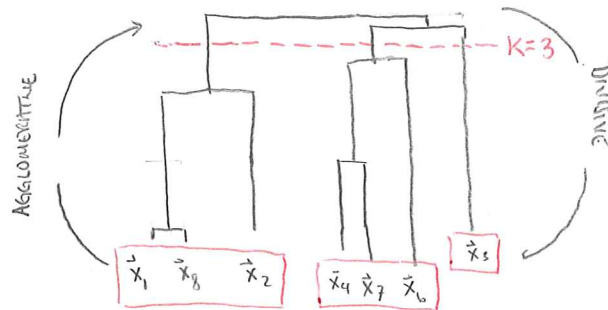


HIERARCHICAL CLUSTERING

SETTING: N OBSERVATIONS $\vec{x}_1, \dots, \vec{x}_N \in \mathbb{R}^d$

GOAL: FIND CLUSTERS (GROUPING/SUBSET) OF POINTS WHICH ARE SIMILAR

OUTPUT: DENDROGRAM (HIERARCHICAL TREE DIAGRAM)



- AGGLOMERATIVE: K -CLUSTER SOLUTION GENERATED BY MERGING TWO CLUSTERS FROM $K+1$ CLUSTER SOLUTION (BOTTOM UP)
- DIVISIVE: $(k+1)$ -CLUSTER SOLUTION GENERATED BY SPLITTING ONE CLUSTER FROM K -CLUSTER SOLUTION (TOP DOWN)
- GIVEN A DENDROGRAM, WE CAN DRAW A HORIZONTAL LINE CROSSING K VERTICAL LINES TO GENERATE K CLUSTERS OF DATA (BRANCHES)
 - ITEMS WHICH ARE SIMILAR ARE COMBINED AT LOW HEIGHTS
 - DIFFERENCE IN HEIGHTS DEFINE HOW CLOSE POINTS ARE TOGETHER
 - \vec{x}_1, \vec{x}_8 ARE MORE SIMILAR THAN \vec{x}_4, \vec{x}_7
 - \vec{x}_1, \vec{x}_8 CLUSTER IS MORE SIMILAR TO \vec{x}_2 THAN \vec{x}_4, \vec{x}_7 IS TO \vec{x}_6
 - HORIZONTAL DISTANCES ARE IRRELEVANT

DISSIMILARITIES / DISTANCE

CLUSTERS ARE MERGED / SPLIT BASED ON DISSIMILARITIES OF ITEMS IN CLUSTERS

• For $\vec{x}_i, \vec{x}_j \in \mathbb{R}^d$ DISTANCE IS A NATURAL CHOICE OF DISSIMILARITY

• Minkowski: $\|\vec{x}_i - \vec{x}_j\|_p = \left(\sum_{k=1}^d |x_{i,k} - x_{j,k}|^p \right)^{1/p} \quad p \geq 1$

• $p=1$: MANHATTAN

• $p=2$: EUCLIDEAN

• 1-CORRELATION

$$d(\vec{x}_i, \vec{x}_j) = 1 - \rho_{ij} = 1 - \frac{s_{ij}}{s_i s_j}$$

← SAMPLE COVARIANCE

← SAMPLE CORRELATION OF OBSERVATIONS OF \vec{x}_i, \vec{x}_j

← SAMPLE STD. DEV.

* DISTANCE MEASURES COMPUTED OVER VARIABLES, NOT OBSERVATIONS

• CHOICE OF DISTANCE / DISSIMILARITY GIVES RISE TO DISTANCE MATRIX $D \in \mathbb{R}^{N \times N}$

• HOW DO WE DEFINE DISTANCES BETWEEN CLUSTERS

AGGLOMERATIVE NESTING

- 1) INIT: DISTANCE / DISSIMILARITY MATRIX D . INITIALLY HAVE N "CLUSTERS"
- 2) FIND SMALLEST DISSIMILARITY, SAY D_{IJ} IN $D = D^{(1)}$. MAKE I, J INTO A NEW CLUSTER IT
- 3) COMPUTE DISSIMILARITIES BETWEEN IT AND ALL OTHER CLUSTERS $K \neq IT$

