

# Chapter 12 - Cluster Analysis

12.1

- Goal: Group individual observations into a set of "natural" sets. Unlike Discriminant Analysis, groups are not "known" in advance. Grouping made on similarities/distances between individual cases.

## 12.2 - Similarity Measures

Distances for Pairs of Items

- Euclidean:  $d(\underline{x}, \underline{y}) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2} = \sqrt{(\underline{x} - \underline{y})'(\underline{x} - \underline{y})}$

- Statistical:  $d(\underline{x}, \underline{y}) = \sqrt{(\underline{x} - \underline{y})' A (\underline{x} - \underline{y})}$   $A \equiv S^{-1}$  typically  
Problematic when groups are not defined.

- Minkowski:  $d(\underline{x}, \underline{y}) = \left[ \sum_{i=1}^p |x_i - y_i|^m \right]^{1/m}$

### Binary outcomes

	Variable				
	1	2	3	4	5
item i	1	0	0	1	1
item k	1	1	0	1	0

$x_{ij} \equiv$  score for  $j^{\text{th}}$  var on item  $i$

$x_{kj} \equiv$  " " " " " "  $k$

$$\sum_{j=1}^p (x_{ij} - x_{kj})^2 = (1-1)^2 + (0-1)^2 + (0-0)^2 + (1-1)^2 + (1-0)^2 = 2$$

$\equiv$  # mismatches (weights 1-1 and 0-0 matches equally)

Contingency table for items i, k

		Item k		
		1	0	Totals
Item i	1	a	b	a+b
	0	c	d	c+d
Totals		a+c	b+d	a+b+c+d=p

Table 12.1 Similarity Coefficients for Clustering Items (binaries)

- 1)  $\frac{a+d}{p}$  (Equal weights for 1-1, 0-0 matches)
- 2)  $\frac{2(a+d)}{2(a+d)+b+c}$  (Double weight for 1-1, 0-0 matches)
- 3)  $\frac{a+d}{a+d+2(b+c)}$  (Double weight for mis-matches)
- 4)  $\frac{a}{p}$  (No 0-0 matches in numerator)
- 5)  $\frac{a}{a+b+c}$  (0-0 matches removed)
- 6)  $\frac{2a}{2a+b+c}$  (0-0 removed, double weight for 1-1)
- 7)  $\frac{a}{a+2(b+c)}$  (0-0 removed, double weight for mis-matches)
- 8)  $\frac{a}{b+c}$  (0-0 removed, ratio of matches/mis-matches).

Example - Diversity of Artifacts of 8 forts in Western Canada

~~Total Artifact Diversity~~ ~~Lux goods div.~~ ~~# of artifacts~~ ~~occupation years~~

- Vars: (1) Total Artifact Diversity (# of Types)  $X_1 = \begin{cases} 1 & \text{if } \geq 48 \\ 0 & \text{if } < 48 \end{cases}$
- (2) Lux Goods Diversity  $X_2 = \begin{cases} 1 & \text{if } \geq 10 \\ 0 & \text{if } < 10 \end{cases}$ , (3) # artifacts  $X_3 = \begin{cases} 1 & \text{if } \geq 1500 \\ 0 & \text{if } < 1500 \end{cases}$
- (4) Occupation years  $X_4 = \begin{cases} 1 & \text{if } \geq 10 \\ 0 & \text{if } < 10 \end{cases}$  (5) Ethnicity  $X_5 = \begin{cases} 1 & \text{if } \text{Orkney} \\ 0 & \text{if } \text{French Canadian} \end{cases}$  (6) Region  $X_6 = \begin{cases} 1 & \text{if } \text{Saskatchewan} \\ 0 & \text{if } \text{Alberta} \end{cases}$

Data

Fort	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>
1	0	0	1	0	0	0
2	0	0	0	0	0	0
3	1	1	1	0	0	0
4	0	1	0	0	1	0
5	1	0	0	1	1	1
6	0	0	1	0	0	1
7	1	1	1	1	1	0
8	0	0	0	1	0	0

	Fort 1	Fort 2	Total
Fort 1	3	1	4
Fort 2	0	2	2
Total	3	3	6

Similarity index 1:  $\frac{a+d}{p} = \frac{3+2}{6} = \frac{5}{6}$

Similarity matrix (measure 1)

