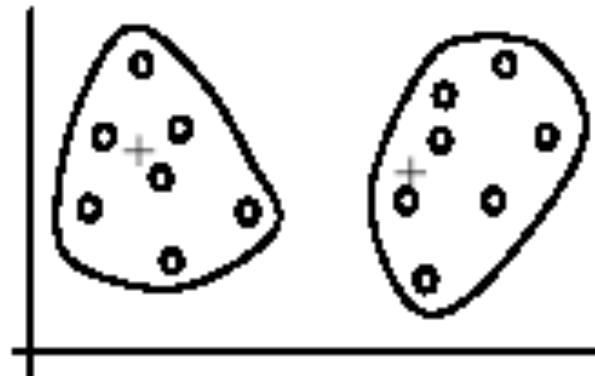


Iteration 1: (B). Cluster assignment

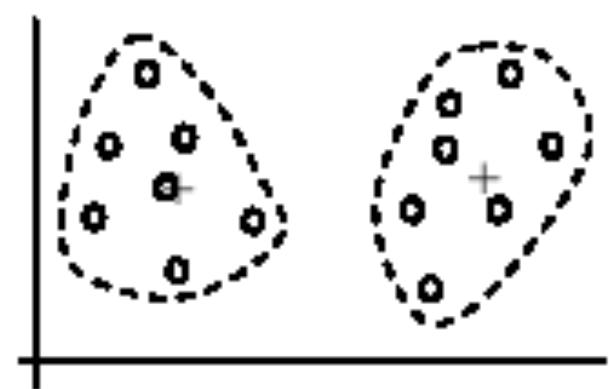


(C). Re-compute centroids

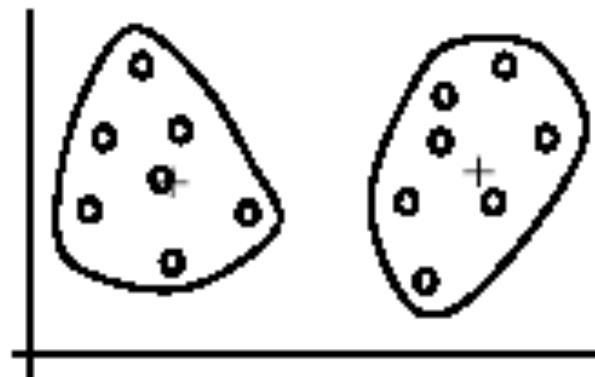
An example



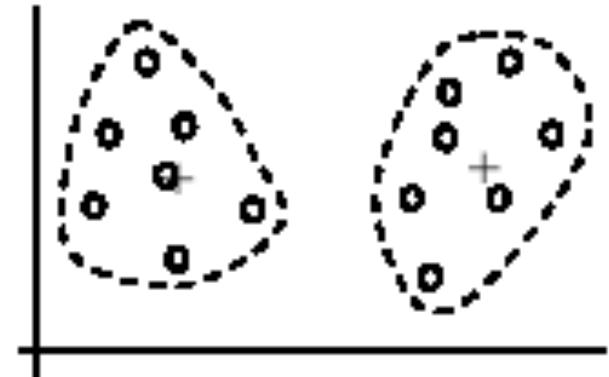
Iteration 2: (D). Cluster assignment



(E). Re-compute centroids



Iteration 3: (F). Cluster assignment

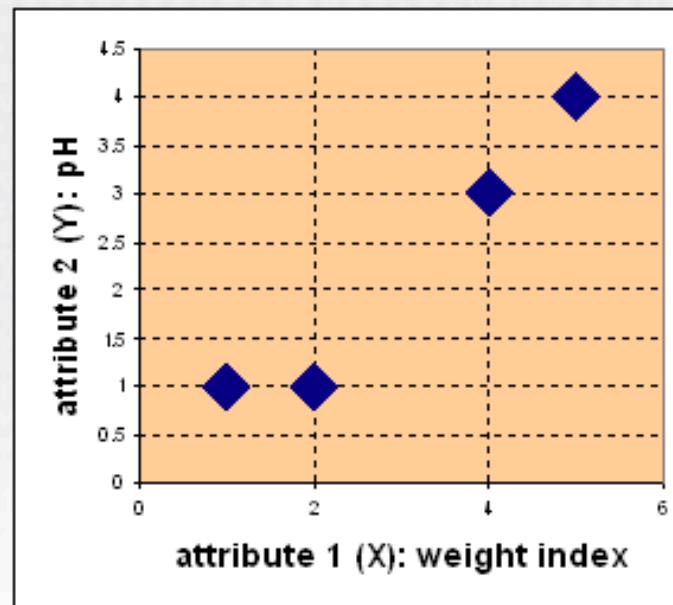


(G). Re-compute centroids

A numeric example

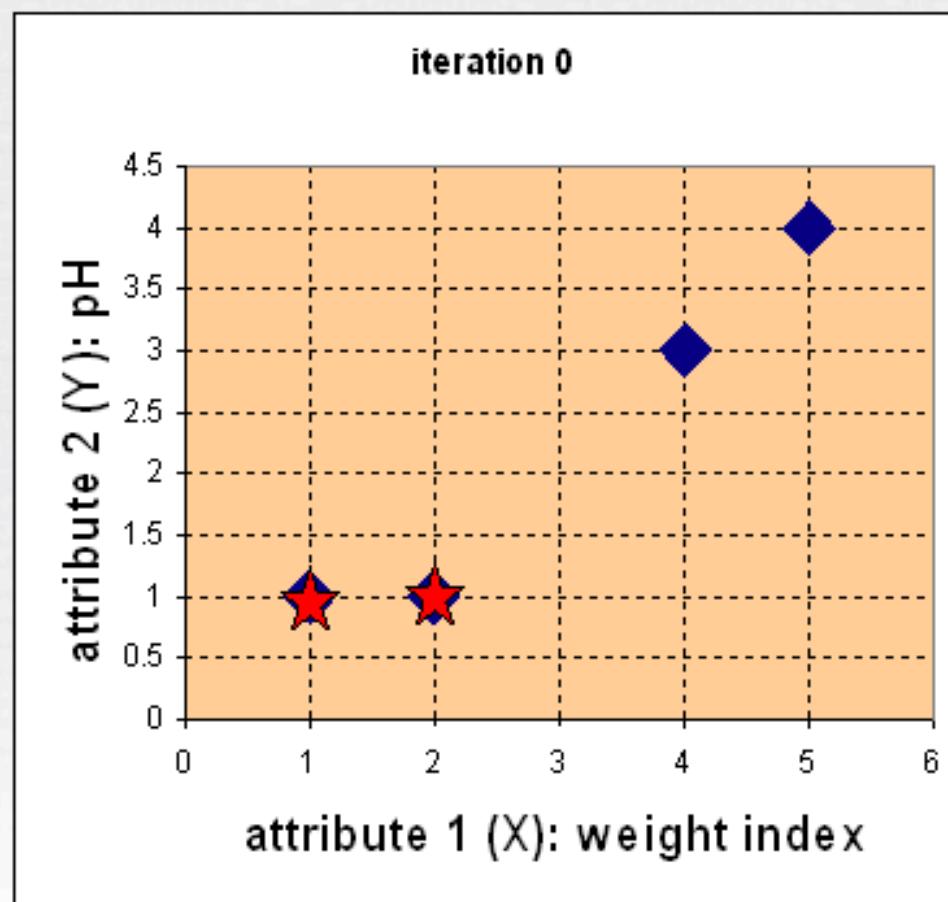
- 4 data points: A, B, C, D

Data point	Attribute 1 (weight)	Attribute 2 (ph index)
A	1	1
B	2	1
C	4	3
D	5	4



Initial centroid selection

- K=2
- Initial centroid (A, B): $c_1=(1,1)$, $c_2=(2,1)$



Iteration 0

- Calculate distance between the centroid and each data point
 - Distance of data point c to the first centroid and the second centroid

$$\mathbf{c}_1 = (1, 1) \quad \sqrt{(4-1)^2 + (3-1)^2} = 3.61$$

$$\mathbf{c}_2 = (2, 1) \quad \sqrt{(4-2)^2 + (3-1)^2} = 2.83$$

$$\mathbf{D}^0 = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 1 & 0 & 2.83 & 4.24 \end{bmatrix} \quad \mathbf{c}_1 = (1, 1) \quad group - 1 \\ \mathbf{c}_2 = (2, 1) \quad group - 2$$

