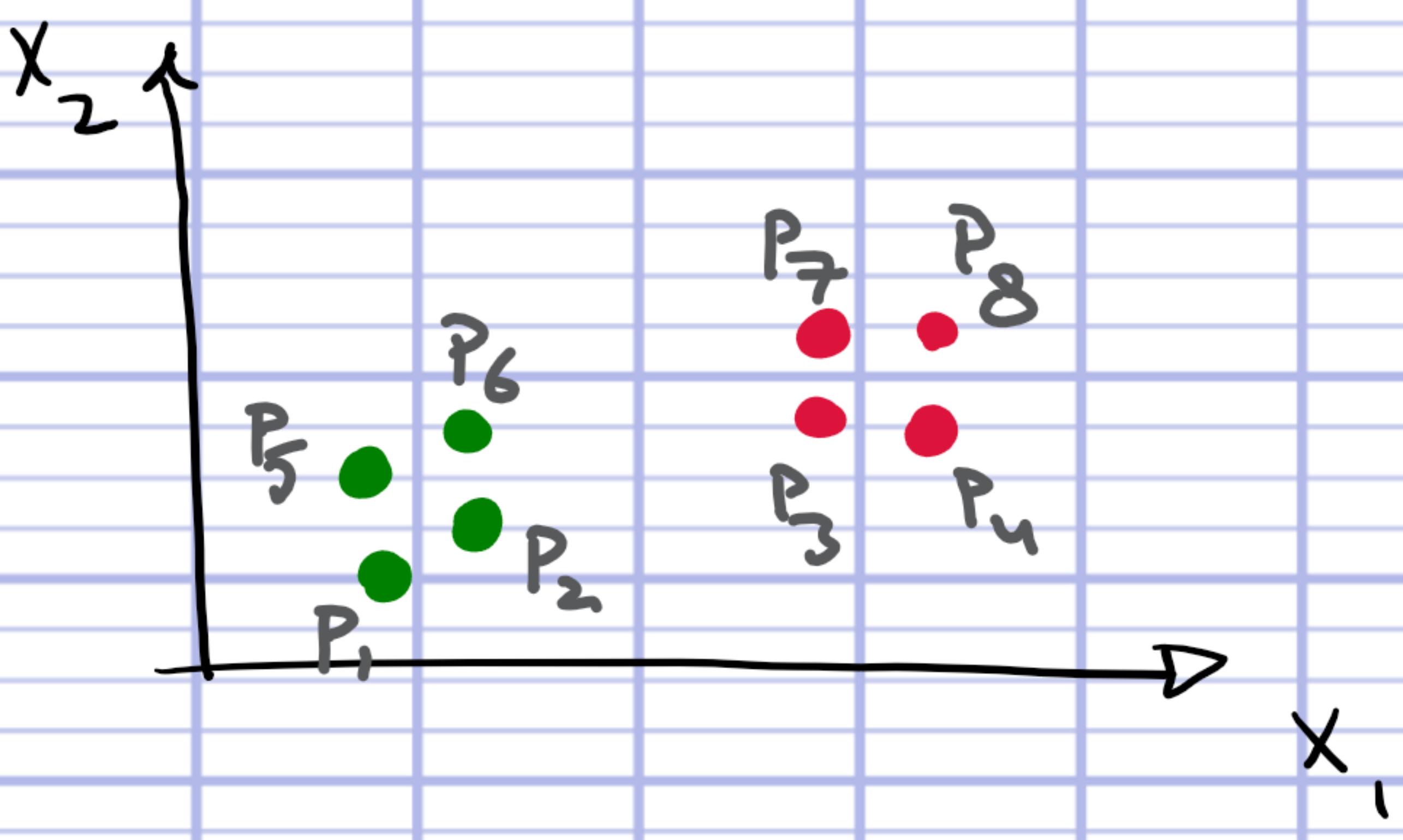


Clustering



The problem is to find a Partitioning of the points (samples) that is — ?.

$$\begin{array}{l} \text{Input} \\ \hline \end{array} \times = \left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_{n_s} \end{array} \right] \rightarrow \begin{array}{l} \text{Sample 1} \\ z \\ \vdots \end{array}$$

