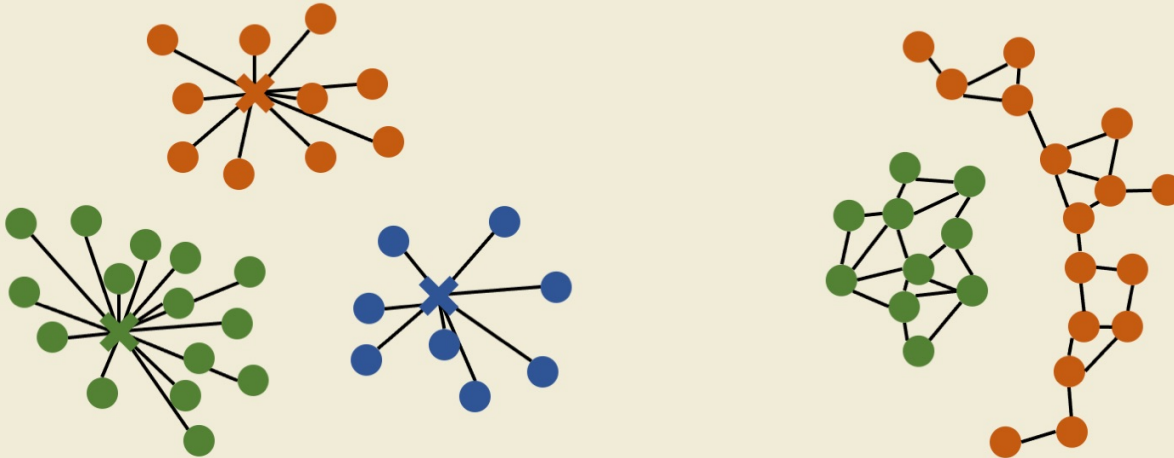# Clustering
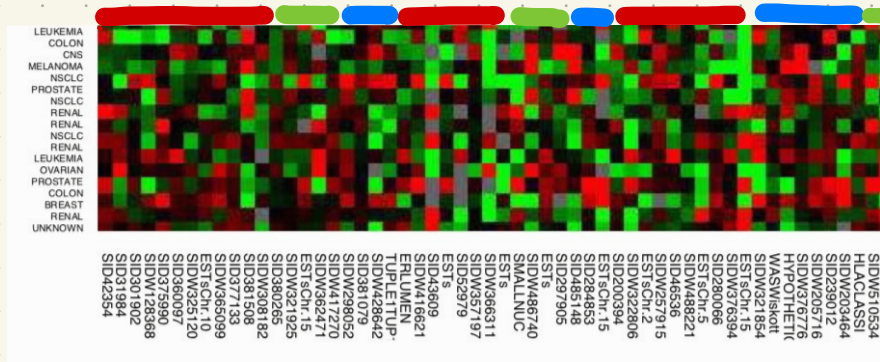
# Clustering

Split data into groups of similar points.



Example: Group genes displaying simalarities.

# K-means clustering

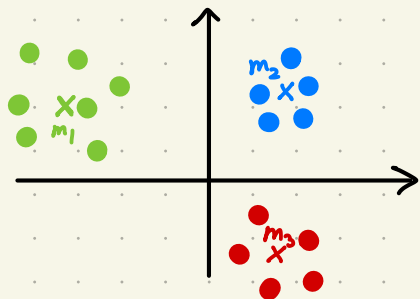Given $x_1, \ldots x_n \in \mathbb{R}^p$ assign them to $K$ clusters.

$$x_1 \quad x_2 \quad \ldots \ldots \quad x_{n-1} \quad x_n$$
$$\downarrow \qquad \downarrow \qquad \qquad \qquad \downarrow \qquad \downarrow$$
$$3 \qquad 2 \qquad \qquad \qquad 2 \qquad K$$

Assume that points in the same cluster are close is Euclidean distance.