- Cluster assignment

$$\mathbf{G}^0 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix} \quad group - 1 \\ \quad group - 2$$

A B C D

<i>A B C D</i>	$\begin{bmatrix} 1 & 2 & 4 & 5 \end{bmatrix}$	<i>X</i>
	$\begin{bmatrix} 1 & 1 & 3 & 4 \end{bmatrix}$	<i>Y</i>

Iteration 0

- Cluster assignment

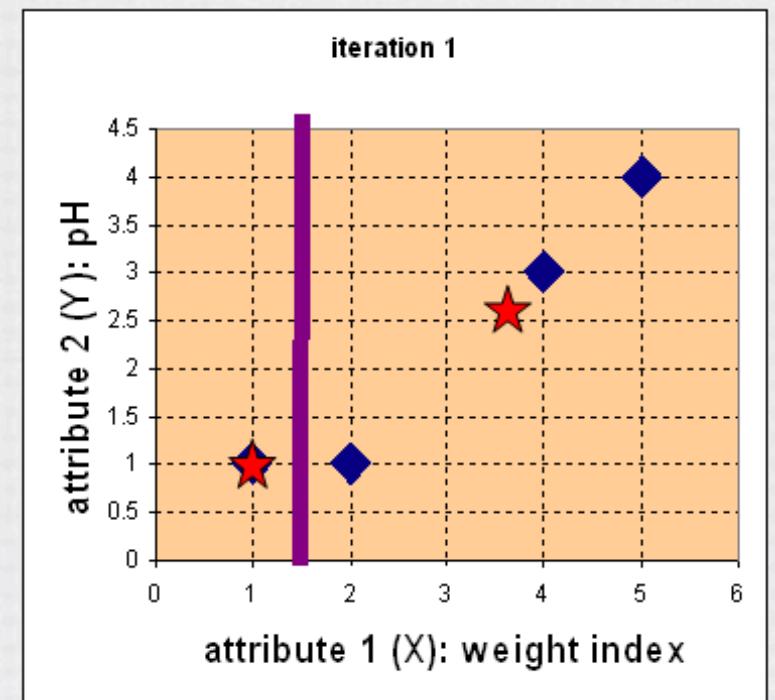
$$\mathbf{G}^0 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix} \quad \begin{array}{l} \text{group -1} \\ \text{group -2} \end{array}$$

A B C D

- New centroid

$$\mathbf{c}_1 = (1, 1)$$

$$\mathbf{c}_2 = \left(\frac{2+4+5}{3}, \frac{1+3+4}{3} \right) = \left(\frac{11}{3}, \frac{8}{3} \right)$$



Iteration 1

- Calculate distance between the centroid and each data point

$$\mathbf{D}^1 = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 3.14 & 2.36 & 0.47 & 1.89 \end{bmatrix} \quad \mathbf{c}_1 = (1, 1) \text{ group - 1} \\ \mathbf{c}_2 = \left(\frac{11}{3}, \frac{8}{3}\right) \text{ group - 2}$$

A	B	C	D	
[1]	2	4	5	X
[1]	1	3	4	Y

- Cluster assignment

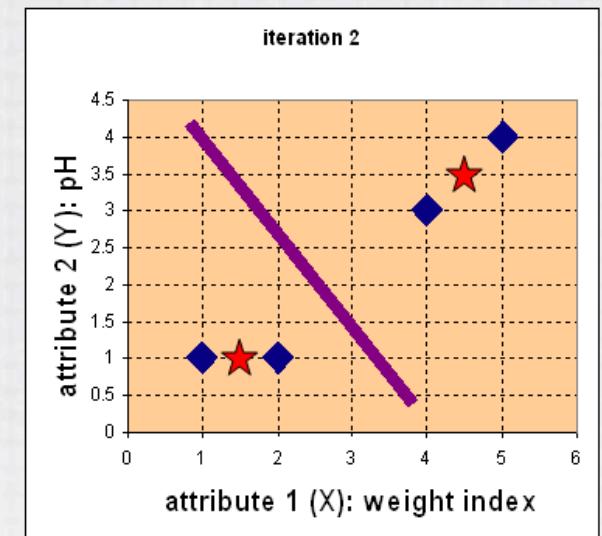
$$\mathbf{G}^1 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \text{group - 1} \\ \text{group - 2}$$

A	B	C	D
---	---	---	---

- New centroid

$$\mathbf{c}_1 = \left(\frac{1+2}{2}, \frac{1+1}{2}\right) = \left(1\frac{1}{2}, 1\right)$$

$$\mathbf{c}_2 = \left(\frac{4+5}{2}, \frac{3+4}{2}\right) = \left(4\frac{1}{2}, 3\frac{1}{2}\right)$$



Iteration 2

- Calculate distance between the centroid and each data point

$$\mathbf{D}^2 = \begin{bmatrix} 0.5 & 0.5 & 3.20 & 4.61 \\ 4.30 & 3.54 & 0.71 & 0.71 \end{bmatrix} \quad \mathbf{c}_1 = (1\frac{1}{2}, 1) \text{ group - 1}$$
$$\quad \quad \quad \mathbf{c}_2 = (4\frac{1}{2}, 3\frac{1}{2}) \text{ group - 2}$$
$$\begin{array}{cccc} A & B & C & D \end{array}$$
$$\begin{bmatrix} 1 & 2 & 4 & 5 \\ 1 & 1 & 3 & 4 \end{bmatrix} \quad X$$
$$Y$$

- Cluster assignment

$$\mathbf{G}^2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \text{group - 1}$$
$$\quad \quad \quad \text{group - 2}$$
$$\begin{array}{cccc} A & B & C & D \end{array}$$

$$\mathbf{G}^2 = \mathbf{G}^1$$



- New centroid G²

$$\mathbf{c}_1 = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = (1\frac{1}{2}, 1)$$

$$\mathbf{c}_2 = \left(\frac{4+5}{2}, \frac{3+4}{2} \right) = (4\frac{1}{2}, 3\frac{1}{2})$$

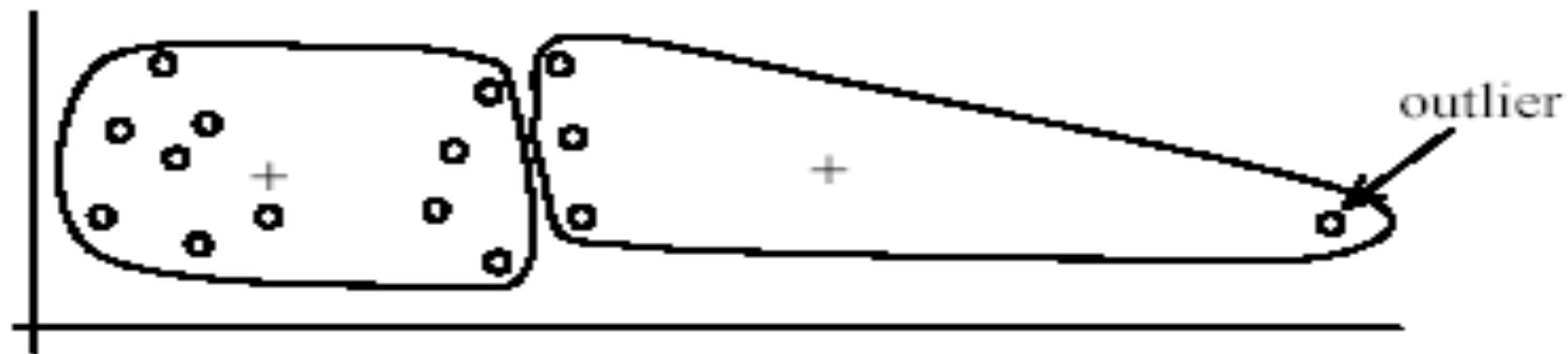
Strengths of k-means

- Strengths:
 - ❑ Simple: easy to understand and to implement
 - ❑ Efficient: time complexity: $O(tkn)$,
where n is the number of data points,
 k is the number of clusters, and
 t is the number of iterations.
 - ❑ Since both k and t are small. k-means is considered a linear algorithm.
- K-means is the most popular clustering algorithm.
- Note that: it terminates at a **local optimum** if SSE is used. The **global optimum** is hard to find due to complexity.

