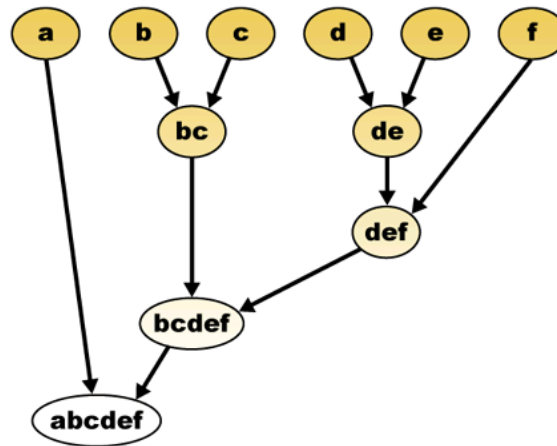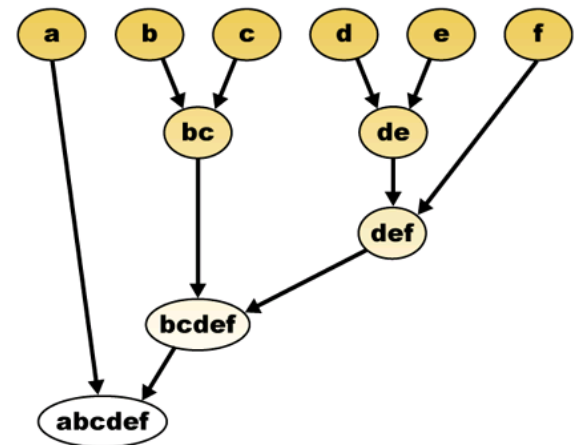# Hierarchical Clustering

# Hierarchical Clustering Algorithm

Start with *n* clusters (record = cluster)

Step 1: two closest records are merged into one cluster

At every step, pair of clusters with *smallest distance* are merged (either single record added to existing cluster, or two existing clusters are combined)

Requires a definition of **distance**

# Pairwise distance between records

Single measurement case:  Each record has 1 value.

Multiple measurement case : Each record has a multiple values.

$$d_{ij} = \text{distance between observations } i \text{ and } j$$
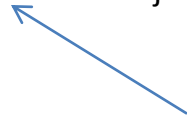
Distance Requirements:
- Non-negative ( $d_{ij} > 0$ )
- $d_{ii} = 0$
- Symmetry ($d_{ij} = d_{ji}$ )
- Triangle inequality ( $d_{ij} + d_{jk} \geq d_{ik}$ )

Distance between any pair cannot exceed the sum of distances between the other two pairs

# UG Business Programs
## Universities Clustering.xls

Data for 25 undergraduate programs at business schools in US universities in 1995.



This dataset excludes **image variables** (student satisfaction, employer satisfaction, deans' opinions)