Fort	1	2	3	4	5	6	7	8
1	1							
2	$\frac{5}{6}$	1						
3	$\frac{3}{6}$	$\frac{2}{6}$	1					
4	$\frac{2}{6}$	$\frac{3}{6}$	$\frac{1}{6}$	1				
5	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{2}{6}$	$\frac{3}{6}$	1			
6	$\frac{3}{6}$	$\frac{4}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{2}{6}$	1		
7	$\frac{2}{6}$	$\frac{1}{6}$	$\frac{5}{6}$	$\frac{4}{6}$	$\frac{3}{6}$	$\frac{1}{6}$	1	
8	$\frac{4}{6}$	$\frac{5}{6}$	$\frac{0}{6}$	$\frac{2}{6}$	$\frac{3}{6}$	$\frac{3}{6}$	$\frac{2}{6}$	1

Highest Similarities {1,2}, {1,6}, {2,8}, {3,4}, {3,7}

Possible Subgroups {1,2,6,8}, {3,4,7}, {5}

Similarity Measures:  $\tilde{S}_{ik} = \frac{1}{1+d_{ik}}$   $d_{ik} \equiv$  distance measure

$$0 < \tilde{S}_{ik} \leq 1$$

If Similarity matrix is nnd, max Similarity  $\tilde{S}_{ii} = 1$ ,

then  $d_{ik} = \sqrt{2(1-\tilde{S}_{ik})}$  has properties of distance

### Similarity and Association Measures for Pairs of Variables

Binary Case

		Variable k		
		1	0	
Variable i	1	a	b	a+b
	0	c	d	c+d
		a+c	b+d	a+b+c+d = n

$$r = \frac{ad-bc}{\sqrt{(a+b)(c+d)(a+c)(b+d)}}$$

~~For example on variables~~ For example on variables

		$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$
Variable	$X_1$	8					
	$X_2$	7	8				
	$X_3$	4	5	8			
	$X_4$	5	4	3	8		
	$X_5$	8	7	4	5	8	
	$X_6$	4	3	4	5	4	8

## 12.3 Hierarchical Clustering Methods

• Agglomerative - Start w/ individual objects (units or variables) and combine until a single cluster formed.

### Linkage Methods for Combining clusters

- Single Linkage - Minimum distance between objects in clusters.
- Complete Linkage - Maximum distance between objects in clusters
- Average Linkage - Mean distance between objects in clusters.

• Divisive - start w/ single cluster and split off until each object is a cluster

Dendrogram - 2-Dimensional Diagram of the Clustering process.

EXAMPLE - AVERAGE POINTS, REBOUNDS, ASSISTS of  $n=5$  WNBA

Players (2014): Maya Moore, Skylar Diggins, Candace Parker, Angel McCoughtry, Tina Charles

$$\text{Distance} = \sqrt{\sum_{j=1}^3 (X_{ij} - Y_j)^2}$$

	MM	SD	CP	AM	TC
MM	0				
SD	6.92	0			
CP	4.63	4.71	0		
AM	6.09	3.40	2.21	0	
TC	6.78	7.93	3.73	4.59	0

## Single Linkage

- 2 Closest Items: C, A  $d_{CA} = 2.21$

$\Rightarrow$  Merge C, A into cluster (CA)

$$d_{(CA)M} = \min\{d_{CM}, d_{AM}\} = \min\{4.63, 6.09\} = 4.63$$

$$d_{(CA)S} = \min\{d_{CS}, d_{AS}\} = \min\{4.71, 3.40\} = 3.40$$

$$d_{(CA)T} = \min\{d_{CT}, d_{AT}\} = \min\{3.73, 4.59\} = 3.73$$

$$D^* =$$

	(CA)	M	S	T
(CA)	0			
M	4.63	0		
S	3.40	6.92	0	
T	3.73	6.78	7.93	0

Smallest distance in  $D^* = d_{(CA)S} = 3.40$

$\Rightarrow$  merge (CA, S) into (CAS)

$$d_{(CAS)M} = \min\{d_{(CA)M}, d_{SM}\} = \min\{4.63, 6.92\} = 4.63$$

$$d_{(CAS)T} = \min\{d_{(CA)T}, d_{ST}\} = \min\{3.73, 7.93\} = 3.73$$

$$D^{**}$$

	(CAS)	M	T
(CAS)	0		
M	4.63	0	
T	3.73	6.78	0

Smallest distance in

$$D^{**} = d_{(CAS)T} \Rightarrow \text{merge}$$

(CAS, T) into (CAST)