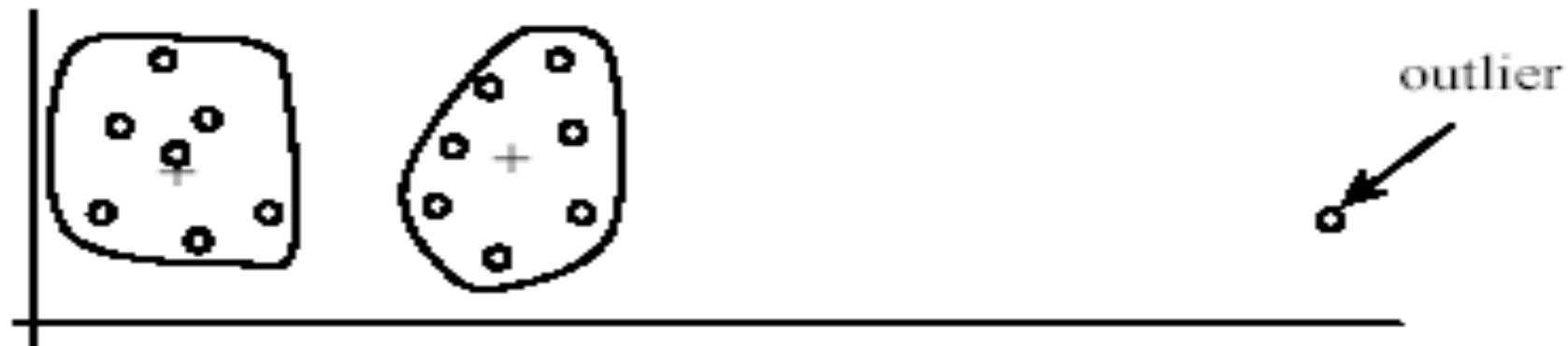
Weaknesses of k-means

- The algorithm is only applicable if the **mean** is defined.
 - For categorical data, k -mode - the centroid is represented by most frequent values.
- The user needs to specify **k** .
- The algorithm is sensitive to **outliers**
 - Outliers are data points that are very far away from other data points.
 - Outliers could be errors in the data recording or some special data points with very different values.

Weaknesses of k-means: Problems with outliers



(A): Undesirable clusters



(B): Ideal clusters

Weaknesses of k-means: To deal with outliers



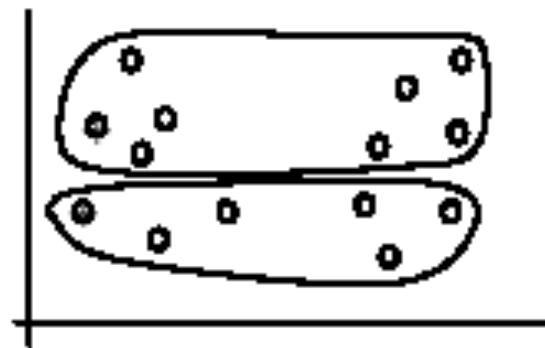
- One method is to remove some data points in the clustering process that are much further away from the centroids than other data points.
- Another method is to perform random sampling. Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
 - Assign the rest of the data points to the clusters by distance or similarity comparison, or classification

Weaknesses of k-means

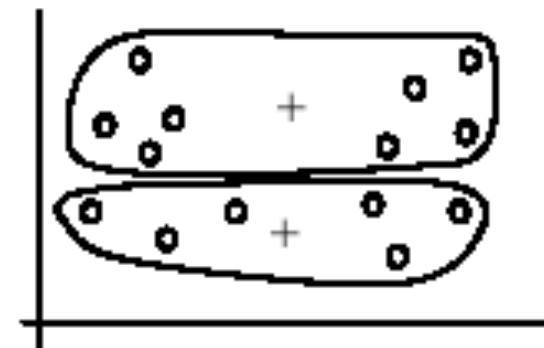
- The algorithm is sensitive to **initial seeds**.



(A). Random selection of seeds (centroids)



(B). Iteration 1

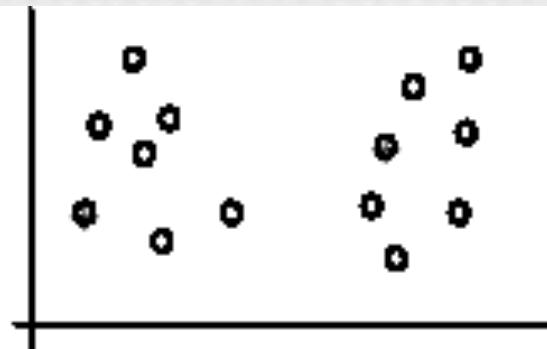


(C). Iteration 2

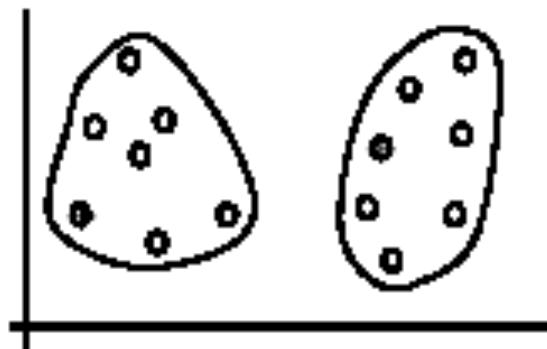
Weaknesses of k-means

- If we use **different seeds**: good results

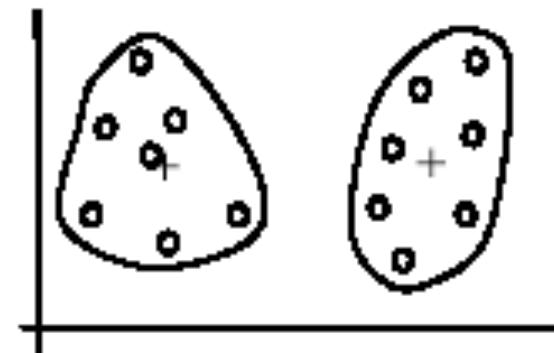
There are some methods to help choose good seeds



(A). Random selection of k seeds (centroids)



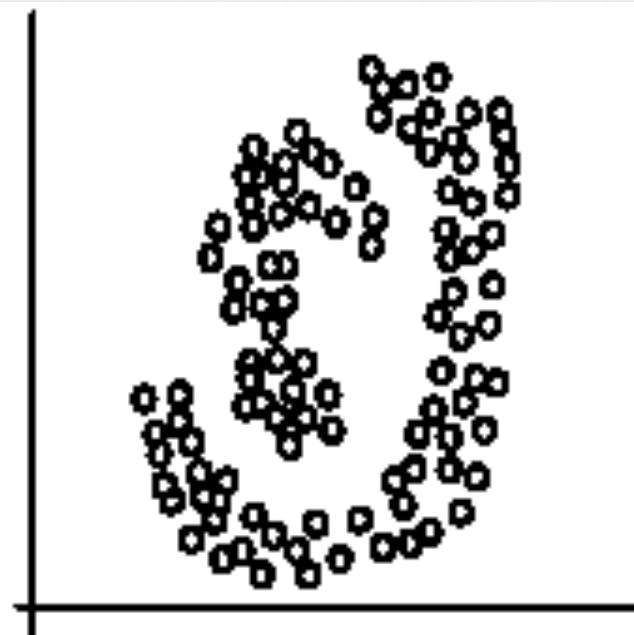
(B). Iteration 1



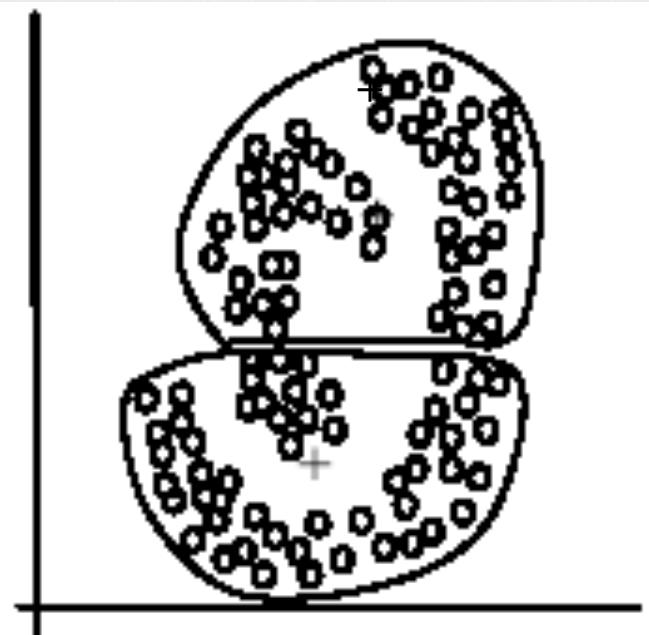
(C). Iteration 2

Weaknesses of k-means

- The k -means algorithm is not suitable for discovering clusters that are not hyper-ellipsoids (or hyper-spheres).



