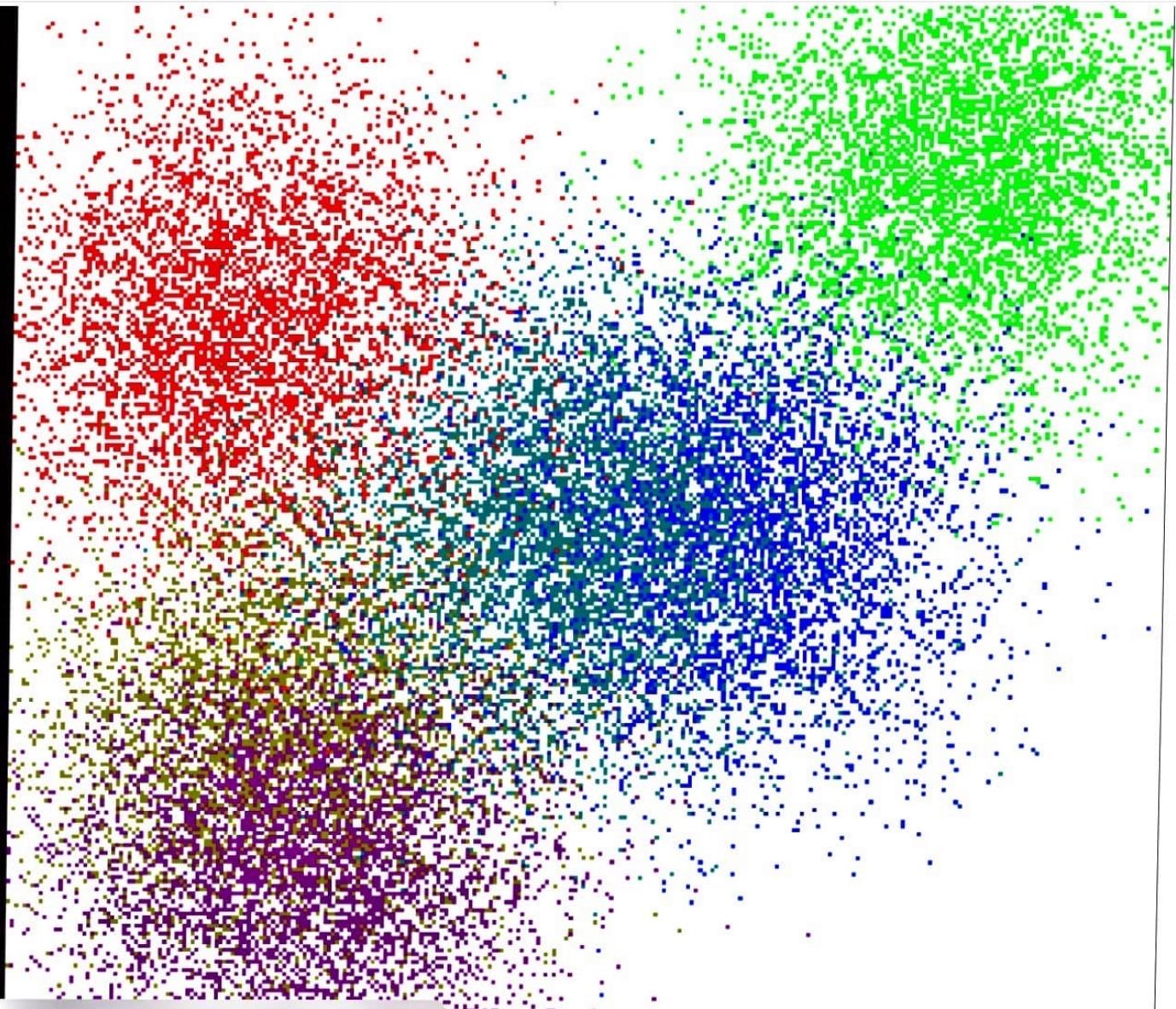


# Data Mining : Clustering





# What is clustering?

Cluster Analysis (“data segmentation”) is an exploratory method for identifying homogenous groups (“clusters”) of records

Similar records should belong to the same cluster

Dissimilar records should belong to different clusters

# So What and Why Clustering?

**Clustering = Grouping “Similar” things together!**

- **Understand/Discover** structure in data
- **Summarize** data points by their “**Cluster center**”