SINGLE LINKAGE $d_{IT,K} = \min_{s \in IT, t \in K} D_{st}$

COMPLETE $d_{IT,K} = \max_{s \in IT} \max_{t \in K} D_{st}$

AVERAGE $d_{IT,K} = \frac{1}{N_{IT} N_K} \sum_{s \in IT} \sum_{t \in K} D_{st}$

- 4) FORM A NEW $(N-1) \times (N-1)$ DISTANCE MATRIX $D^{(2)}$ BY REMOVING ROWS AND COLUMNS I AND J , AND ADDING NEW ROW, COLUMN FOR IT
- 5) REPEAT 2-4 $(N-1)$ TIMES
- 6) OUTPUT: LIST OF MERGED CLUSTERS AT EACH STEP w/ VALUE (HEIGHT) OF DISSIMILARITY AT MERGE

EXAMPLE:

$$x_1 = (1, 3)^T$$

$$x_2 = (2, 4)^T$$

$$x_3 = (1, 5)^T$$

$$x_4 = (5, 5)^T$$

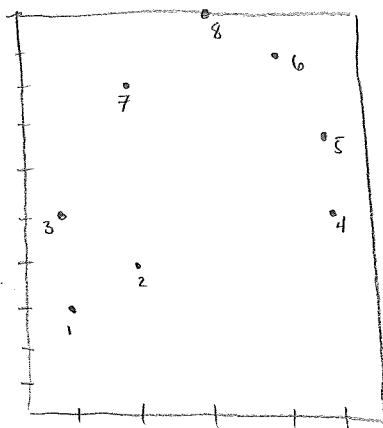
$$x_5 = (5, 7)^T$$

$$x_6 = (4, 9)^T$$

$$x_7 = (3, 8)^T$$

$$x_8 = (3, 10)^T$$

$$D^{(1)} = \begin{bmatrix} 0 & 1.4 & 2 & 4.8 & 5.7 & 6.7 & 5.1 & 7.3 \\ & 0 & 1.4 & 3.2 & 4.2 & 5.4 & 4 & 6.1 \\ & & 0 & 4 & 4.5 & 5 & 3.2 & 5.4 \\ & & & 0 & 2 & 4.1 & 4.2 & 5.4 \\ & & & & 0 & 2.2 & 3.2 & 3.6 \\ & & & & & 0 & 2.2 & 1.4 \\ & & & & & & 0 & 2.2 \\ & & & & & & & 0 \end{bmatrix}$$



AVG. LINKAGE (PRIORITIZES AVERAGE DISTANCE BETWEEN CLUSTERS)

$$D^{(2)} = \begin{matrix} & 12 & 3 & 4 & 5 & 6 & 7 \\ \begin{matrix} 12 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \end{matrix} & \begin{bmatrix} 0 & 1.7 & 3.8 & 5.0 & 6.0 & 4.6 & 6.7 \\ & 0 & & & & & \\ & & 0 & & & & \\ & & & 0 & & & \\ & & & & 0 & & \\ & & & & & 0 & \\ & & & & & & 0 \end{bmatrix} \end{matrix}$$

MERGED 12 DISTANCE = 1.4

$$D^{(5)} = \begin{matrix} & 12 & 3 & 4 & 5 & 68 & 7 \\ \begin{matrix} 12 \\ 3 \\ 4 \\ 5 \\ 68 \\ 7 \end{matrix} & \begin{bmatrix} 0 & & & & & \\ & 0 & & & & \\ & & 0 & & & \\ & & & 0 & & \\ & & & & 0 & \\ & & & & & 0 & 2.2 \\ & & & & & & 0 \end{bmatrix} \end{matrix}$$

MERGED 68 DISTANCE = 1.4

$$D^{(4)} =$$

| | 123 | 4 | 5 | 68 | 7 |
|-----|-----|-----|-----|-----|-----|
| 123 | 0 | 3.9 | 4.8 | 6.0 | 4.1 |
| 4 | | | | | |
| 5 | | | | | |
| 68 | | | | | |
| 7 | | | | | |