& implicitly we assume a distance d

$$d(x_1, x_2)$$

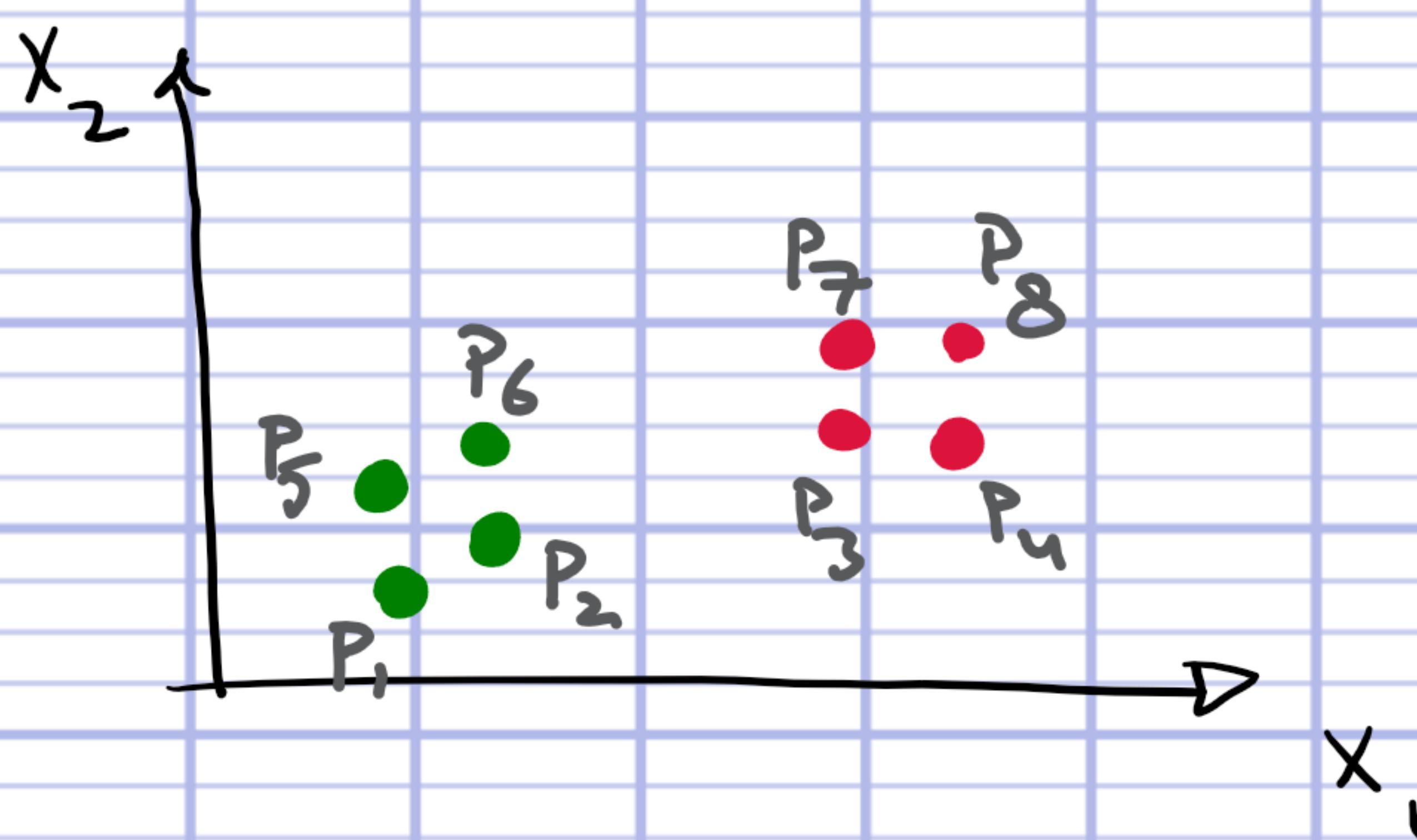
* Sometimes it is ok if it does not satisfy $d(x, y) \leq d(x, z) + d(z, y)$

Output Partitioning P such that

if $x_i \in X_j$ are in the same cluster,

then $P(x_i) = P(x_i)$.

Example



① $P(x_i) = 1 \rightarrow$ All in the same partition.

② $P(x_i) = i \rightarrow$ Each partitioned individually.