# Nature is “naturally” organized

五

## **Periodic Table of the Elements**

## Lanthanide Series

Actinide  
Series



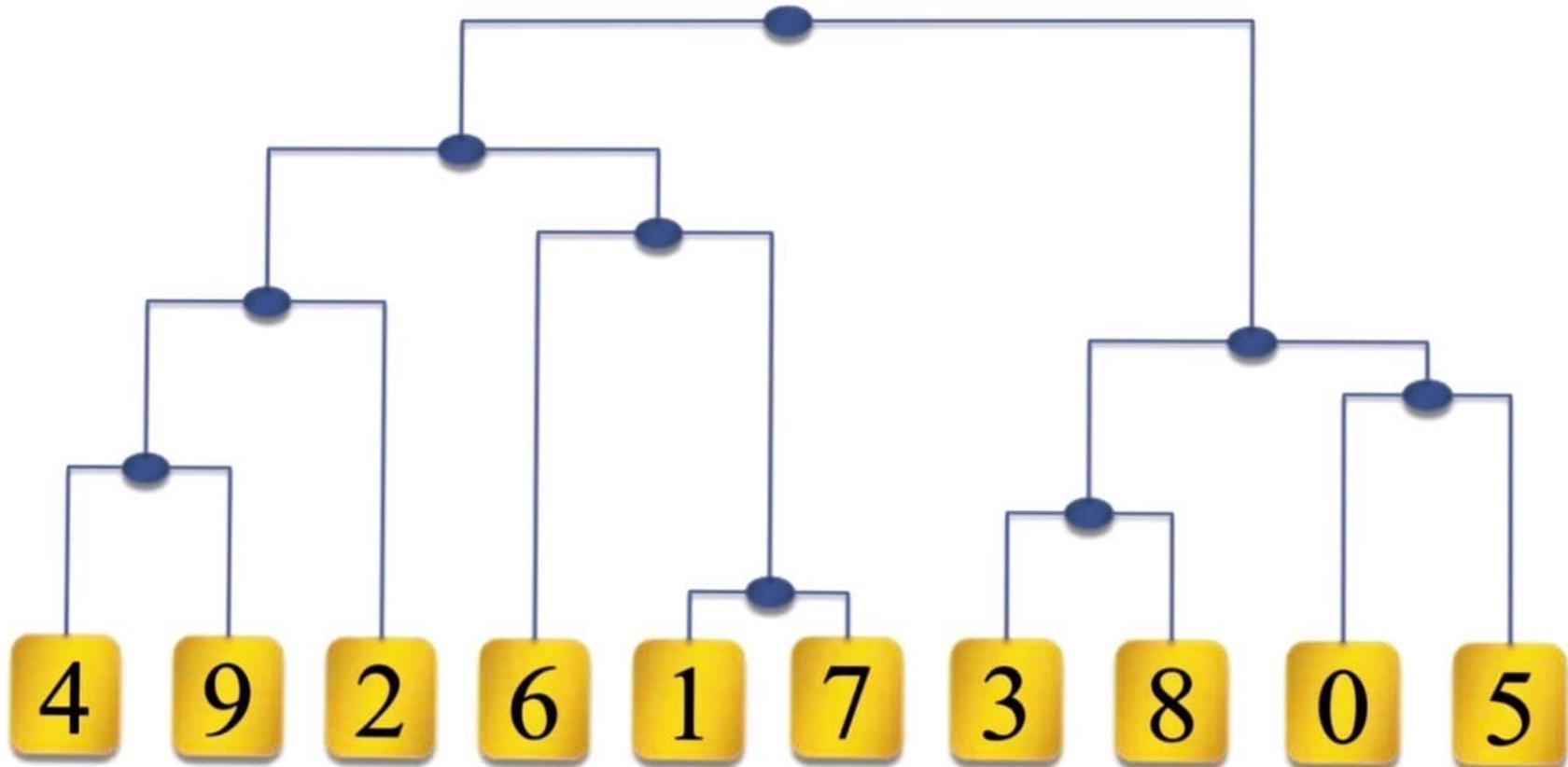
# What Business questions can be answered?

- Types of **pages** are there on the **Web**?
- Types of **customers** are there in my **market**?
- Types of **people** are there on a **Social network**?
- Types of **E-mails** in my **Inbox**?
- Types of **Genes** the **human genome** has?