Goal: find $m_1 \ldots m_K \in \mathbb{R}^p$ cluster centroids

and $z_1, \ldots z_n \in \{1 \ldots K\}$ cluster assignment

mimimizing the distortion function

$$J\left(\{z_i\}_{i=1}^n, \{m_k\}_{k=1}^K\right) = \sum_{i=1}^n \|x_i - m_{z_i}\|^2$$

$\underbrace{\phantom{\{z_i\}}}_{z} \qquad \underbrace{\phantom{\{m_k\}}}_{m}$

## Algorithm

Input: points $x_1, \ldots x_n \in \mathbb{R}^p$, number of clusters $k$

arbitrary centroids $m_1, \ldots m_k \in \mathbb{R}^p$
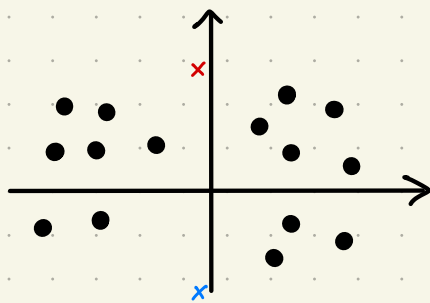
Step 1: Assign each point to the closest centroid

$$z_i = \underset{k=1,\ldots k}{\arg\min} \|x_i - m_k\|^2 \qquad \text{for} \quad i = 1, \ldots n$$

Step 2: Update centroids with cluster means

$$m_k = \frac{\sum_{i=1}^{n} I(z_i = k) \, x_i}{\sum_{i=1}^{n} I(z_i = k)} \qquad \text{for} \quad k = 1, \ldots k$$
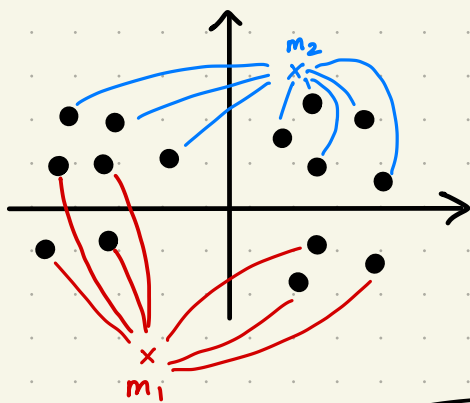
Output: cluster assignments $z_1, \ldots z_n$
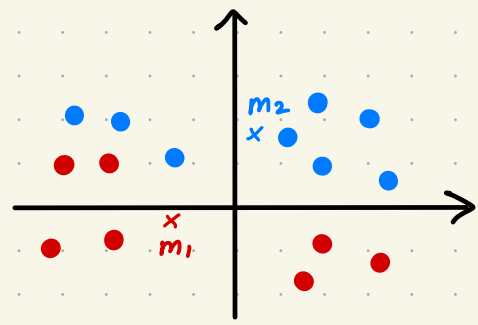
cluster centers $m_1, \ldots m_k$.
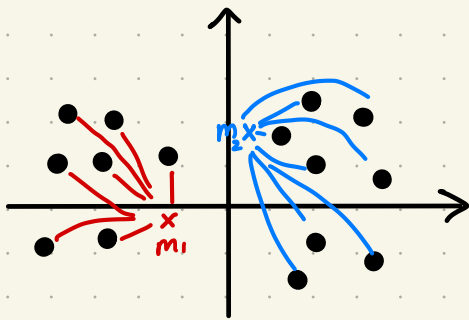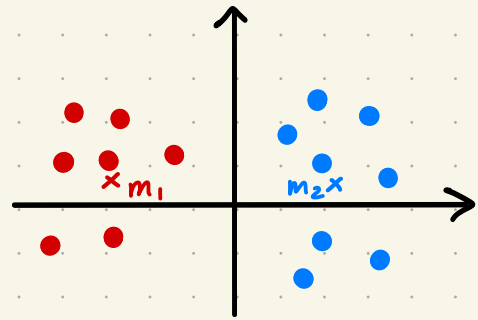
Input
(k=2)