③ $P(x_1) = P(x_2) = P(x_5) = P(x_6) = 1$

$P(x_3) = P(x_4) = P(x_7) = P(x_8) = 2$.

This problem does not have a unique answer.

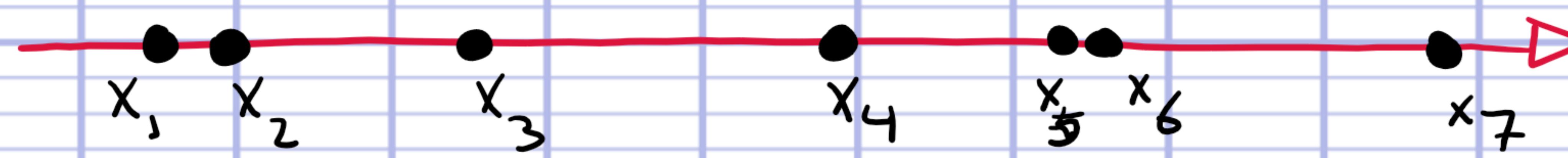
We need to think about what our expectations are.

We will go through some of the algorithms.

Single Linkage

For simplicity take a 1-D example

e.g. x_i is a number:



* SL (m) mis an input.

Algorithm:

① $P(x_i) = i$

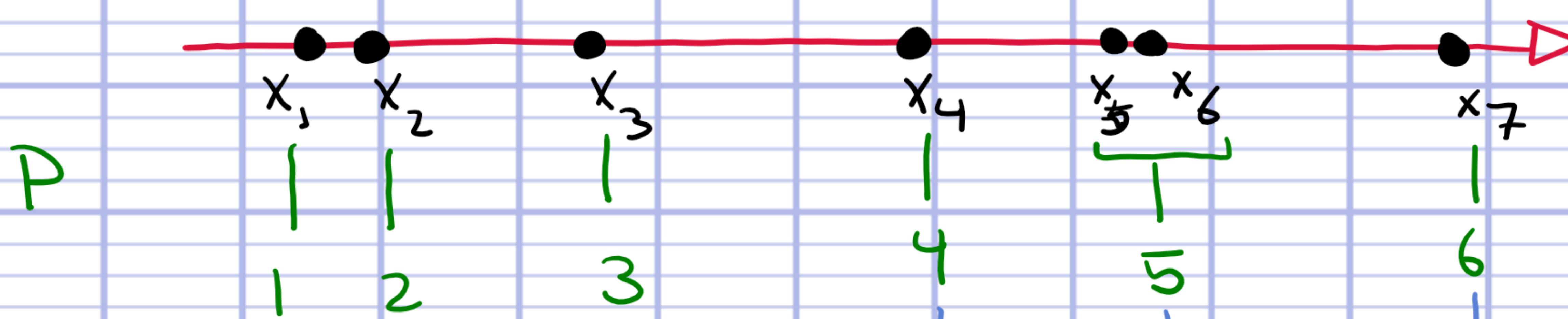
② Merge the two closest samples into one and $m=m-1$