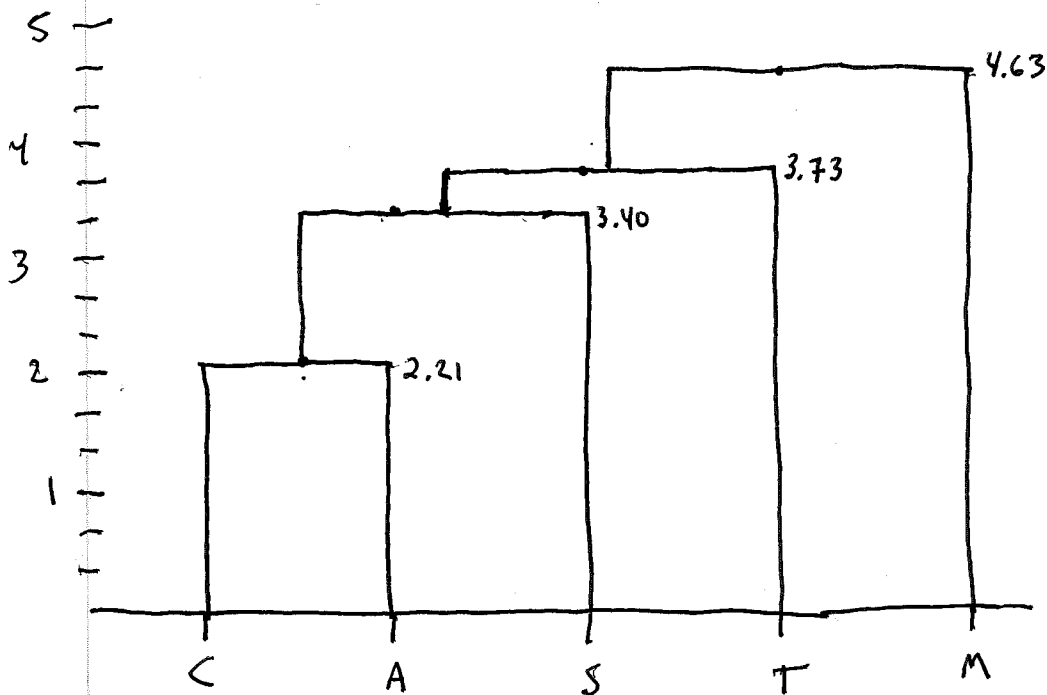
$$D^{***}$$

	(CAST)	M
(CAST)	0	
M	4.63	0

$$d_{(CAST)M} = \min\{d_{(CAS)M}, d_{TM}\} = \min\{4.63, 6.78\} = 4.63$$

# Dendrogram

12.7



## Complete Linkage

	M	S	C	A	T
M	0				
S	6.92	0			
C	4.63	4.71	0		
A	6.09	3.40	2.21	0	
T	6.78	7.93	3.73	4.59	0

Step 1  $\Rightarrow$  CA (2.21)

(Same as single linkage)

$$d_{CA(M)} = \max \{d_{CM}, d_{AM}\} = \max \{4.63, 6.09\} = 6.09$$

$$d_{CA(S)} = \max \{d_{CS}, d_{AS}\} = \max \{4.71, 3.40\} = 4.71$$

$$d_{CA(T)} = \max \{d_{CT}, d_{AT}\} = \max \{3.73, 4.59\} = 4.59$$

$\Rightarrow D^* =$

	(CA)	M	S	T
(CA)	0			
M	6.09	0		
S	4.71	6.92	0	
T	4.59	6.78	7.93	0

$\Rightarrow$  Combine CA, T into (CAT)