MERGED 12, 3 DISTANCE 1.7

$$D^{(5)} =$$

| | 123 | 45 | 68 | 7 |
|-----|-----|----|-----|-----|
| 123 | | | 6.0 | |
| 45 | | 0 | 3.8 | 3.7 |
| 68 | | | | |
| 7 | | | | |

MERGED 4, 5 DISTANCE 2

$$D^{(6)} =$$

| | 123 | 45 | 678 |
|-----|-----|-----|-----|
| 123 | 0 | 4.3 | 5.3 |
| 45 | | 0 | 3.8 |
| 678 | | | 0 |

MERGED 6, 7 DISTANCE 2.2

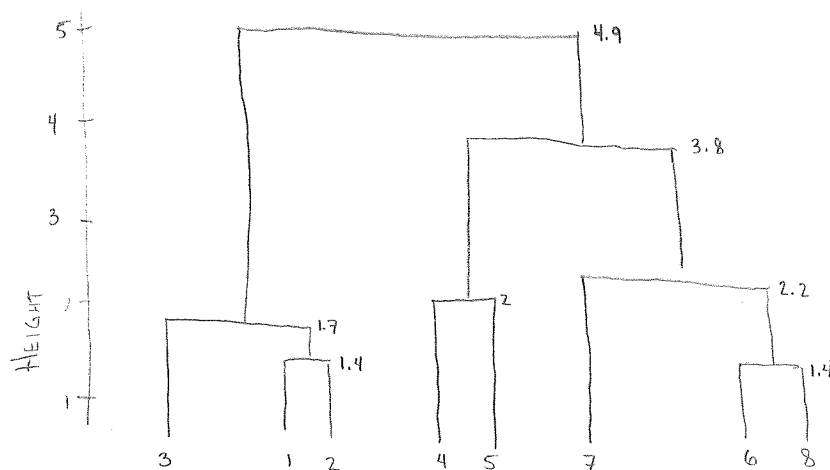
$$D^{(7)} =$$

| | 123 | 45, 678 |
|---------|-----|---------|
| 123 | 0 | 4.94 |
| 45, 678 | | 0 |