(A): Two natural clusters



(B): k -means clusters



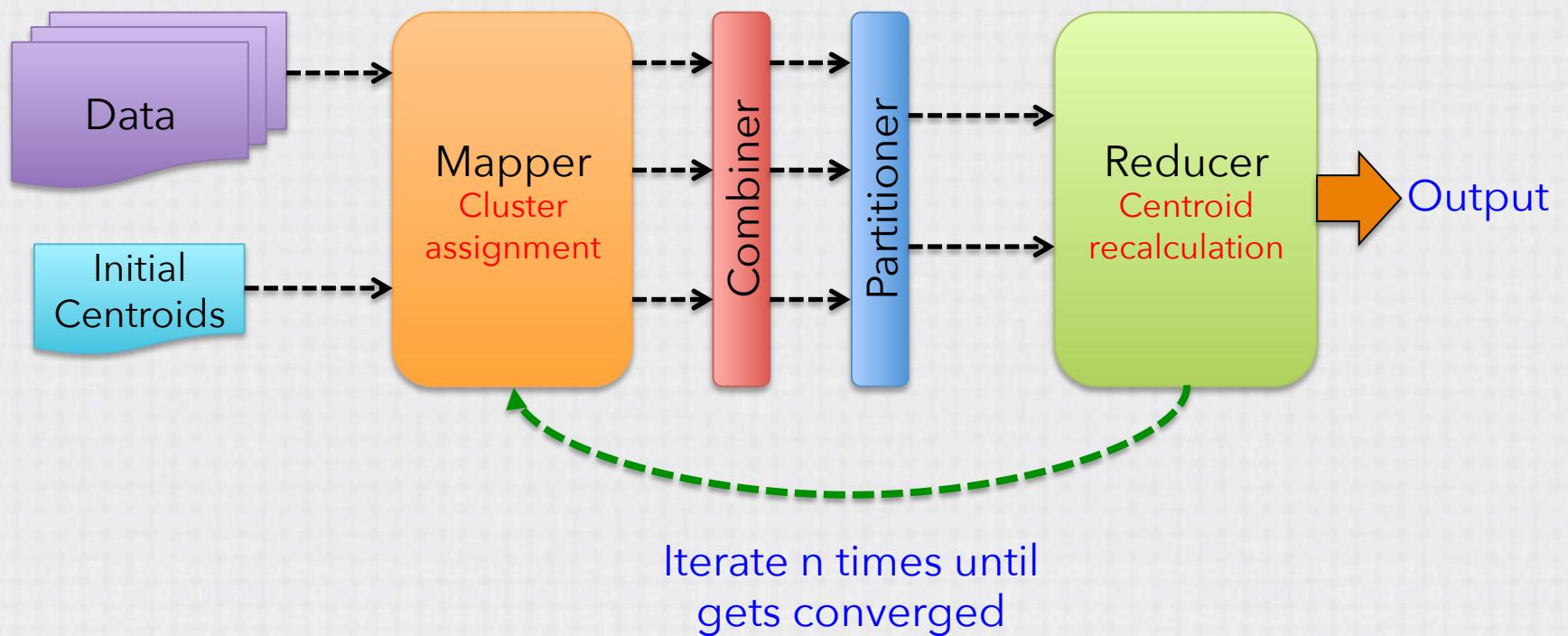
UNIVERSITY OF
MARYLAND

ROBERT H. SMITH

SCHOOL OF BUSINESS

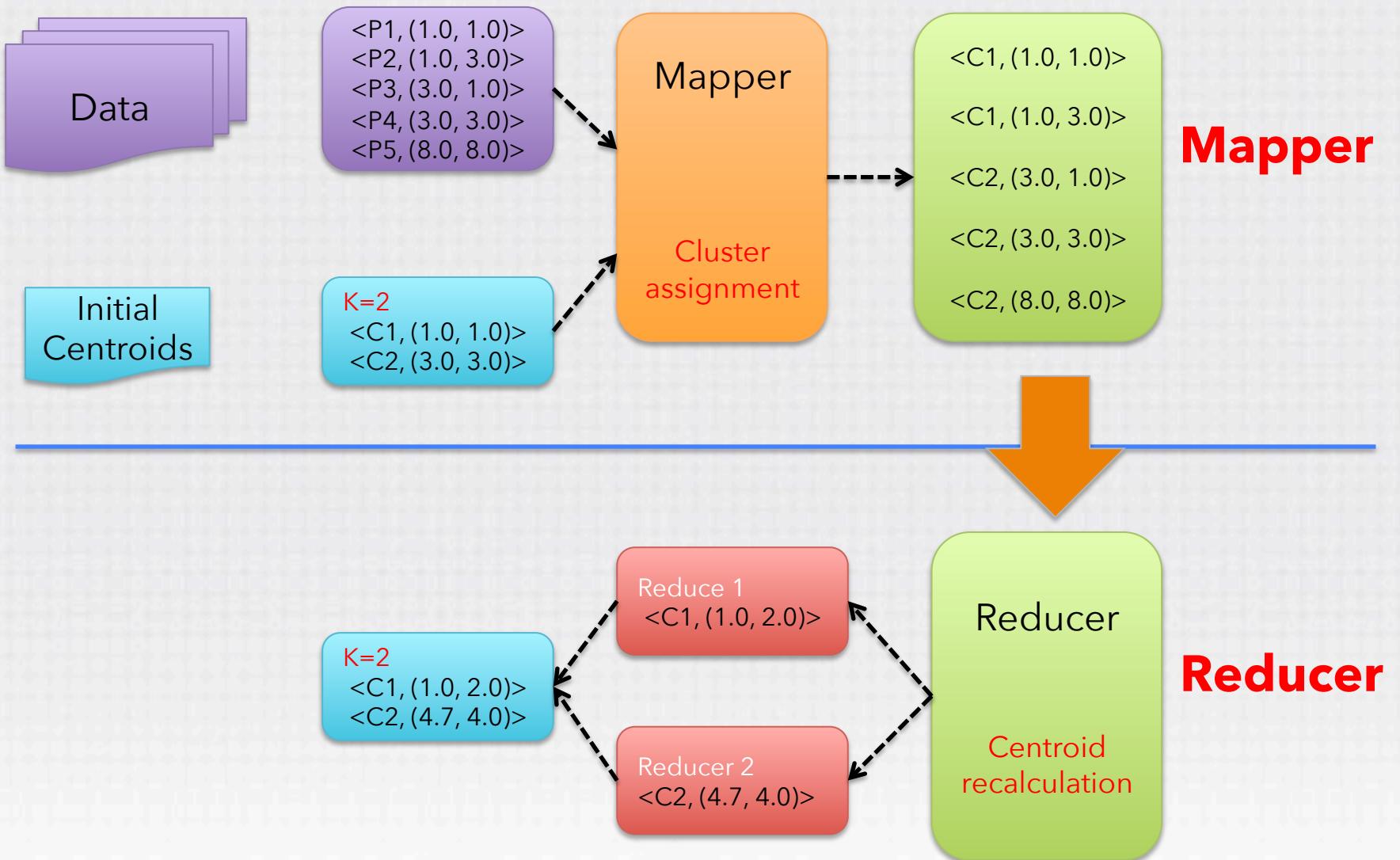
K-Means in MapReduce

K-Means under MapReduce



Data never changes in the MapReduce process

Mapper and Reducer



K-means summary

- Despite weaknesses, *k*-means is still the most popular algorithm due to its simplicity, efficiency and
 - ❑ other clustering algorithms have their own lists of weaknesses.
- No clear evidence that any other clustering algorithm performs better in general
 - ❑ although they may be more suitable for some specific types of data or applications.
- Comparing different clustering algorithms is a difficult task. No one knows the correct clusters!