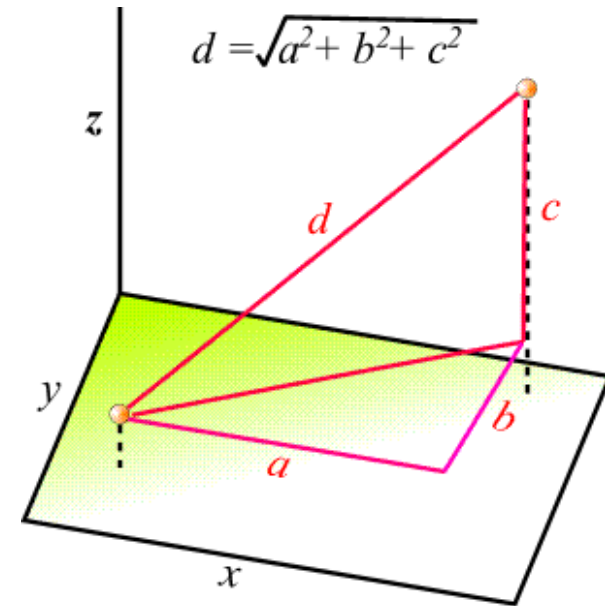| | | Student quality | | Program | | | Placement |
|---|---|---|---|---|---|---|---|
| Univ | SAT | Top10 | Accept | SFRatio | Expenses | GradRate |
|---|---|---|---|---|---|---|
| Brown | 1310 | 89 | 22 | 13 | 22,704 | 94 |
| CalTech | 1415 | 100 | 25 | 6 | 63,575 | 81 |
| CMU | 1260 | 62 | 59 | 9 | 25,026 | 72 |
| Columbia | 1310 | 76 | 24 | 12 | 31,510 | 88 |
| Cornell | 1280 | 83 | 33 | 13 | 21,864 | 90 |
| Dartmouth | 1340 | 89 | 23 | 10 | 32,162 | 95 |
| Duke | 1315 | 90 | 30 | 12 | 31,585 | 95 |
| Georgetown | 1255 | 74 | 24 | 12 | 20,126 | 92 |
| Harvard | 1400 | 91 | 14 | 11 | 39,525 | 97 |
| JohnsHopkins | 1305 | 75 | 44 | 7 | 58,691 | 87 |
| MIT | 1380 | 94 | 30 | 10 | 34,870 | 91 |
| Northwestern | 1260 | 85 | 39 | 11 | 28,052 | 89 |
| NotreDame | 1255 | 81 | 42 | 13 | 15,122 | 94 |
| PennState | 1081 | 38 | 54 | 18 | 10,185 | 80 |
| Princeton | 1375 | 91 | 14 | 8 | 30,220 | 95 |
| Purdue | 1005 | 28 | 90 | 19 | 9,066 | 69 |
| Stanford | 1360 | 90 | 20 | 12 | 36,450 | 93 |
| TexasA&M | 1075 | 49 | 67 | 25 | 8,704 | 67 |
| UCBerkeley | 1240 | 95 | 40 | 17 | 15,140 | 78 |
| UChicago | 1290 | 75 | 50 | 13 | 38,380 | 87 |
| UMichigan | 1180 | 65 | 68 | 16 | 15,470 | 85 |
| UPenn | 1285 | 80 | 36 | 11 | 27,553 | 90 |
| UVA | 1225 | 77 | 44 | 14 | 13,349 | 92 |
| UWisconsin | 1085 | 40 | 69 | 15 | 11,857 | 71 |
| Yale | 1375 | 95 | 19 | 11 | 43,514 | 96 |

Notation: $\quad x_i = (x_{i1}, x_{i2}, \ldots, x_{ip})$

$$x_j = (x_{j1}, x_{j2}, \ldots, x_{jp})$$

Caltech = (1415, 100, 25, 6, 63575, 81)

Cornell = (1280, 83, 33, 13, 21864, 90)

# Euclidean Distance

$$d = \sqrt{a^2 + b^2 + c^2}$$

$$d_{ij} = \sqrt{\left(x_{i1} - x_{j1}\right)^2 + \left(x_{i2} - x_{j2}\right)^2 + \cdots + \left(x_{ip} - x_{jp}\right)^2}$$

6-dimensional Euclidean distance between Caltech and Cornell:

$\sqrt{}$ [ $(1415-1280)^2$ + $(100-83)^2$ + $(25-33)^2$ + $(6-13)^2$ +

+ $(63575-21864)^2$ + $(81-90)^2$] = 41,711.22

6

# Standardize if multiple variables (p>1)

Euclidean distance is influenced by the **units** of the different measurements

Solution: standardize (=normalize) each variable before measuring distances

# Standardizing: Example

$$Z\_SAT = \frac{SAT - mean(SAT)}{std(SAT)}$$

| Univ | Z_SAT | Z_Top10 | Z_Accept | Z_SFRatio | Z_Expenses | Z_GradRate |
|------|-------|---------|----------|-----------|------------|------------|
| Brown | 0.401994 | 0.644235 | -0.871888 | 0.068840897 | -0.32471667 | 0.80372917 |
| CalTech | 1.370988 | 1.210256 | -0.719814 | -1.65218153 | 2.508651168 | -0.631501491 |
| CMU | -0.059432 | -0.74509 | 1.003685 | -0.91460049 | -0.16374483 | -1.625122718 |
| Columbia | 0.401994 | -0.024699 | -0.770506 | -0.17701945 | 0.285756214 | 0.141315019 |
| Cornell | 0.125139 | 0.335496 | -0.314285 | 0.068840897 | -0.38294938 | 0.362119736 |
| Dartmouth | 0.67885 | 0.644235 | -0.821197 | -0.66874014 | 0.330955887 | 0.914131529 |
| Duke | 0.448137 | 0.695691 | -0.466359 | -0.17701945 | 0.290955563 | 0.914131529 |
| Georgetown | -0.105574 | -0.127612 | -0.770506 | -0.17701945 | -0.50343562 | 0.582924453 |

Euclidean distance between standardized
    Caltech and Cornell:


    $\sqrt{[(1.371-1.125)^2 + (1.210-0.335)^2 + \ldots + (-.632-.362)^2]}$
    = 3.84

# Lots of other distance metrics

**Statistical (Mahalanobis) distance**

Uses correlation matrix

**Manhattan distance**

$$d_{ij} = \left| x_{i1} - x_{j1} \right| + \left| x_{i2} - x_{j2} \right| + \cdots + \left| x_{ip} - x_{jp} \right|$$

# Distances for Binary Data

Similarity-based metrics based on 2x2 table of counts

|   | 0 | 1 |
|---|---|---|
| **0** | a | b |
| **1** | c | d |

| | Married? | Smoker? | Manager? |
|---|---|---|---|
| **Carrie** | Y | Y | Y |
| **Sam** | N | Y | N |
| **Miranda** | N | N | Y |

| | | Miranda | |
|---|---|---|---|
| **Carrie** | | **N** | **Y** |
| | **N** | 0 | 0 |
| | **Y** | 2 | 1 |

- Binary Euclidean Distance: *(b+c)/(a+b+c+d)*
- Simple matching Coefficient:  *(a+d)/(a+b+c+d)*
- Jaquard's coefficient: *d/(b+c+d)*

For >2 categories, distance =0 only if both items have same category. Otherwise =1.