$$d_{(CAT)M} = \max \{ d_{(CAT)M}, d_{TM} \} = \max \{ 6.09, 6.78 \} = 6.78$$

$$d_{(CAT)S} = \max \{ d_{(CAT)S}, d_{TS} \} = \max \{ 4.71, 7.93 \} = 7.93$$

$D^{**}$

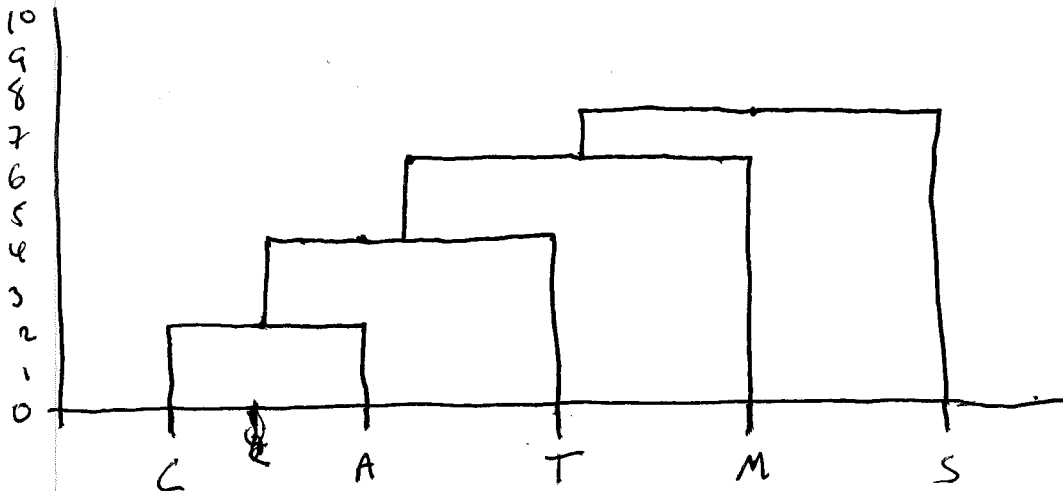
	(CAT)	M	S
(CAT)	0		
M	6.78	0	
S	7.93	6.92	0

⇒ Combine CAT, M into CATM

$$d_{(CATM)S} = \max \{ d_{(CAT)S}, d_{MS} \} = \max \{ 7.93, 6.92 \} = 7.93$$

$D^{***}$

	(CATM)	S
(CATM)	0	
S	7.93	0





Average Linkage

	M	S	C	A	T
D					
M	0				
S	6.92	0			
C	4.63	4.71	0		
A	6.09	3.40	2.21	0	
T	6.78	7.93	3.73	4.59	0

Again choose to combine C, A into (CA)

$$d_{(CA)M} = \frac{d_{cm} + d_{am}}{2} = \frac{4.63 + 6.09}{2} = 5.36$$

$$d_{(CA)S} = \frac{d_{cs} + d_{as}}{2} = \frac{4.71 + 3.40}{2} = 4.06$$

$$d_{(CA)T} = \frac{d_{ct} + d_{at}}{2} = \frac{3.73 + 4.59}{2} = 4.16$$

	(CA)	M	S	T
D'				
(CA)	0			
M	5.36	0		
S	4.06	6.92	0	
T	4.16	6.78	7.93	0

Combine (CA), S into (CAS)

$$d_{(CAS)M} = \frac{d_{(CA)M} + d_{sm}}{2} = \frac{5.36 + 6.92}{2} = 6.14$$

$$d_{(CAS)T} = \frac{d_{(CA)T} + d_{st}}{2} = \frac{4.16 + 7.93}{2} = 6.05$$

	(CAS)	M	T
D''			
(CAS)	0		
M	6.14	0	
T	6.05	6.78	0

Combine (CAS), T into (CAST)

Average Linkage

	M	S	C	A	T
D					
M	0				
S	6.92	0			
C	4.63	4.71	0		
A	6.09	3.40	2.21	0	
T	6.78	7.93	3.73	4.59	0

Again choose to combine C, A into (CA)

$$d_{(CA)M} = \frac{d_{CM} + d_{AM}}{2} = \frac{4.63 + 6.09}{2} = 5.36$$

$$d_{(CA)S} = \frac{d_{CS} + d_{AS}}{2} = \frac{4.71 + 3.40}{2} = 4.06$$

$$d_{(CA)T} = \frac{d_{CT} + d_{AT}}{2} = \frac{3.73 + 4.59}{2} = 4.16$$

	(CA)	M	S	T
D*				
(CA)	0			
M	5.36	0		
S	4.06	6.92	0	
T	4.16	6.78	7.93	0

Combine (CA), S into (CAS)