UNIVERSITY OF
MARYLAND

ROBERT H. SMITH

SCHOOL OF BUSINESS

Hierarchical Clustering

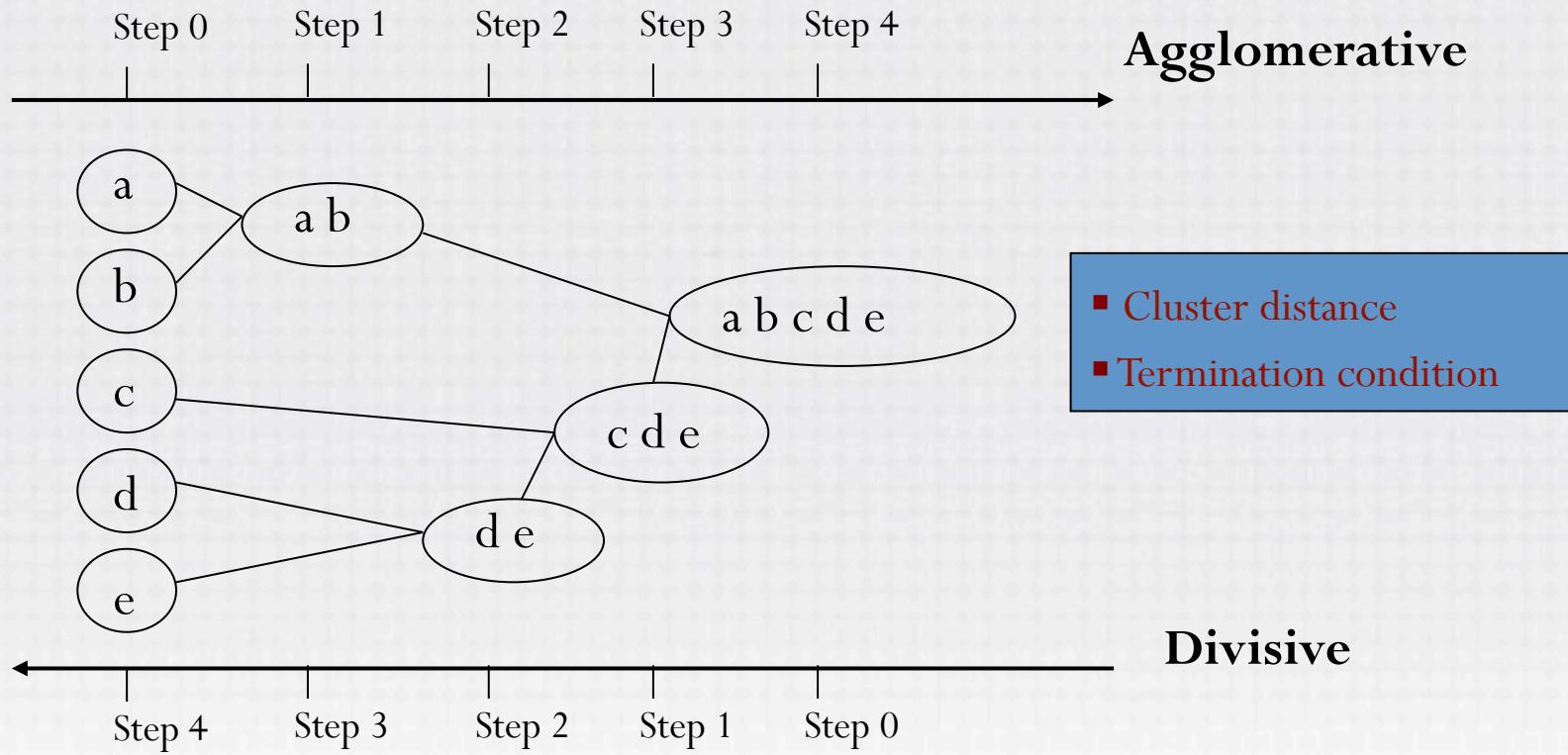
Hierarchical clustering

- Hierarchical Clustering Approach
 - A typical clustering analysis approach via partitioning data set **sequentially**
 - Construct nested partitions layer by layer via grouping objects into a tree of clusters (**without the need to know the number of clusters in advance**)
 - Uses distance matrix as clustering criteria
- Agglomerative vs. Divisive
 - Two sequential clustering strategies for constructing a tree of clusters
 - **Agglomerative: a bottom-up strategy**
 - Initially each data object is in its own (atomic) cluster
 - Then merge these atomic clusters into larger and larger clusters
 - **Divisive: a top-down strategy**
 - Initially all objects are in one single cluster
 - Then the cluster is subdivided into smaller and smaller clusters

Introduction

- Illustrative Example

Agglomerative and divisive clustering on the data set {a, b, c, d ,e }

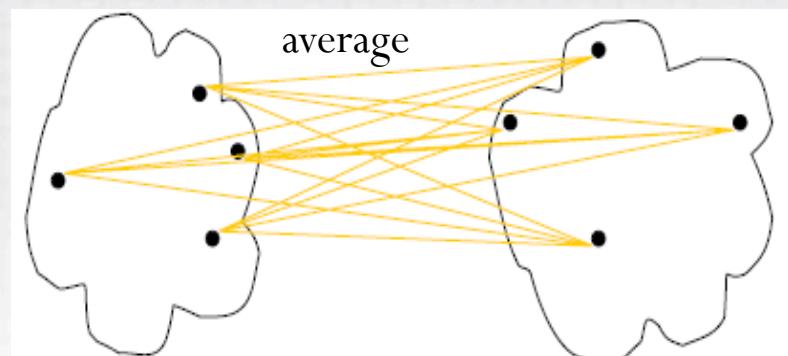
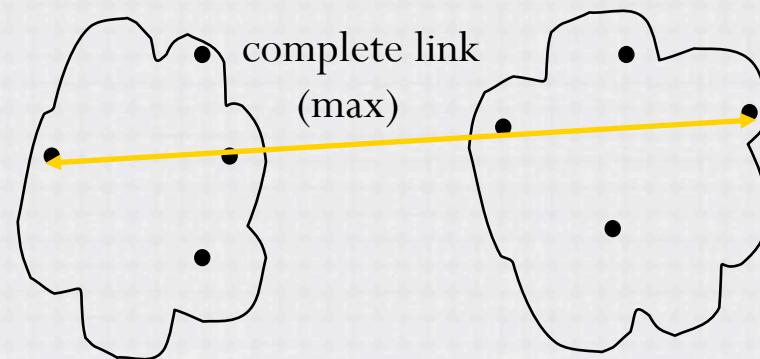
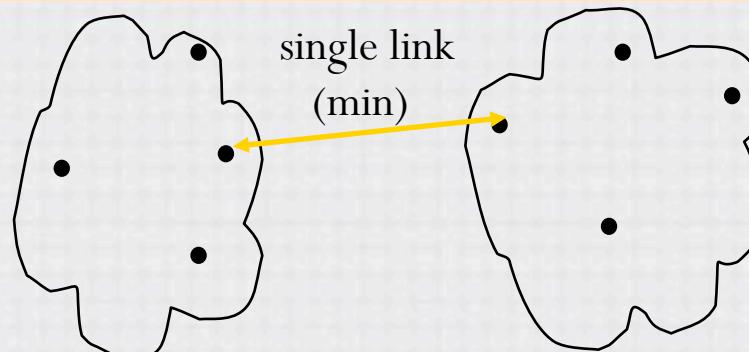


Cluster distance measures



ROBERT H. SMITH
SCHOOL OF BUSINESS

- **Single link:** smallest distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \min \{d(x_{ip}, x_{jq})\}$
- **Complete link:** largest distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \max \{d(x_{ip}, x_{jq})\}$
- **Average:** avg distance between elements in one cluster and elements in the other, i.e.,
$$d(C_i, C_j) = \text{avg} \{d(x_{ip}, x_{jq})\}$$