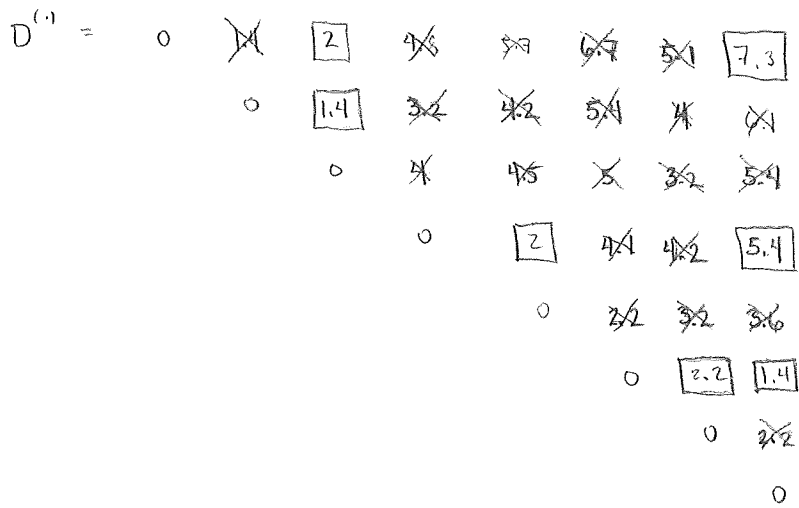
MERGED

45, 678 DISTANCE 3.8

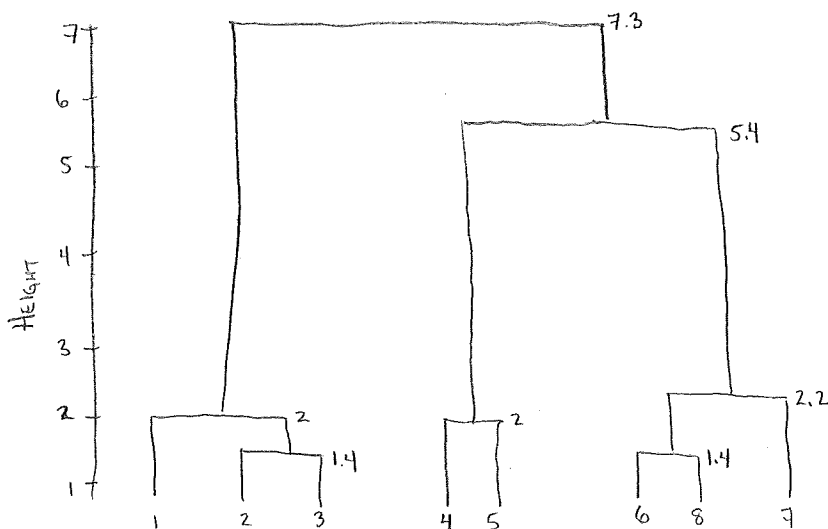
MERGE 123, 45678 DISTANCE



COMPLETE (REQUIRES ALL ELEMENTS BETWEEN CLUSTERS TO BE CLOSE



- MERGE 2, 3 AT HEIGHT 1.4
- MERGE 6, 8 AT HEIGHT 1.4
- MERGE 1, 2, 3 AT HEIGHT 2
- MERGE 4, 5 AT HEIGHT 2
- MERGE 6, 8 AT HEIGHT 2.2
- MERGE 4, 5, 6, 8 AT HEIGHT 5.4
- MERGE 1, 2, 3, 4, 5, 6, 8 AT HEIGHT 7.3

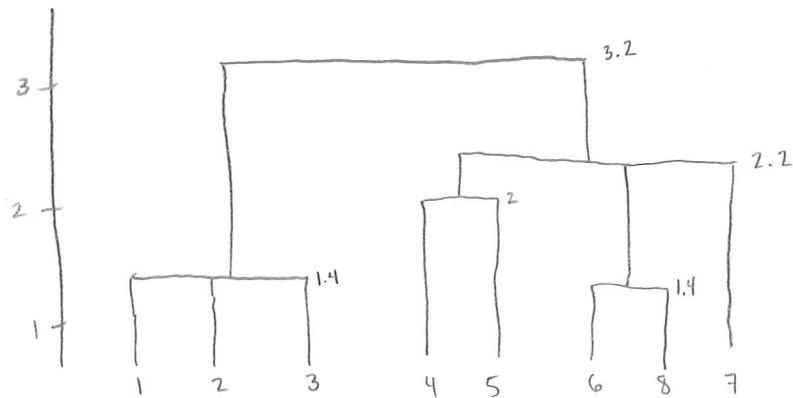


SINGLE LINKAGE (ACCEPTS LOCAL, SMALL CONNECTIONS)

(1)

| | | | | | | | | | |
|---|---|---|---------------------------------------------------|--------------|-------------------------------------------------|---------------------------------------------------|---------------------------------------------------|---------------------------------------------------|----------------|
| D | = | 0 | 1.4 | 2 | 4.8 | 5.7 | 6.7 | 5.1 | 7.3 |
| | | 0 | 1.4 | 3.2 | 4.2 | 5.4 | 4 | 6.1 | |
| | | | 0 | 4 | 4.5 | 5 | 3.2 | 5.4 | |
| | | | | 0 | 2 | 4.1 | 4.2 | 5.4 | |
| | | | | | 0 | 2.2 | 3.2 | 3.6 | |
| | | | | | | 0 | 2.2 | 1.4 | |
| | | | | | | | 0 | 2 | |
| | | | | | | | | 0 | |

- MERGE 23 AT 1.4
- MERGE 1,23 AT 1.4
- MERGE 68 AT 1.4
- MERGE 45 AT 2
- MERGE 68,7 AT 2.2
- MERGE 45,678 AT 2.2
- MERGE 123, 45678 AT 3.2