# Hierarchical Clustering

Agglomerative Clustering

# Hierarchical Clustering :Bottom-Up (Agglomerative)



- **SIMILARITY** between two **DATA POINTS?**  $\text{Sim}(0, 5)$
- **SIMILARITY** between **CLUSTERS?**  $\text{Sim}(\{3, 8\}, \{0, 5\})$

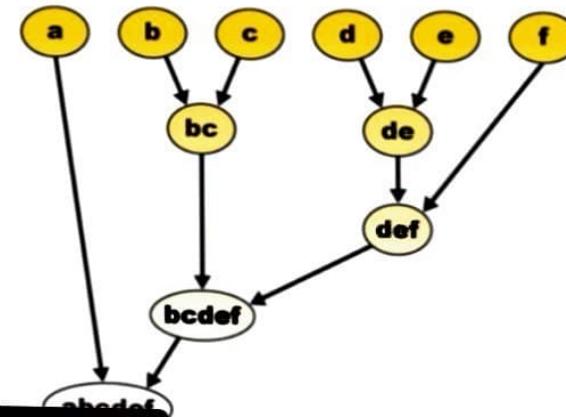
# Hierarchical Clustering Algorithm

Start with  $n$  clusters (record = cluster)

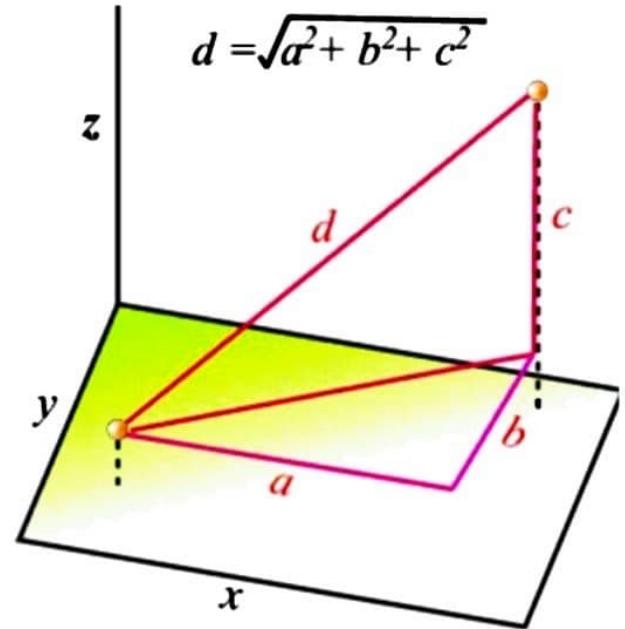
Step 1: two closest records are merged into one cluster

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**



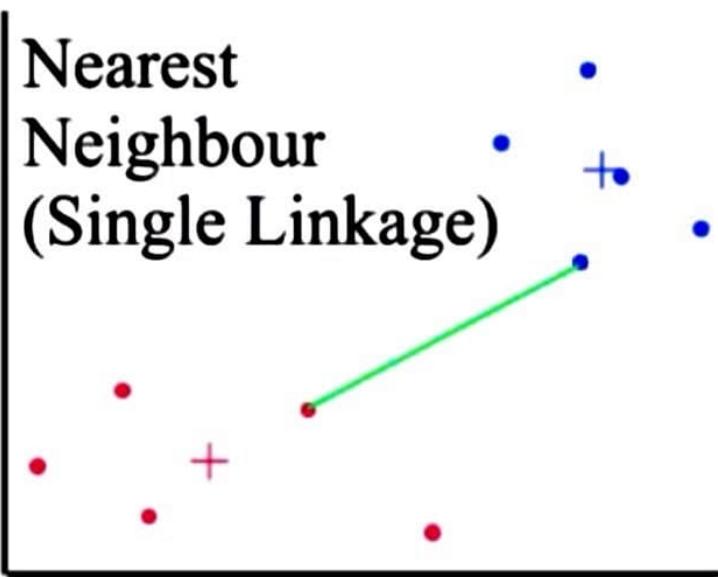
# Euclidean Distance



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

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

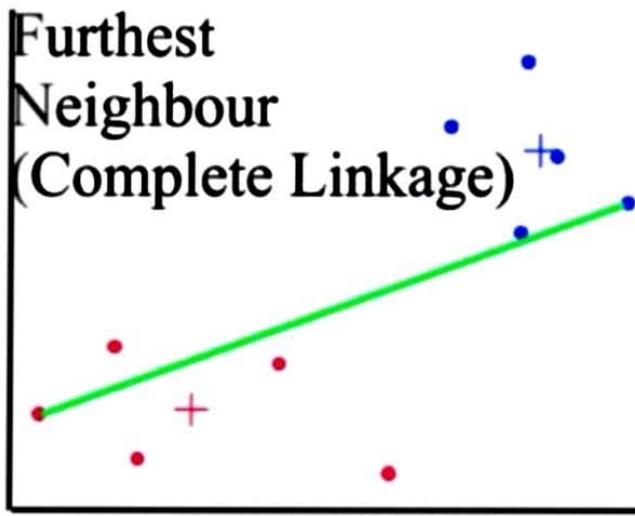
Distance between 2 clusters = **minimum distance** between members of the two clusters



$$\Delta(C_\alpha, C_\beta) = \min_{x \in C_\alpha, y \in C_\beta} \{ \Delta(x, y) \}$$