Cluster distance measures



ROBERT H. SMITH
SCHOOL OF BUSINESS

Example: Given a data set of five objects characterised by a single feature, assume that there are two clusters: $C_1: \{a, b\}$ and $C_2: \{c, d, e\}$.

	a	b	c	d	e
Feature	1	2	4	5	6

1. Calculate the distance matrix.

	a	b	c	d	e
a	0	1	3	4	5
b	1	0	2	3	4
c	3	2	0	1	2
d	4	3	1	0	1
e	5	4	2	1	0

2. Calculate three cluster distances between C_1 and C_2 .

Single link

$$\begin{aligned} \text{dist}(C_1, C_2) &= \min\{d(a, c), d(a, d), d(a, e), d(b, c), d(b, d), d(b, e)\} \\ &= \min\{3, 4, 5, 2, 3, 4\} = 2 \end{aligned}$$

Complete link

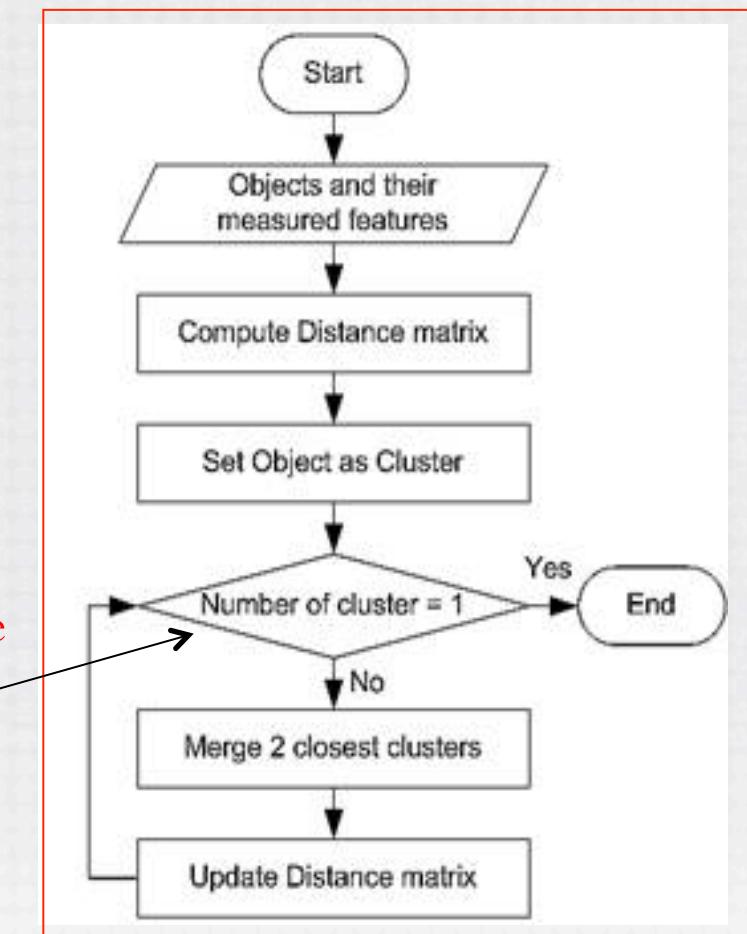
$$\begin{aligned} \text{dist}(C_1, C_2) &= \max\{d(a, c), d(a, d), d(a, e), d(b, c), d(b, d), d(b, e)\} \\ &= \max\{3, 4, 5, 2, 3, 4\} = 5 \end{aligned}$$

Average

$$\begin{aligned} \text{dist}(C_1, C_2) &= \frac{d(a, c) + d(a, d) + d(a, e) + d(b, c) + d(b, d) + d(b, e)}{6} \\ &= \frac{3 + 4 + 5 + 2 + 3 + 4}{6} = \frac{21}{6} = 3.5 \end{aligned}$$

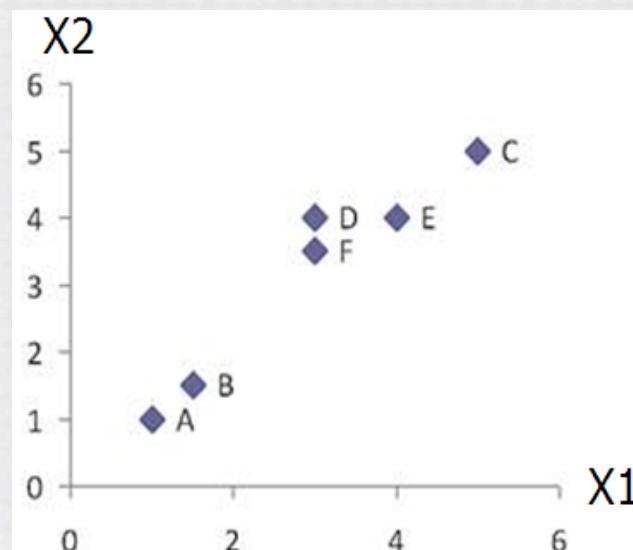
Agglomerative algorithm

- The *Agglomerative* algorithm is carried out in three steps:
 - Convert object attributes to distance matrix
 - Set each object as a cluster (thus if we have N objects, we will have N clusters at the beginning)
 - Repeat until number of cluster is one (or known # of clusters)
 - Merge two closest clusters
 - Update distance matrix



Example

- Problem: clustering analysis with agglomerative algorithm



$$d_{AB} = \sqrt{(1-1.5)^2 + (1-1.5)^2} = \sqrt{\frac{1}{2}} = 0.7071$$

$$d_{DF} = \sqrt{(3-3)^2 + (4-3.5)^2} = 0.5$$

Euclidean distance

	X1	X2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5

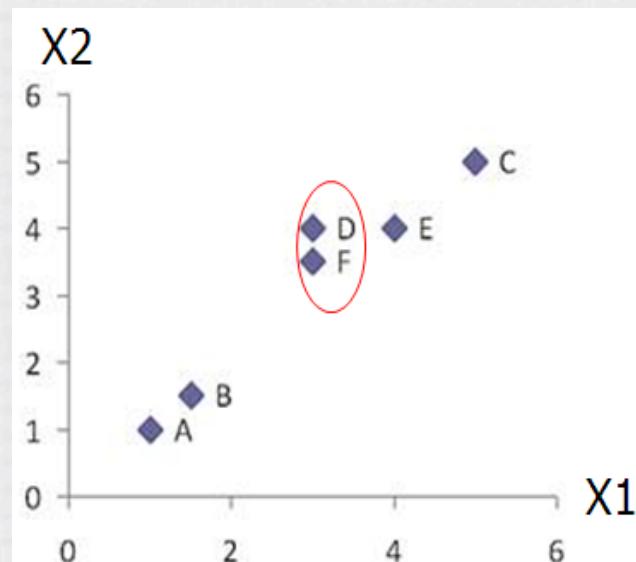
data matrix

Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

distance matrix

Example

- Merge two closest clusters (iteration 1)



Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	?	4.24
B	0.71	0.00	4.95	?	3.54
C	5.66	4.95	0.00	?	1.41
D, F	?	?	?	0.00	?
E	4.24	3.54	1.41	?	0.00

Example

- Update distance matrix (iteration 1)

Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B	0.71	0.00	4.95	2.92	3.54	2.50
C	5.66	4.95	0.00	2.24	1.41	2.50
D	3.61	2.92	2.24	0.00	1.00	0.50
E	4.24	3.54	1.41	1.00	0.00	1.12
F	3.20	2.50	2.50	0.50	1.12	0.00

$$d_{(D,F) \rightarrow A} = \min(d_{DA}, d_{FA}) = \min(3.61, 3.20) = 3.20$$

$$d_{(D,F) \rightarrow B} = \min(d_{DB}, d_{FB}) = \min(2.92, 2.50) = 2.50$$

$$d_{(D,F) \rightarrow C} = \min(d_{DC}, d_{FC}) = \min(2.24, 2.50) = 2.24$$

$$d_{E \rightarrow (D,F)} = \min(d_{ED}, d_{EF}) = \min(1.00, 1.12) = 1.00$$

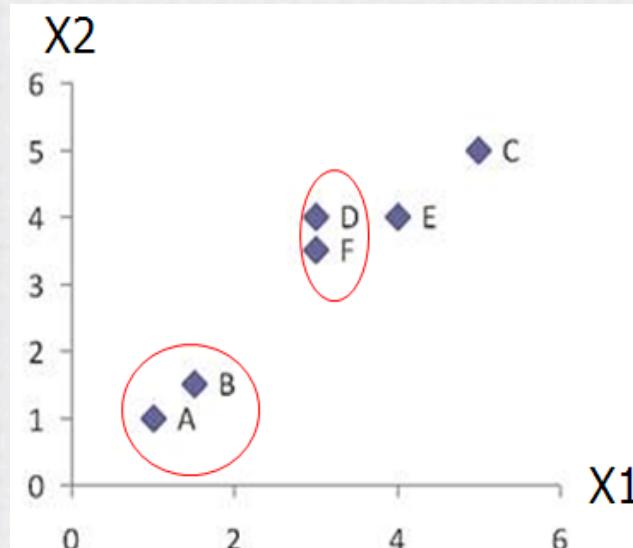
Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	?	4.24
B	0.71	0.00	4.95	?	3.54
C	5.66	4.95	0.00	?	1.41
D, F	?	?	?	0.00	?
E	4.24	3.54	1.41	?	0.00

Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00

Example

- Merge two closest clusters (iteration 2)



Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00

Dist	A,B	C	(D, F)	E
A,B	0	?	?	?
C	?	0	2.24	1.41
(D, F)	?	2.24	0	1.00
E	?	1.41	1.00	0

Example

- Update distance matrix (iteration 2)

Min Distance (Single Linkage)

Dist	A	B	C	D, F	E
A	0.00	0.71	5.66	3.20	4.24
B	0.71	0.00	4.95	2.50	3.54
C	5.66	4.95	0.00	2.24	1.41
D, F	3.20	2.50	2.24	0.00	1.00
E	4.24	3.54	1.41	1.00	0.00

$$d_{C \rightarrow (A,B)} = \min(d_{CA}, d_{CB}) = \min(5.66, 4.95) = 4.95$$

$$\begin{aligned} d_{(D,F) \rightarrow (A,B)} &= \min(d_{DA}, d_{DB}, d_{FA}, d_{FB}) \\ &= \min(3.61, 2.92, 3.20, 2.50) = 2.50 \end{aligned}$$

$$d_{E \rightarrow (A,B)} = \min(d_{EA}, d_{EB}) = \min(4.24, 3.54) = 3.54$$

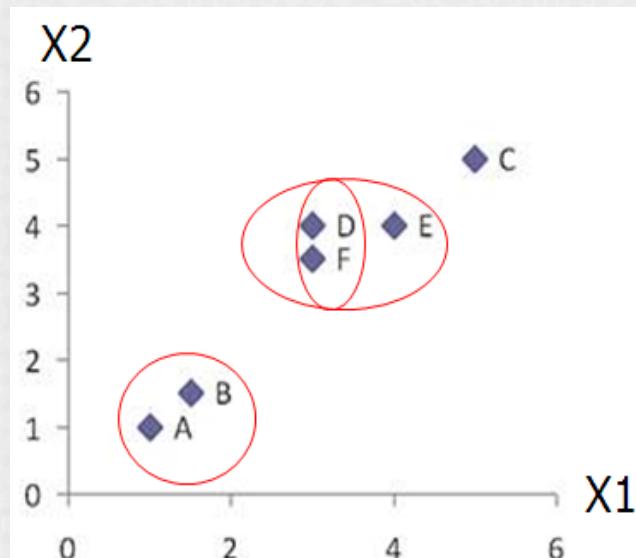
Dist	A,B	C	(D, F)	E
A,B	0	?	?	?
C	?	0	2.24	1.41
(D, F)	?	2.24	0	1.00
E	?	1.41	1.00	0

Min Distance (Single Linkage)

Dist	A,B	C	(D, F)	E
A,B	0	4.95	2.50	3.54
C	4.95	0	2.24	1.41
(D, F)	2.50	2.24	0	1.00
E	3.54	1.41	1.00	0

Example

- Merge two closest clusters/update distance matrix
(iteration 3)



Min Distance (Single Linkage)

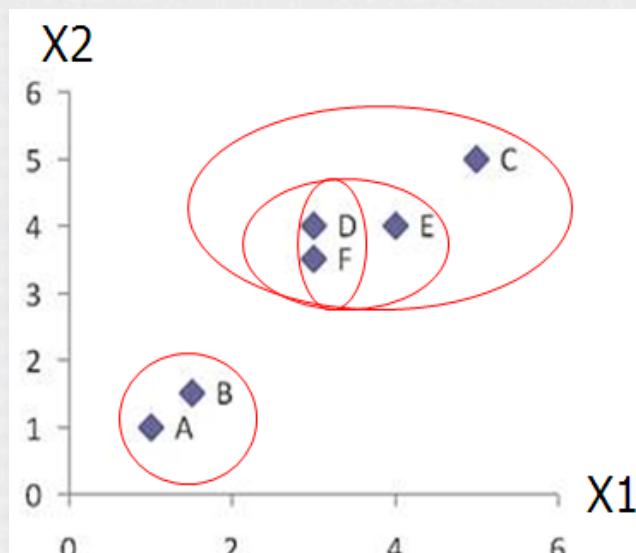
Dist	A,B	C	(D, F)	E
A,B	0	4.95	2.50	3.54
C	4.95	0	2.24	1.41
(D, F)	2.50	2.24	0	1.00
E	3.54	1.41	1.00	0

Min Distance (Single Linkage)

Dist	(A,B)	C	(D, F), E
(A,B)	0.00	4.95	2.50
C	4.95	0.00	1.41
(D, F), E	2.50	1.41	0.00

Example

- Merge two closest clusters/update distance matrix
(iteration 4)



Min Distance (Single Linkage)

Dist	(A,B)	C	(D, F), E
(A,B)	0.00	4.95	2.50
C	4.95	0.00	1.41
(D, F), E	2.50	1.41	0.00

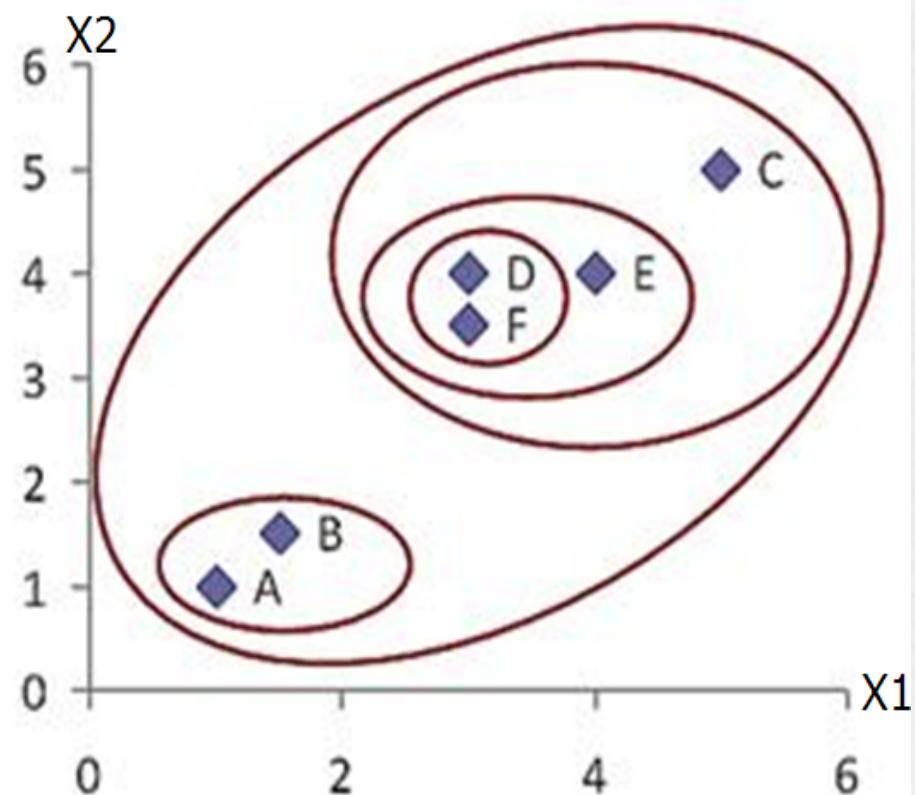
Min Distance (Single Linkage)