Iteration 1

Step 1

$m_2$

Step 2

$m_2$
×

×
$m_1$

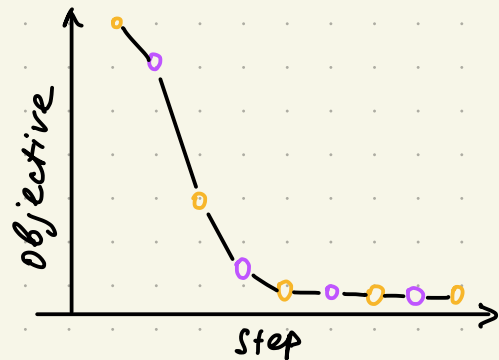Iteration 2

Step 1

$m_2$×

×
$m_1$

Step 2

×$m_1$
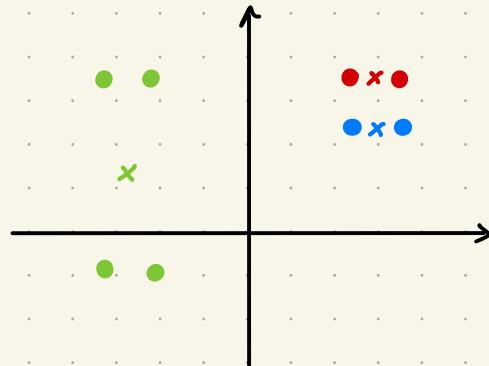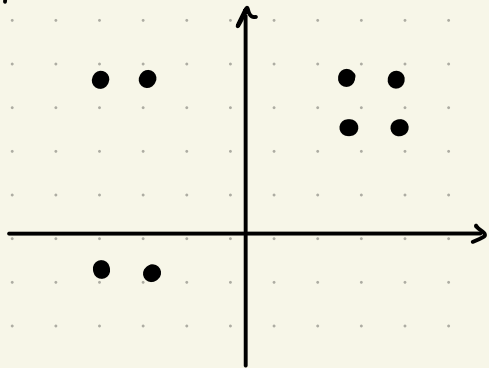
$m_2$×

You can show that
- $Step\,1$ decreases the objective $J(z,m)$
- $Step\,2$ decreases the objective $J(z,m)$

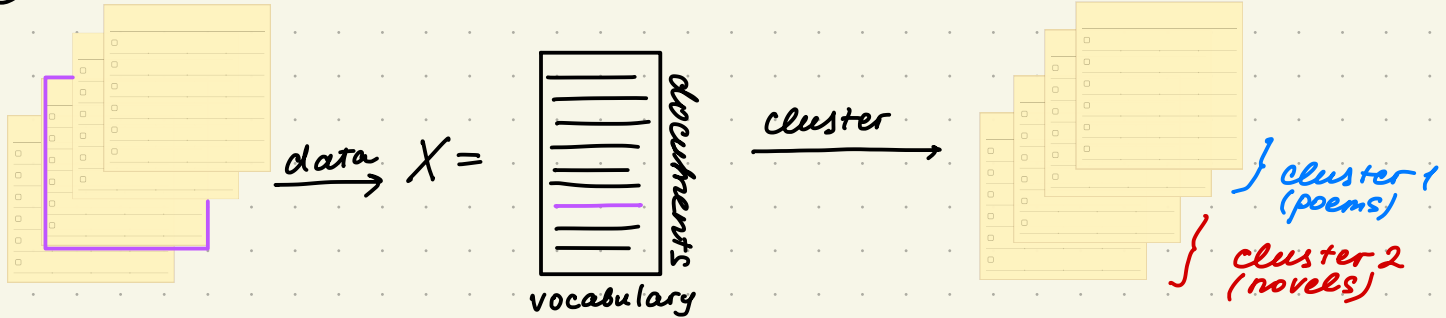So the objective $J(z,m)$
will converge (HW3)


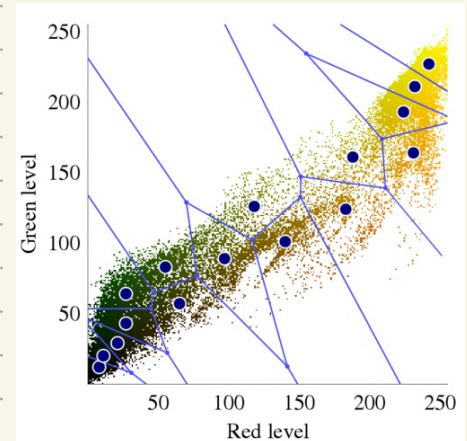
It is possible to get stuck at local minima

Input
(k = 3)

# Applications

① Documentation classification



$\xrightarrow{\text{data}} X = $ 

documents / vocabulary

$\xrightarrow{\text{cluster}}$

} cluster 1 (poems)

} cluster 2 (novels)

② Image compression (vector quantization)
reduce color pallette to K colors



$\xrightarrow{\text{data}} X = $ 

pixels

R G B

$\xrightarrow{\text{clustering}}$

## Limitations:

- K-means is sensitive to outliers
- Needs random restarts.