# Distances for Mixed
# (numerical + categorical) Data

**Simple:** standardize numerical variables to [0,1], then use Euclidian distance for all

Gower's General Dissimilarity Coefficient

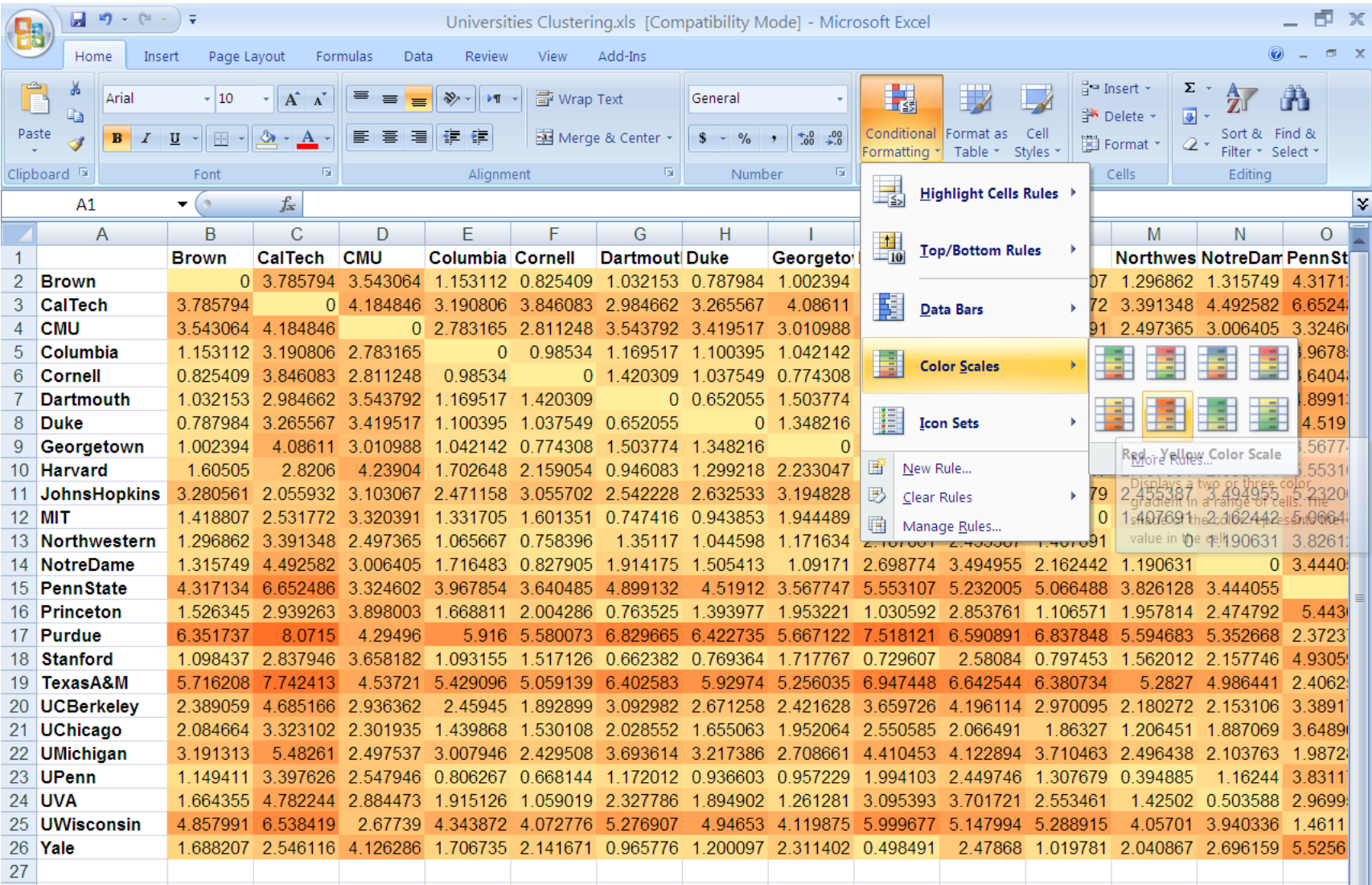$$d_{ij} = \Sigma_k\, w_{ijk}d_{ijk}\, /\, \Sigma_k\, w_{ijk}$$

$d_{ijk}$ = distance provided by $k$th variable.

$w_{ijk}$ = usually 1 or 0 depending whether or not the comparison is valid for the $k$th variable.

More on Gower's measure for mixed data:
www.soziologie.wiso.unierlangen.de/koeln/script/chap6.pdf

**Worksheet: Euclid Distance Matrix**

# Distance Matrix



Feed matrix into hierarchical clustering algorithm
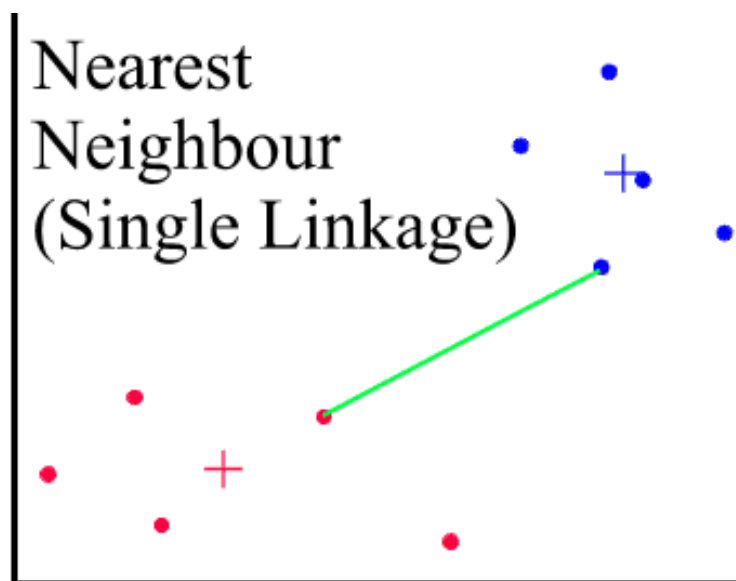
# Once Again: The Hierarchical Clustering Algorithm