Dist	(A,B)	((D, F), E), C
(A,B)	0.00	2.50
((D, F), E), C	2.50	0.00

Example

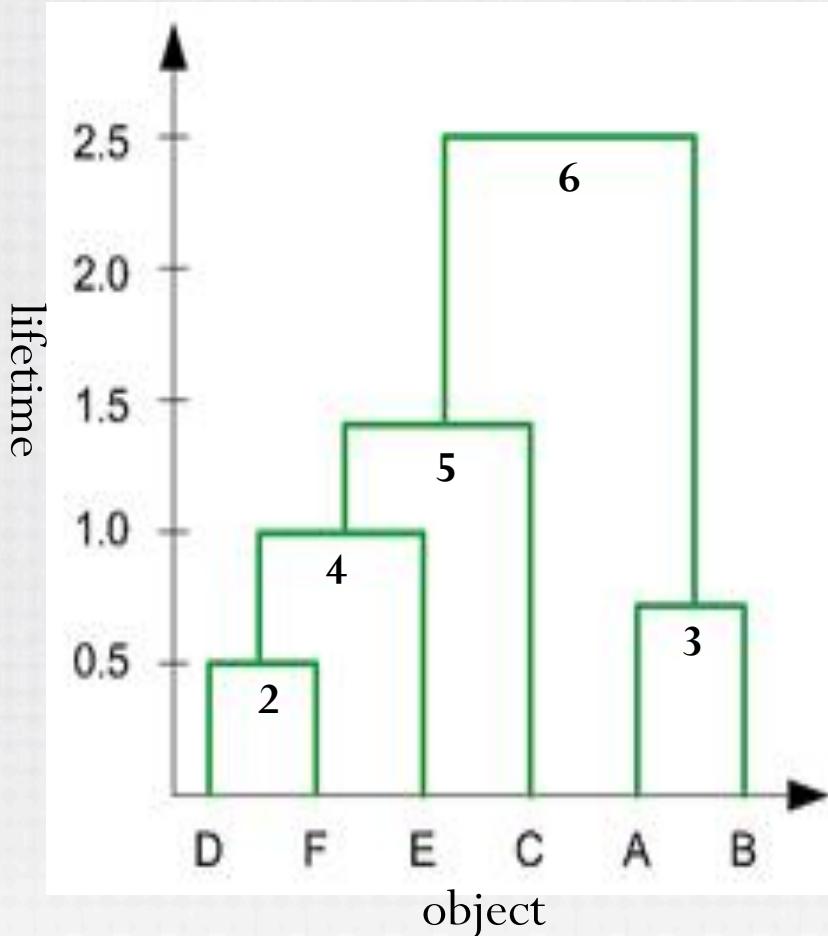
- Final result (meeting termination condition)

	X1	X2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5



Example

- Dendrogram tree representation



- In the beginning we have 6 clusters: A, B, C, D, E and F
- We merge clusters D and F into cluster (D, F) at distance 0.50
- We merge cluster A and cluster B into (A, B) at distance 0.71
- We merge clusters E and (D, F) into ((D, F), E) at distance 1.00
- We merge clusters ((D, F), E) and C into (((D, F), E), C) at distance 1.41
- We merge clusters (((D, F), E), C) and (A, B) into ((((D, F), E), C), (A, B)) at distance 2.50
- The last cluster contain all the objects, thus conclude the computation

Major issue - labeling

- After clustering algorithm finds clusters - how can they be useful to the end user?
- Need pithy label for each cluster
 - ❑ In search results, say “Animal” or “Car” in the *jaguar* example.
 - ❑ In topic trees, need navigational cues.
 - Often done by hand.

How to label clusters

- Show titles of typical documents
 - Titles are easy to scan
 - Authors create them for quick scanning!
 - But you can only show a few titles which may not fully represent cluster
- Show words/phrases prominent in cluster
 - More likely to fully represent cluster
 - Use distinguishing words/phrases
 - Differential labeling

Labeling

- Common heuristics - list 5-10 most frequent terms in the centroid vector.
 - Drop stop-words; stemming.
- Differential labeling by frequent terms
 - Within a collection “Computers”, clusters all have the word *computer* as frequent term.
- Perhaps better: distinctive noun phrase

What is a good clustering?

- *Internal criterion:* A good clustering will produce high quality clusters in which:
 - ❑ the intra-class (that is, intra-cluster) similarity is high
 - ❑ the inter-class similarity is low
 - ❑ The measured quality of a clustering depends on both the document representation and the similarity measure used

External criteria for clustering quality



- Quality measured by its ability to discover some or all of the hidden patterns or latent classes in gold standard data
- Assesses a clustering with respect to ground truth
- Assume documents with C gold standard classes, while our clustering algorithms produce K clusters, $\omega_1, \omega_2, \dots, \omega_K$ with n_1, n_2, \dots, n_k members.

External evaluation of cluster quality

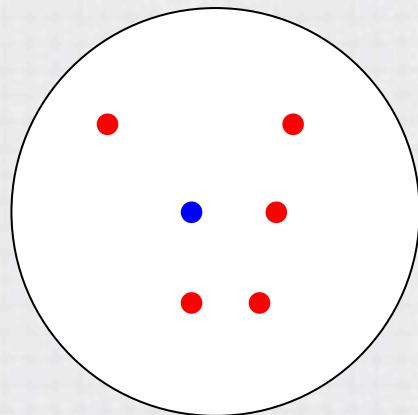


- *Simple measure:* purity, the ratio between the dominant class in the cluster π_i and the size of cluster ω_i

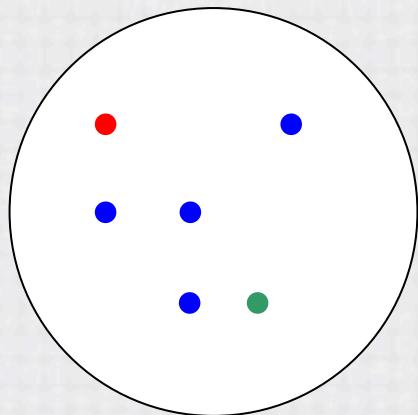
$$Purity(\omega_i) = \frac{1}{n_i} \max_j (n_{ij}) \quad j \in C$$

- Others are entropy of classes in clusters (or mutual information between classes and clusters)

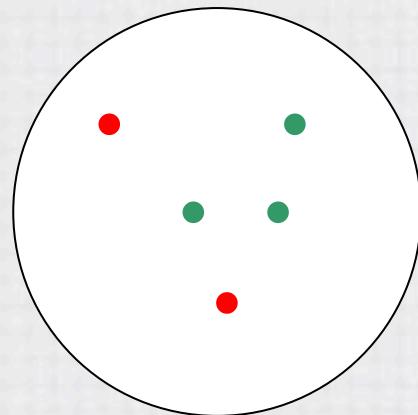
Purity example



Cluster I



Cluster II



Cluster III

Cluster I: Purity = $1/6 \max(5, 1, 0) = 5/6$

Cluster II: Purity = $1/6 \max(1, 4, 1) = 4/6$

Cluster III: Purity = $1/5 \max(2, 0, 3) = 3/5$

Conclusions

- Hierarchical algorithm is a sequential clustering algorithm
 - Use distance matrix to construct a tree of clusters (dendrogram)
 - Hierarchical representation without the need of knowing # of clusters (can set termination condition with known # of clusters)
- Major weakness of agglomerative clustering methods
 - Sensitive to cluster distance measures and noise/outliers
 - Less efficient: $O(n^2)$, where n is the number of total objects

Hierarchical clustering in MapReduce



UNIVERSITY OF
MARYLAND
ROBERT H. SMITH
SCHOOL OF BUSINESS

- Divide this task into following jobs
 - ❑ Calculate the distance matrix
 - ❑ Find the minimum distance entry in the matrix
 - ❑ Merge and update the matrix
- Chaining jobs using their dependence relations

Example: hierarchical agglomerative clustering on documents

1. Build term dictionary
2. Term frequency normalization for each document
3. Construct distance matrix
4. Calculate similarity of all pairs
5. Find maximum similarity and combine these two docs to a “new” doc.