③ Repeat 2 until $m=0$

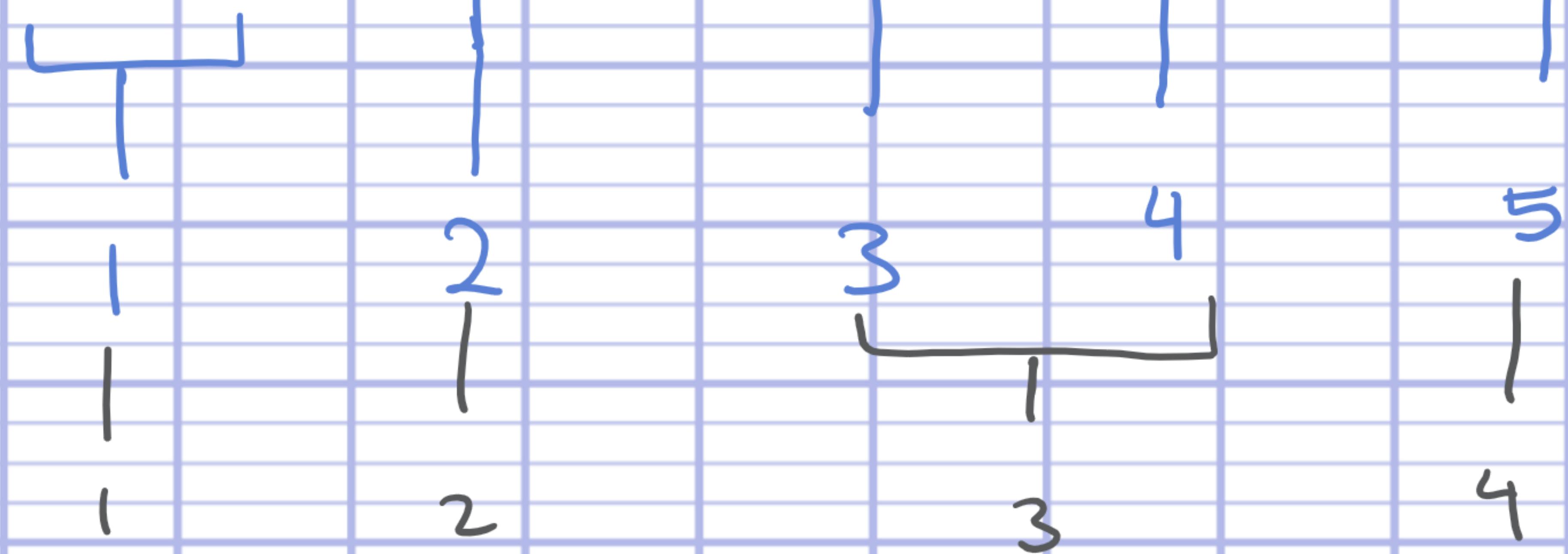
* We need the distance

* For the distance of two clusters, take the closest two points.

$m=1$



$m=2$



$m=3$



:

Some remarks

① $m \leq n_s$

② $m=0 \rightarrow P(x_i) = i$

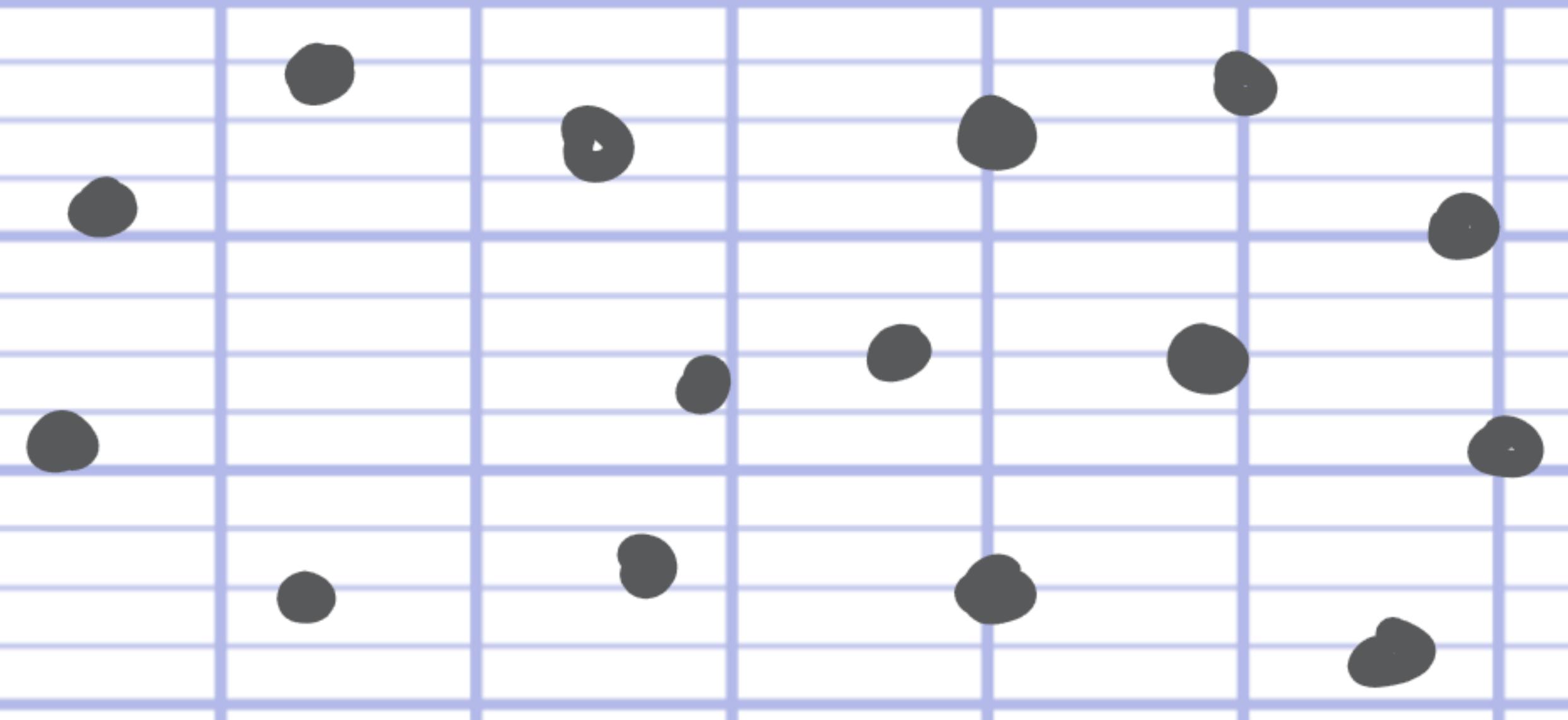
③ $m=n_s - 1 \Rightarrow P(x_i) = 1$