OBSERVATIONS

| <u>SINGLE</u> | <u>AVERAGE</u> | <u>COMPLETE</u> |
|---------------------------------------------------------------------------------------|--------------------------------------------------------------|------------------------------------------------------------------------------|
| JOINS IF <u>ANY</u> PAIR OF POINTS BETWEEN CLUSTERS ARE CLOSE | INTERMEDIATE | JOINS IF <u>ALL</u> PAIRS OF POINTS BETWEEN CLUSTERS ARE CLOSE |
| SHORTER HEIGHTS | | TALLER HEIGHTS |
| LONG CHAINS OF CLUSTERS JOINED BY NEARBY SINGLETONS | | MANY, SMALL COMPACT CLUSTERS |
| IND. OF SIZE OF CLUSTERS | DEPENDS ON SIZE OF CLUSTERS | IND. OF SIZE OF CLUSTERS |
| INVARIANT UNDER MONOTONE TRANSFORMATIONS OF PAIRWISE DISTANCES (STRUCTURE NOT HEIGHT) | MAY BE ALTERED BY NONLINEAR TRANSFORMATIONS TO PAIRWISE DIST | INVARIANT UNDER MONOTONE CHANGE TO PAIRWISE DISTANCES (STRUCTURE NOT HEIGHT) |

COMPARING DENDROGRAMS

LET h_{ij} BE THE HEIGHT THAT \vec{x}_i AND \vec{x}_j WERE MERGED INTO THE SAME CLUSTER

$$H = h_{ij} = \begin{bmatrix} 0 & & & \\ h_{21} & \ddots & & \\ h_{31} & & \ddots & \\ \vdots & & & 0 \end{bmatrix}$$

COPHENETIC CORRELATION IS THE (PEARSON) CORRELATION BETWEEN THE $\frac{N(N-1)}{2}$ PAIRS

$$(h_{ij}, D_{ij}) \quad 1 \leq i < j \leq N$$

1) CAN BE USED TO COMPARE DENDROGRAMS GENERATED BY DIFFERENT ALGORITHMS

2) QUALITY OF THIS MEASURE IS LIMITED BY USEFULNESS OF D_{ij}

* CONSIDER 2 NON-INTERSECTING, BUT CLOSE MANIFOLDS

D BASED ON EUCLIDEAN DISTANCE

-1 BAD 1 GOOD

CHOOSING NUMBER OF CLUSTERS ("BEST TREE CUT")

• MOJENA'S UPPER TAIL RULE

LET $\alpha_0, \dots, \alpha_{N-1}$ ARE THE HEIGHTS CORRESPONDING TO STAGES W/ $n, n-1, \dots, 1$ CLUSTERS. SELECT THE NUMBER OF CLUSTERS, j , CORRESPONDING TO THE FIRST STAGE IN THE DENDROGRAM SATISFYING

$$\alpha_{j+1} > \bar{\alpha} + k S_{\alpha}$$

$$\bar{\alpha} = \frac{1}{N} \sum_j \alpha_j$$

$$S_{\alpha}^2 = \frac{1}{N-1} \sum (\alpha_j - \bar{\alpha})^2$$

* CAN ALSO USE MOVING AVERAGES (MOJENA RULE 2)

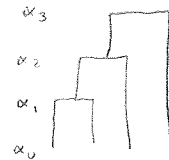
• k IS USER SPECIFIED CUTOFF

• 2.75 - 3.5 (MOJENA 1977)

• 1.25 (MILLIGAN, COOPER 1985)

• t -DISTRIBUTION (ASSUMES NORMALITY OF $\alpha_0, \dots, \alpha_{N-1}$)

• PLOT $(\alpha_j - \bar{\alpha}) / S_{\alpha}$ AND LOOK FOR DELINEATION



Ex. NCI MICROARRAY DATA

SEE CANVAS > FILES > NUMERICAL EXAMPLES > CLUSTERING