$$d_{(CAS)M} = \frac{d_{CAM} + d_{SM}}{2} = \frac{5.36 + 6.92}{2} = 6.14$$

$$d_{(CAS)T} = \frac{d_{CA,T} + d_{ST}}{2} = \frac{4.16 + 7.93}{2} = 6.05$$

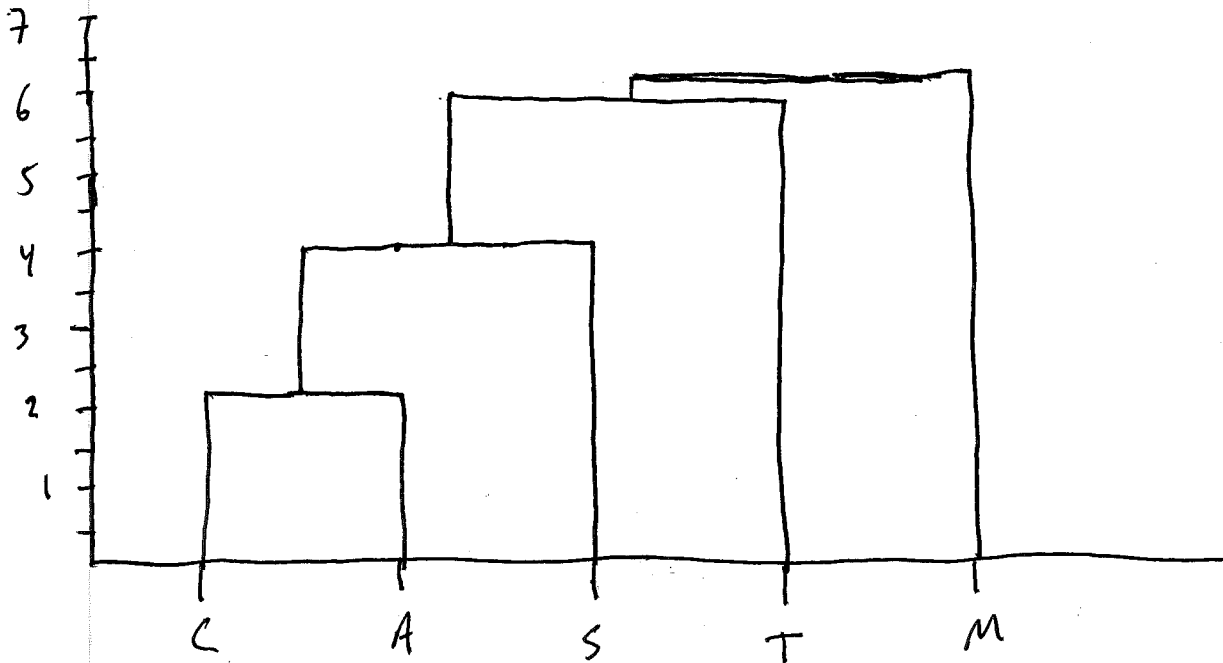
	(CAS)	M	T
D**			
(CAS)	0		
M	6.14	0	
T	6.05	6.78	0

Combine (CAS), T into (CAST)

$$d_{(AST)M} = \frac{d_{(AS)M} + d_{TM}}{2} = \frac{6.14 + 6.78}{2} = 6.46$$

D<sup>\*\*\*</sup>

	(AST)	M
(AST)	0	
M	6.46	0



## 12.4 - Nonhierarchical Clustering Methods

Goal - Group items (not variables) into a group of  $K$  clusters.  $K$  can be specified in advance or determined from the procedure.

- Distance/Similarity matrix does not need to be computed.
- Basic data not stored during run (computationally more efficient than hierarchical methods).

## K-Means Method

- 1) Partition the items into  $K$  initial clusters (random)
- 2) Proceed through list of items, assigning each to the nearest cluster centroid. (Euclidean distance w/ standardized or unstandardized observations).  
Recalculate cluster centroids when observation shifts clusters.
- 3) Repeat Step 2) until no observations are reassigned.

Several initial partitions should be used, to see whether the clusters are stable.

$$\text{Optimization Problem: } \min_{c_1, \dots, c_K} \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}$$

where  $|C_k| \equiv \#$  of items in cluster  $k$ .