## Modifications:

① Replace distances $\|x-y\|^2$ with $d(x,y)$

$$J(z,m) = \sum_{i=1}^{n} d(x_i, m_{z_i})$$

+ May solve the outliers problem

− Step 2 may become very expensive

② k-medoids: pick centroids among the data

i.e. $m_j \in \{x_1, \ldots x_n\}$.

+ : we don't need to recompute $\|x_i - m_k\|^2$

instead we just store $\|x_i - x_j\|^2$.

− Step 2 is more expensive now.

③ ==k-means++== (Arthur and Vassilvitskii, 2008):

improves initialization

<u>Init</u>: pick $m_1$ at random from $x_1 \dots x_n$

for $j = 2 \dots k$

<u>Step 1</u>: compute $D^2(x_i) = \| x_i - m_{k-1} \|^2$

<u>Step 2</u>: Choose $m_k = x_i$ with probability

proportional to $D^2(x_i)$

After selecting the centroids $m_1 \dots m_K$

run k-means.

+ : more accurate and faster than k-means

(4) If we denote $w_{ik} = I(z_i = k)$ in k-means then:

$$J(z, m) = \sum_{i=1}^{n} \| x_i - m_{z_i} \|_2^2$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{k} w_{ik} \| x_i - m_k \|^2 = J(w, m)$$

Step 2 : Update centroids with cluster means

$$m_k = \frac{\sum_{i=1}^{n} I(z_i = k) x_i}{\sum_{i=1}^{n} I(z_i = k)} = \frac{\sum_{i=1}^{n} w_{ik} x_i}{\sum_{i=1}^{n} w_{ik}} \qquad \text{for} \quad k = 1 \dots k$$

Soft-clustering: uses $w_{ij}$ "probability weights" to assign observations to clusters.

minimize $J(w, m) = \sum_{i=1}^{n} \sum_{k=1}^{k} w_{ik} \| x_i - m_k \|^2$

$$0 \leq w_{ik} \leq 1 \quad \forall \, i, k \, ; \quad \sum_{i=1}^{n} w_{ik} > 0 \, ; \quad \sum_{k=1}^{k} w_{ik} = 1$$

# Choosing K

Determining $K$ is a hard problem!

Denote by $C_1 \ldots C_k \subseteq \{1, \ldots, n\}$ the cluster partitioning. Thus $C_k \cap C_{k'} = \emptyset$ for $k \neq k'$ and $C_1 \cup \ldots \cup C_k = \{1 \ldots n\}$. Denote $n_k = |C_k|$

## Example

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$$
$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$
$$1 \quad \ 2 \quad \ 1 \quad \ 1 \quad \ 2$$

$C_1 = \{1, 3, 4\} \qquad C_2 = \{2, 5\} \qquad n_1 = 3 \qquad n_2 = 2$

Denote by $\bar{x}_k = \frac{1}{n_k} \sum_{i \in C_k} x_i$ the cluster centers.

and by $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x$ the mean value across all data points.

Within-cluster scatter (variation) is

$$W(k) = \sum_{k=1}^{k} \sum_{i \in C_k} \| x_i - \bar{x}_k \|^2 \quad - \text{ measures how}$$
tight clusters are
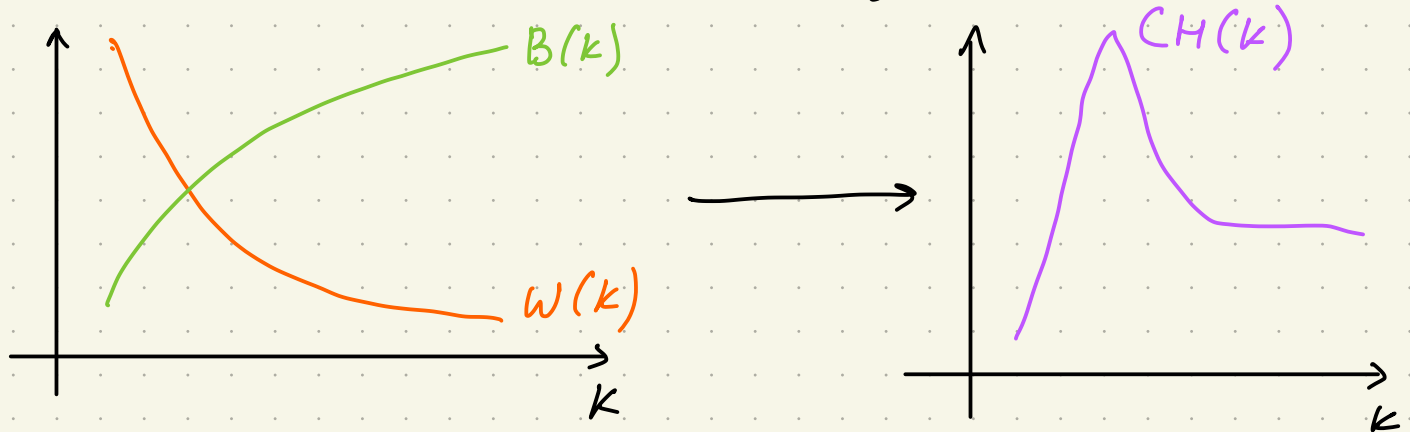
Between cluster scatter (variation) is