✓ Start with *n* clusters (record= cluster)

✓ Step 1: two closest records are merged into one cluster

At every step, pair of clusters with *smallest distance* are merged (either single record added to existing cluster, or two existing clusters are combined)
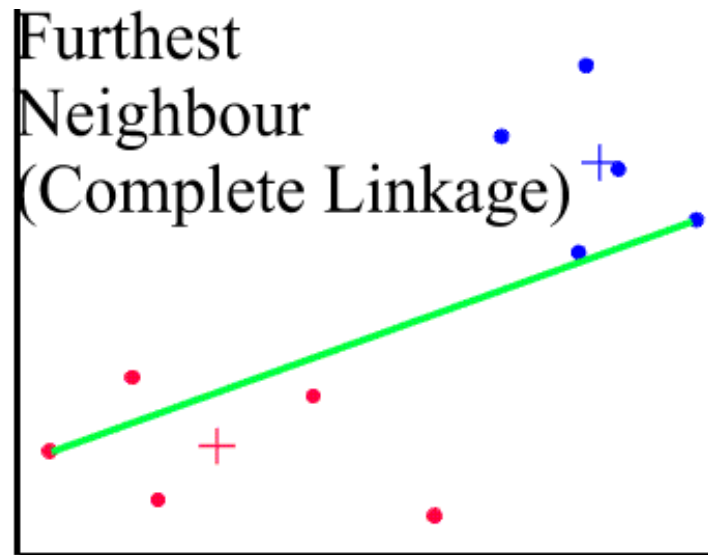
How to measure distances **between clusters**?

# Distances Between Clusters:
# 'single linkage' ('nearest neighbor')

Distance between 2 clusters = **minimum distance** between members of the two clusters
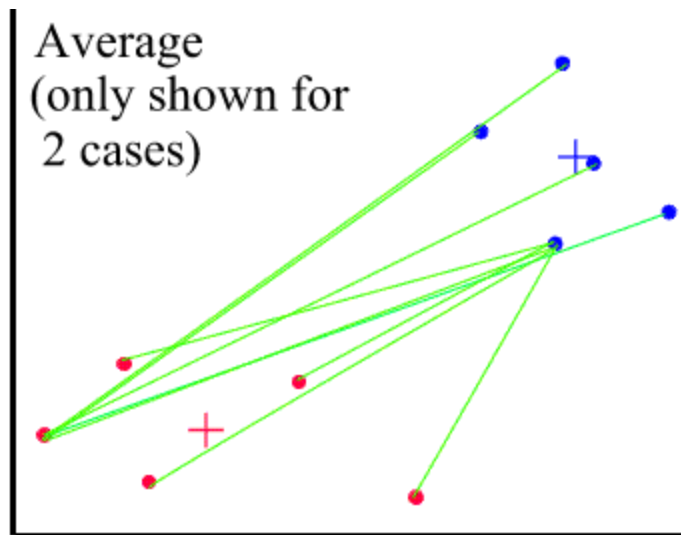


Nearest
Neighbour
(Single Linkage)

# Distances Between Clusters: 'complete linkage' ('farthest neighbor')

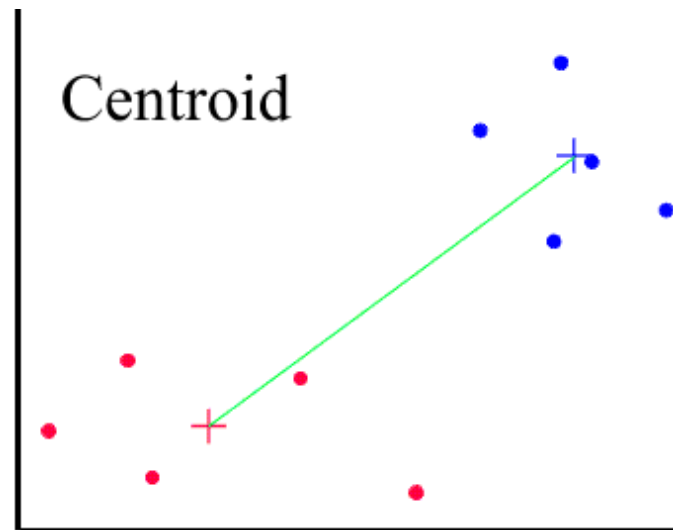Distance between 2 clusters = **greatest distance** between members of the two clusters