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

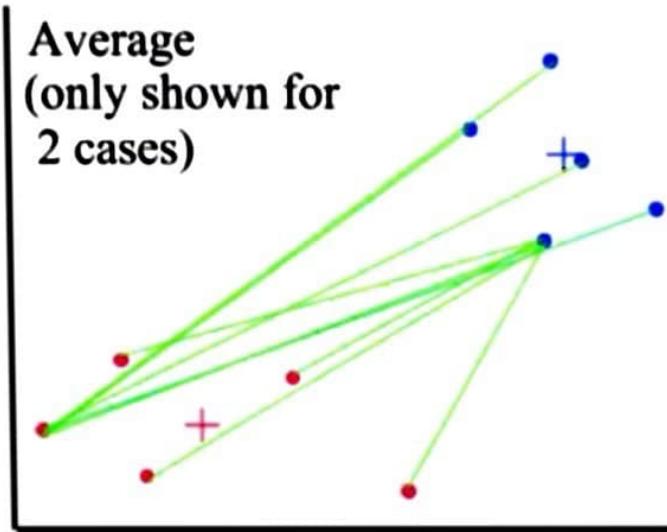
Distance between 2 clusters = **greatest distance** between members of the two clusters



$$\Delta(C_\alpha, C_\beta) = \max_{x \in C_\alpha, y \in C_\beta} \{ \Delta(x, y) \}$$

# Distances Between Clusters: 'average linkage'

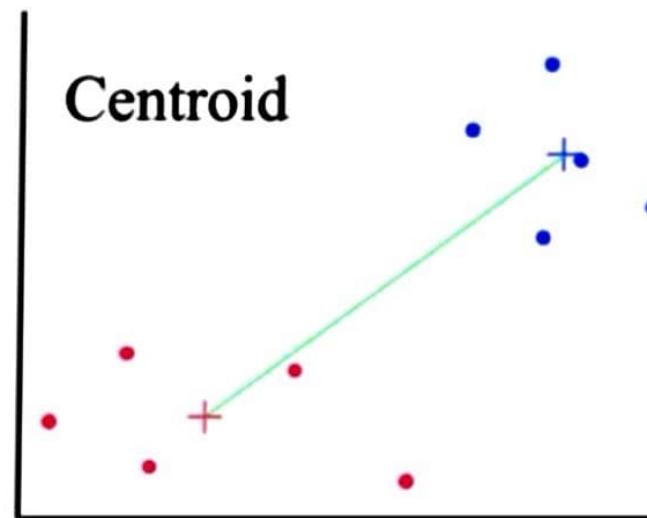
Distance between 2 clusters = **average** of all distances between members of the two clusters



$$\Delta(C_\alpha, C_\beta) = \frac{1}{|C_\alpha||C_\beta|} \sum_{x \in C_\alpha} \sum_{y \in C_\beta} \Delta(x, y)$$

# Distances Between Clusters: 'centroid linkage'

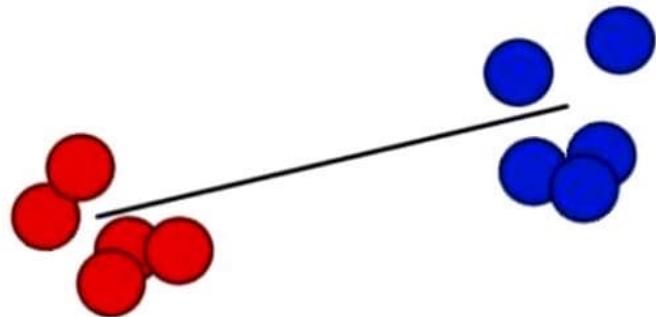
Distance between 2 clusters =  
distance between their **centroids** (centers)