④ # clusters = $n_s - m$

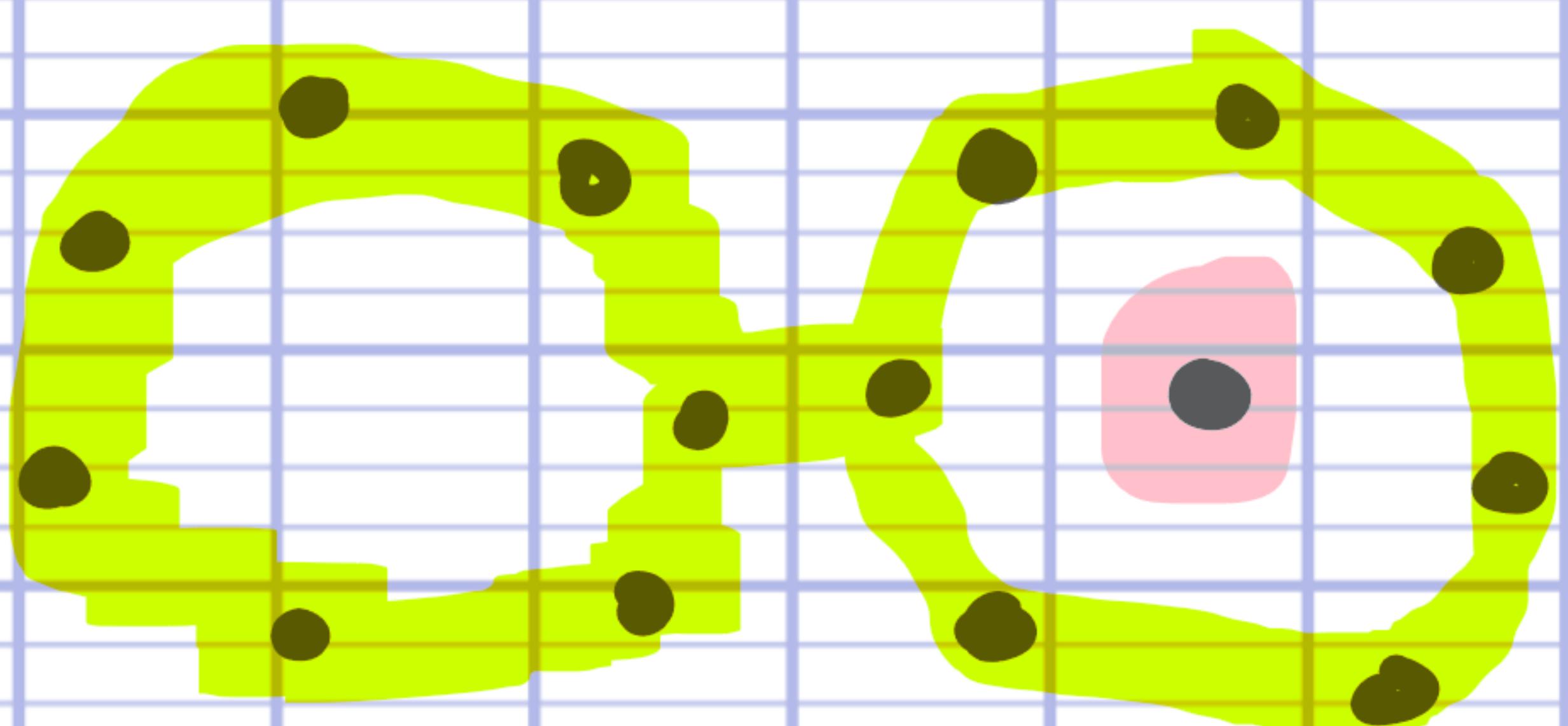
Assignment: What's the scaling of the complexity of the SL?