Furthest
Neighbour
(Complete Linkage)

# Distances Between Clusters: 'average linkage'

Distance between 2 clusters = **average** of all distances between members of the two clusters



Average (only shown for 2 cases)

# Distances Between Clusters: 'centroid linkage'

Distance between 2 clusters =
distance between their **centroids** (centers)

# And Again:
# The Hierarchical Clustering Algorithm

✓ Start with *n* clusters (record = cluster)

✓ Step 1: two closest records are merged into one cluster

At every step, pair of clusters with *smallest distance* are merged.
**At this point the distance matrix is re-computed:**

- **Two rows+columns are merged into single row+column**

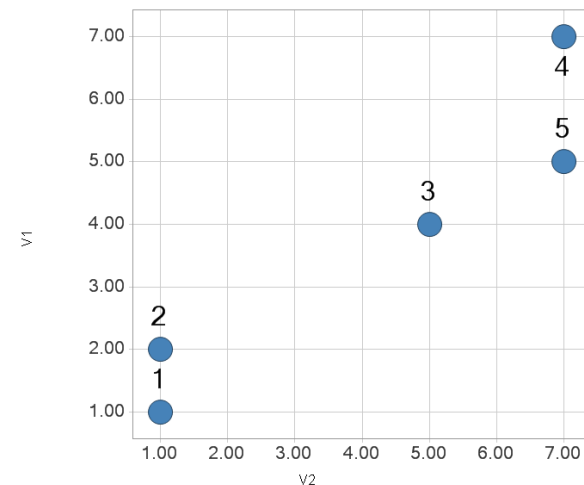- **Distances to the newly merged cluster are recalculated**

Repeat the last step until a single cluster is formed

# The clustering process: example

## Two variables, n=5 items:

| item | v1 | v2 |
|------|----|----|
| 1 | 1 | 1 |
| 2 | 2 | 1 |
| 3 | 4 | 5 |
| 4 | 7 | 7 |
| 5 | 5 | 7 |



## Euclidean distance matrix

|   | 1 | 2 | 3 | 4 | 5 |
|---|-----|-----|-----|-----|-----|
| 1 | 0.0 | | | | |
| 2 | 1.0 | 0.0 | | | |
| 3 | 5.0 | 4.5 | 0.0 | | |
| 4 | 8.5 | 7.8 | 3.6 | 0.0 | |
| 5 | 7.2 | 6.7 | 2.2 | 2.0 | 0.0 |

# What happens next?

- Merge 1&2 into cluster A
- Use single linkage to compute distances from cluster A:

|   | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | 0.0 | | | | |
| 2 | 1.0 | 0.0 | | | |
| 3 | 5.0 | 4.5 | 0.0 | | |
| 4 | 8.5 | 7.8 | 3.6 | 0.0 | |
| 5 | 7.2 | 6.7 | 2.2 | 2.0 | 0.0 |

→

|   | A | 3 | 4 | 5 |
|---|---|---|---|---|
| A | 0.0 | | | |
| 3 | 4.5 | 0.0 | | |
| 4 | 7.8 | 3.6 | 0.0 | |
| 5 | 6.7 | 2.2 | 2.0 | 0.0 |

# What happens next?

Merge 4&5 (cluster B)

Merge 3 & B

|   | A | 3 | B |
|---|---|---|---|
| A | 0.0 | | |
| 3 | 4.5 | 0.0 | |
| B | 6.7 | 2.2 | 0.0 |

→

|   | A | B |
|---|---|---|
| A | 0.0 | |
| B | 4.5 | 0.0 |

# Finally: Summarize process in a **Dendrogram**