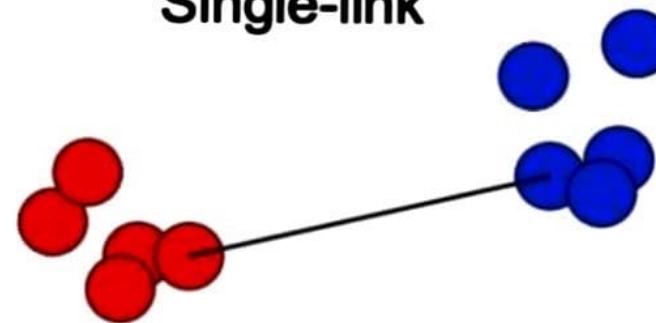
$$\Delta(C_\alpha, C_\beta) = \|\mathbf{m}_\alpha - \mathbf{m}_\beta\|$$

# All linkage methods

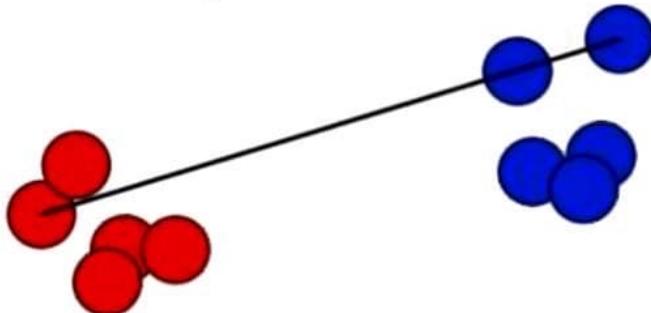
Distance between centroids



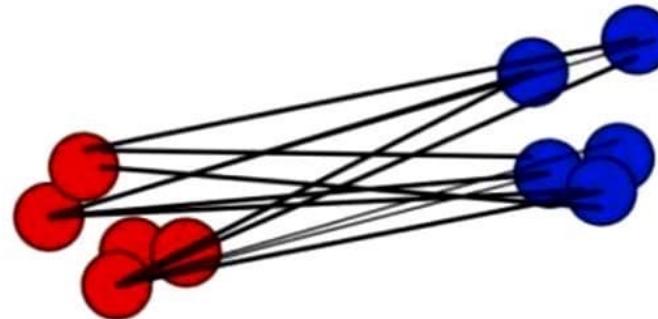
Single-link



Complete-link

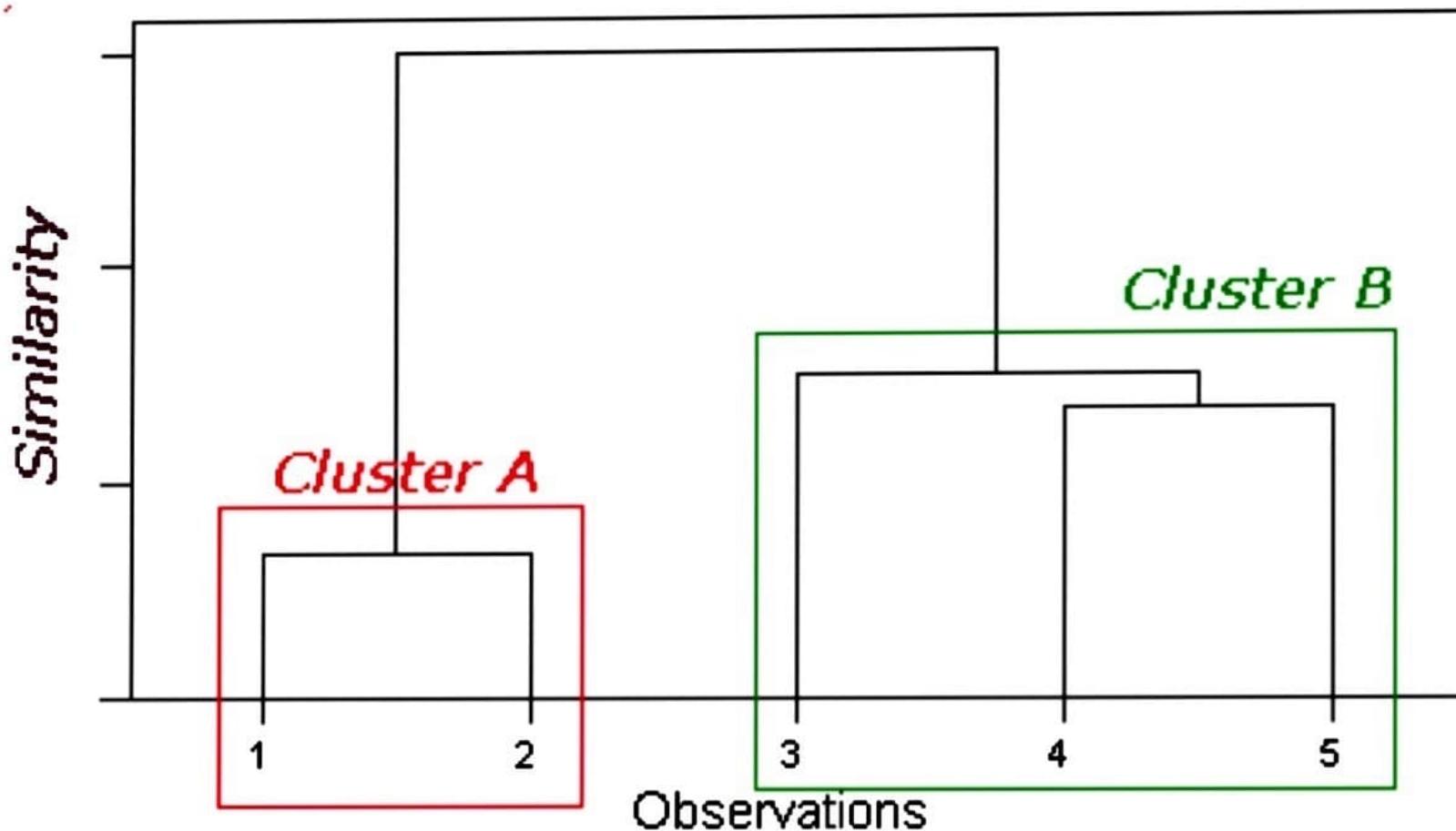


Mean-link



Let us work on simple example

# Finally: Summarize process in a Dendrogram

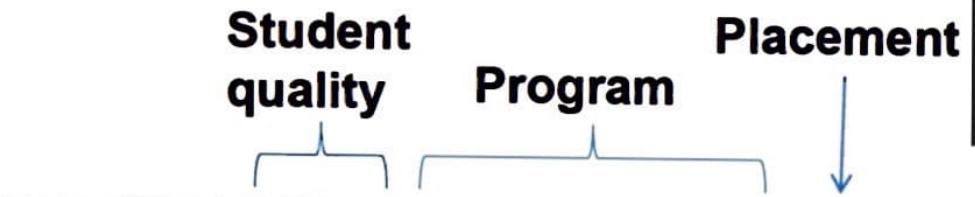


# University Data set

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



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



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
Johns Hopkins	1305	75	44	7	58,691	87
MIT	1380	94	30	10	34,870	91
Northwestern	1260	85	39	11	28,052	89
Notre Dame	1255	81	42	13	15,122	94
Penn State	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
Texas A&M	1075	49	67	25	8,704	67
UC Berkeley	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

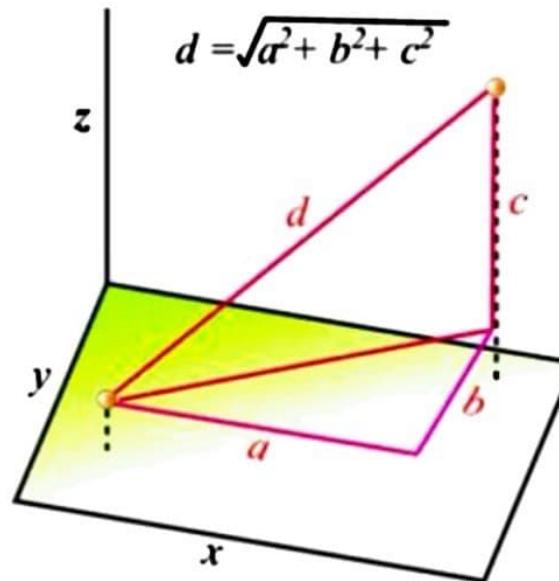
**Notation:**  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$

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

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

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

# Euclidean Distance



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

6-dimensional Euclidean distance between Caltech  
and Cornell:

$$\begin{aligned} & \sqrt{[(1415-1280)^2 + (100-83)^2 + (25-33)^2 + (6-13)^2 + \\ & + (63575-21864)^2 + (81-90)^2]} = 41,711.22 \end{aligned}$$

# 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

$$Z_{-SAT} = \frac{SAT - \text{mean}(SAT)}{\text{std}(SAT)}$$

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:

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

# From Dendrograms to Clusters

- After dendrogram is obtained, **cut** it to create clusters. **How?**
- Examine *distance levels*
  - Cutpoint determines # clusters
  - Obtain statistics on resulting clusters