(Hint: How many times do we need to calculate $d(x_i, x_j)$?).

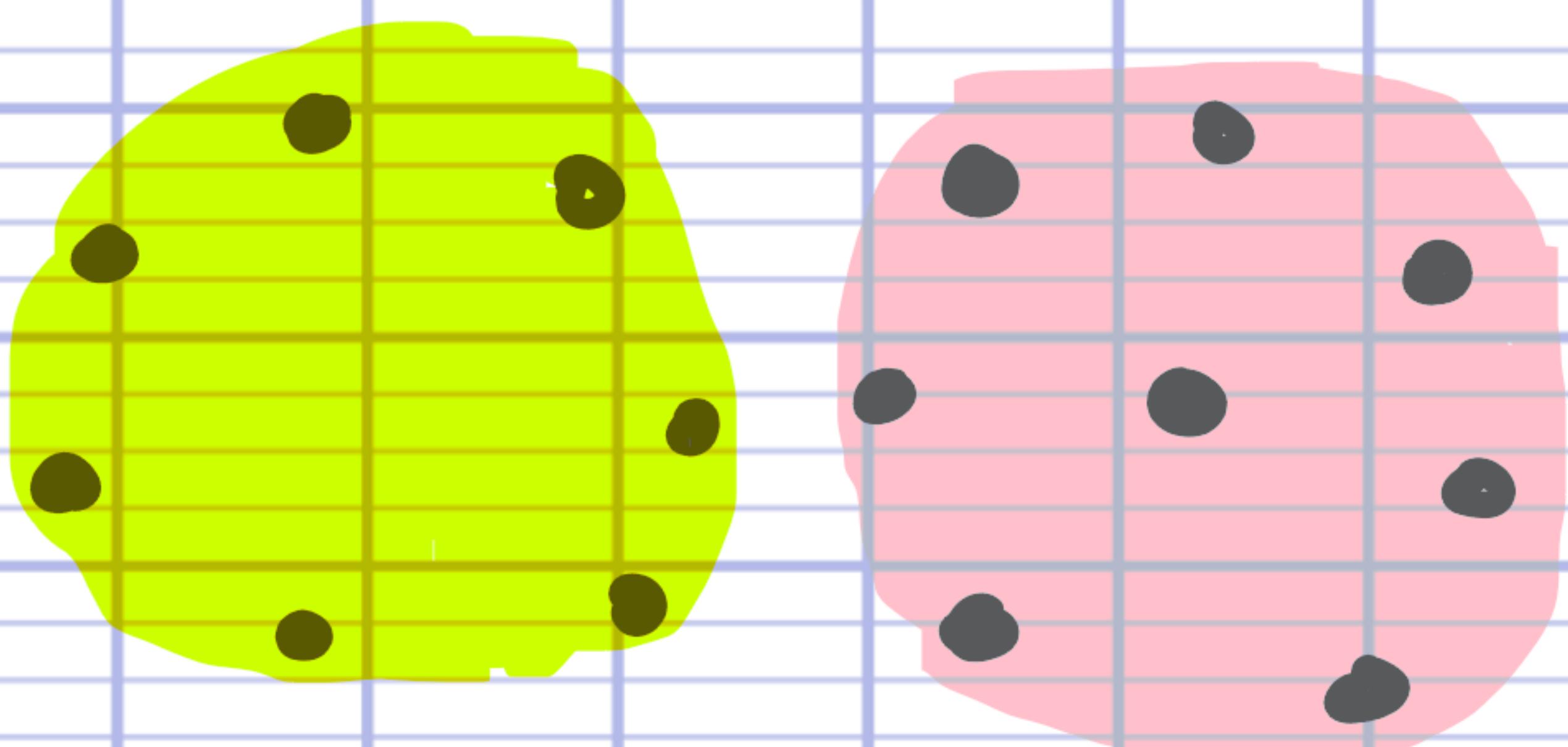
Example 2



Do this for $m=n-2 \rightarrow$ Want two clusters



whereas it may be more sensible (for some applications)
to do it like this



What's the limit of SL?

It does not recognize "Centrality".

It only takes into account the distance.

K-Means

Input: $\{x_i\}$ & $k \rightarrow \# \text{ clusters}$

Algorithm:

* Pick k points @ random & set the centers

$$c_i = x_{\tilde{j}} \quad \tilde{j} \text{ is the random index.}$$

Set clustering as $P(x_i) = \underset{j}{\operatorname{argmin}} d(c_j, x_i)$

Set the centers as

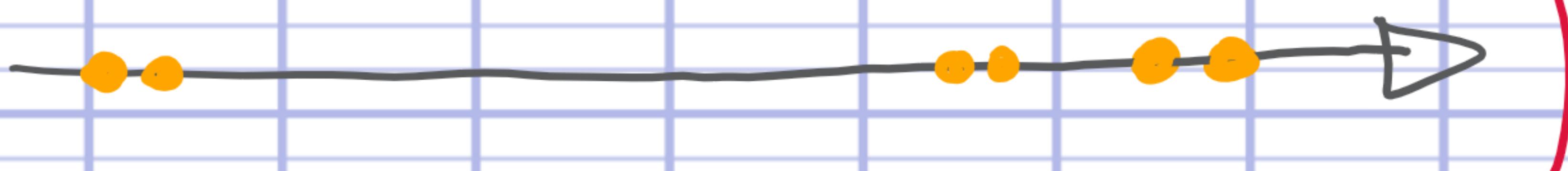
$$c_i = \frac{1}{n_i} \sum_{P(x_j)=i} x_j \rightarrow \text{Average of the points in that cluster i.e. center of mass.}$$

$n_i = \text{Points in } i^{\text{th}}$ cluster.

Repeat until converges.

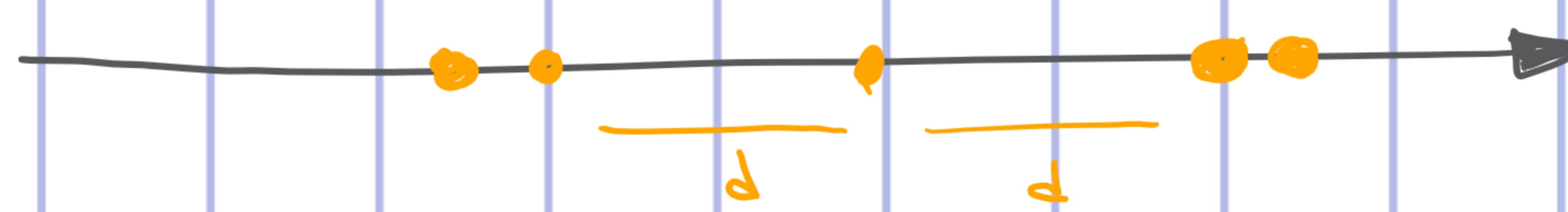
Question: Does it always converge?

Question: Are there local minimum for the optimization?

(hint: Consider 

Question: Complexity?

Question: What if there eq. distances?



Sometimes we have points that don't fit
any cluster or multiple clusters.

Some algorithms
flag these as
anomaly.

Some algorithms
assign a probability
to each sample for
being in any of
the clusters.

"Soft clustering"

[Add DBSCAN] → Anomaly

[Add Gaussian Mixture] → Soft

Assignment

Find a measure that quantifies
how well a clustering algorithm is
doing?

(Hint: First characterize our expectations
from a good clustering)