$$B(k) = \sum_{k=1}^{k} n_k \| \bar{x}_k - \bar{x} \|^2 \quad - \text{ measures how}$$
spread apart clusters are

## CH score (Calinski and Harabasz, 1974)

Low values of $W(k)$ are good.

High values of $B(k)$ are good.



- We can use elbow detection again.
- Alternatively, $$CH(k) = \frac{B(k)/(k-1)}{W(k)/(n-k)}$$

Pick $k$ that has higest CH score.

# Gap statistic (Tibshirani et al., 2001)

- CH cannot be computed for $k=1$.
- Gap statistic uses $W(k)$ and compares it to $W_{unif}(k)$ that is within cluster scatter computed for (simulated) uniform data.

$$Gap(k) = \underbrace{\log W_{unif}(k)}_{\substack{\text{averaged across} \\ \text{several simulations}}} - \log W(k)$$

- Larger Gap means more deviation from the uniform distribution.

We also compute $S(k)$, the standard error of
$\log W_{unif}(k)$ over simulations.

- Find minimum $k$ such that:

$$Gap(k) \geq Gap(k+1) - S(k+1)$$

adjustment for the
variation introduced by
the simulation



$Gap(k)$

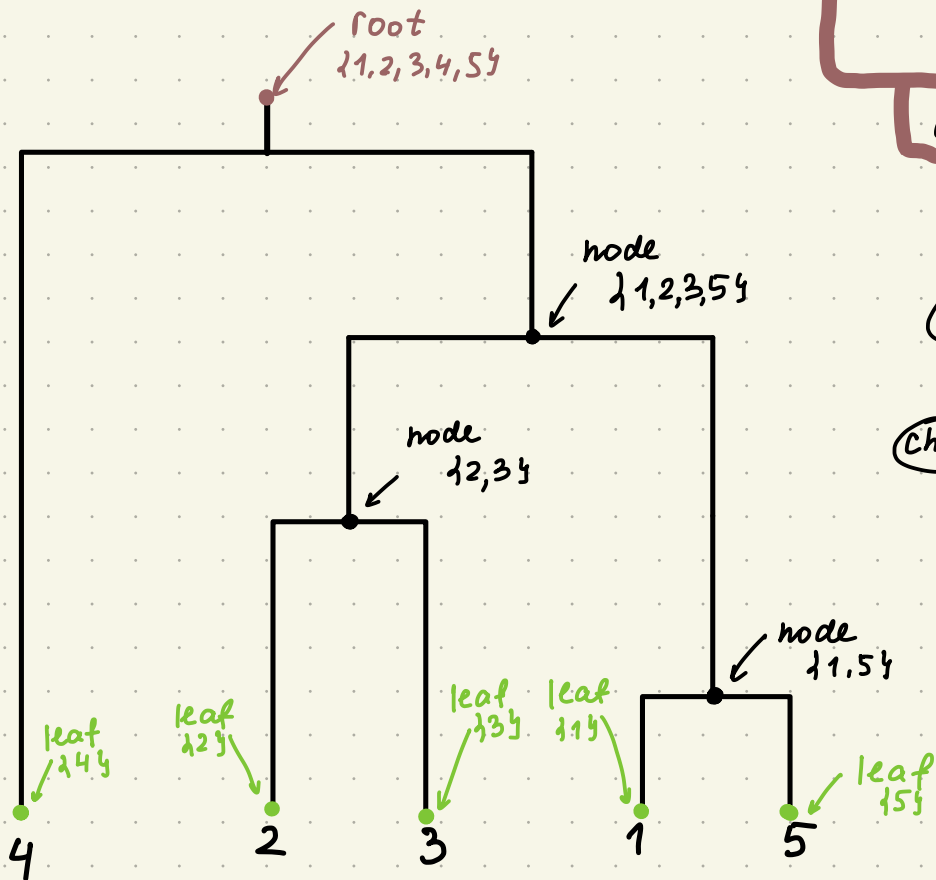$Gap(k) - S(k)$

# Hierarchical clustering

Two types:

k=1 cluster
(all points)

divisive
(top-down)

agglomerative
(bottom-up)

K = n clusters
(1 point per cluster)

# Clustering dendrogram

root
{1,2,3,4,5}

node
{1,2,3,5}

node
{2,3}

node
{1,5}

leaf
{4}

leaf
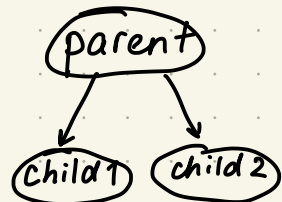{2}

leaf
{3}
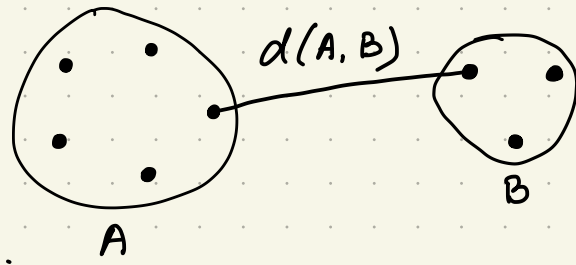
leaf
{1}

leaf
{5}

4   2   3   1   5

a tree :)

parent

child 1   child 2

# Linkages

Linkage defines when to merge/split two clusters.

Given $x_1, \ldots x_n \in \mathbb{R}^p$ define $d_{ij}$ the ==dissimilarity== ==between $x_i$ and $x_j$==, e.g. $d_{ij} = \| x_i - x_j \|$.

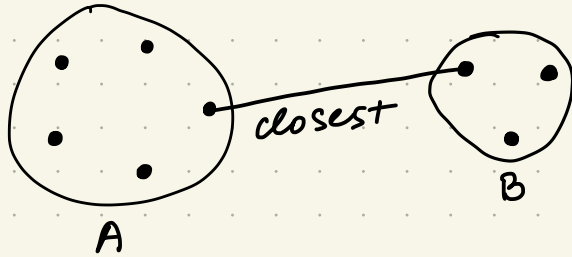Given two groups $A$ and $B \subseteq \{1 \ldots n\}$ denote by $n_A = |A|$ and $n_B = |B|$ the size of $A$ and $B$.

Definde ==dissimilarity== $d(A,B)$ ==between $A$ and $B$==,
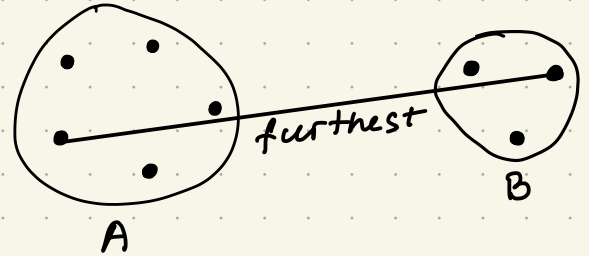
e.g. $\displaystyle d(A,B) = \min_{i \in A, j \in B} d_{ij}$.

## Agglomerative clustering:

· Start with one point per group.
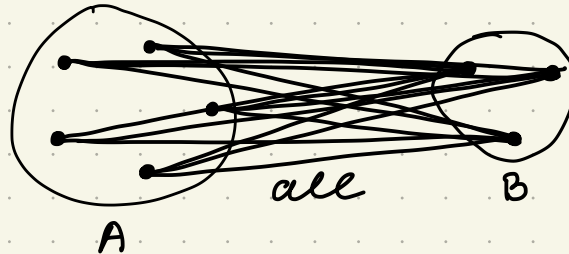
· Merge $A$ and $B$ with the smallest $d(A,B)$.

$$d_{single}(A, B) = \min_{i \in A, j \in B} d_{ij}.$$

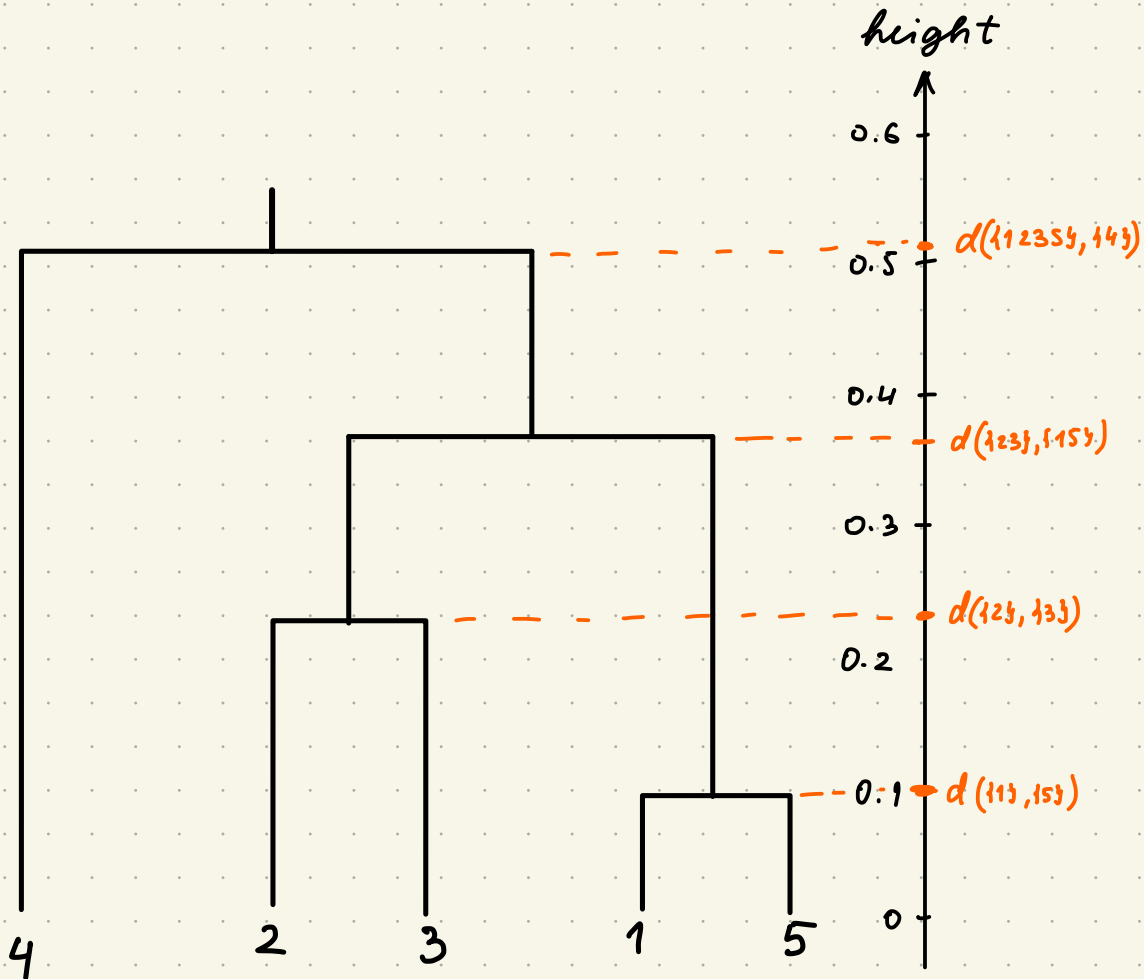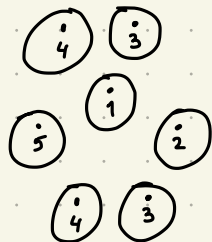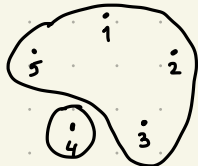$$d_{complete}(A, B) = \max_{i \in A, j \in B} d_{ij}.$$



closest

A    B



furthest

A    B

$$d_{average}(A, B) = \frac{1}{n_A n_B} \sum_{i \in A} \sum_{j \in B} d_{ij}.$$



all

A    B

- Single linkage suffers from chaining
- Complete linkage suffers from crowding

# Given a dendrogram you can decide on the dissimilarity cutoff. ($h = 0.3$)