# Why is this reasonable? (Lloyd's Algo)

$$\hat{G}_1, \cdots, \hat{G}_K = \underset{G_1, \cdots, G_K}{\text{argmin}} \sum_k N_k \sum_{i \in G_k} \|X_i - \bar{X}_k\|^2$$

$\underbrace{\hspace{3cm}}_{W}$

## Generalize:

$$\hat{G}_1, \cdots, \hat{G}_K, \hat{m}_1, \cdots, \hat{m}_k = \underset{G_k, m_k}{\text{argmin}} \boxed{\sum_k N_k \sum_{i \in G_k} \|X_k - m_k\|^2}$$

## Fact: 

① given $G_k$s the best values for $m_k$

is $m_k = \bar{X}_k$

② Given fixed values for $m_k$, the

best $G_k$ is to assign pts to cluster w/

closest $m_k$
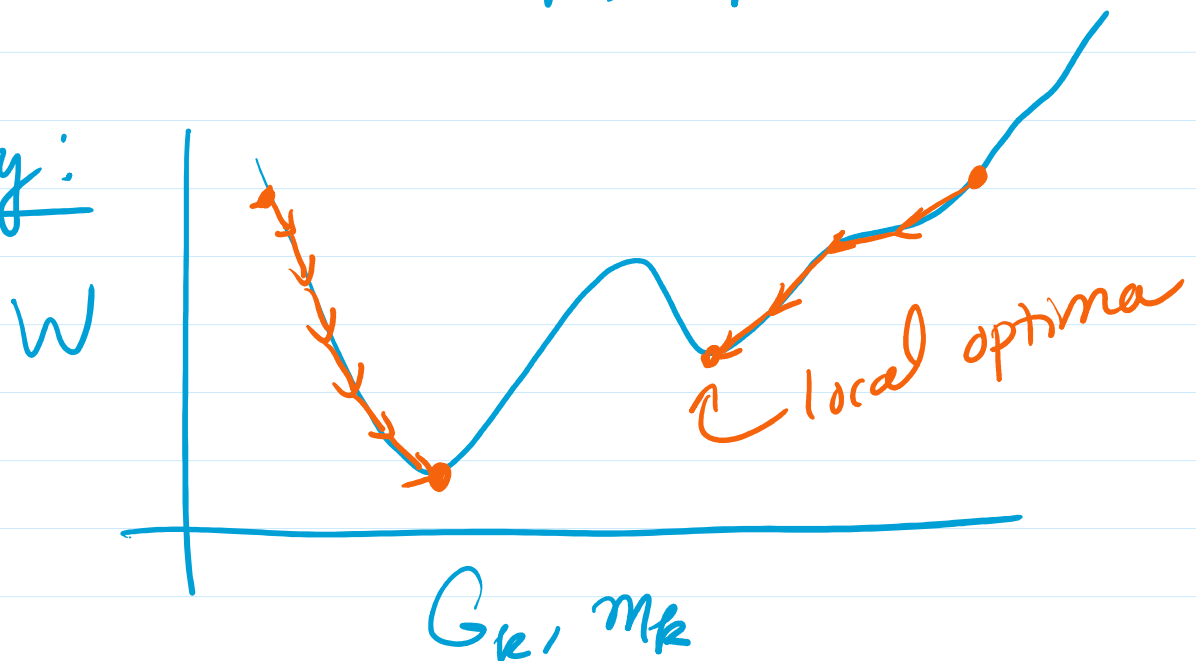
closest $m_k$

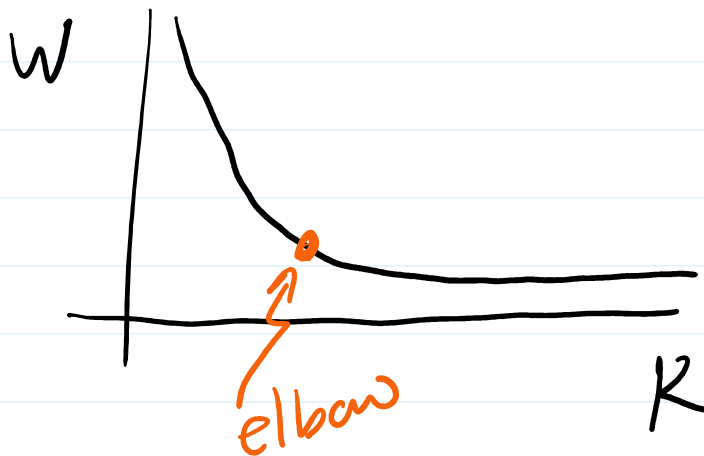Because of these two facts, each step of Lloyd's algo will decrease $W$

Ideally:



$W$

$G_k, m_k$

Reality:



$W$

$\leftarrow$ local optima

$G_k, m_k$

$G_k, m_k$

To avoid local minima try multiple random initializations and choose clustering w/ lowest val of W.

How do I choose K?

Can't choose K to minimize W, increasing K will always decrease W:



What about non-numeric data or non-euclidean dissim metric.

# All you need is D

## K-mediods

Step 0: initialization

choose $K$ points $i_k^*$  $k = 1, ..., K$
randomly in my data

↳ "representatives" of my clusters —
called medioids

For $t = 1, 2, 3, ...$

Step 1: assignment

assign each data point $i$ to cluster $G_k$
if the dissim btwn $i$ and $i_k^*$ is the
smallest among all choices of $k$

assign $i$ to $G_k$ if $D_{i i_k^*} \leq D_{i i_{k'}^*} \; \forall k'$

assign ... -k ... -u_k ... u_k' ...

## Step 2 : Update

Choose new medioids for each group as pt w/ least total dissim to all other pts in cluster

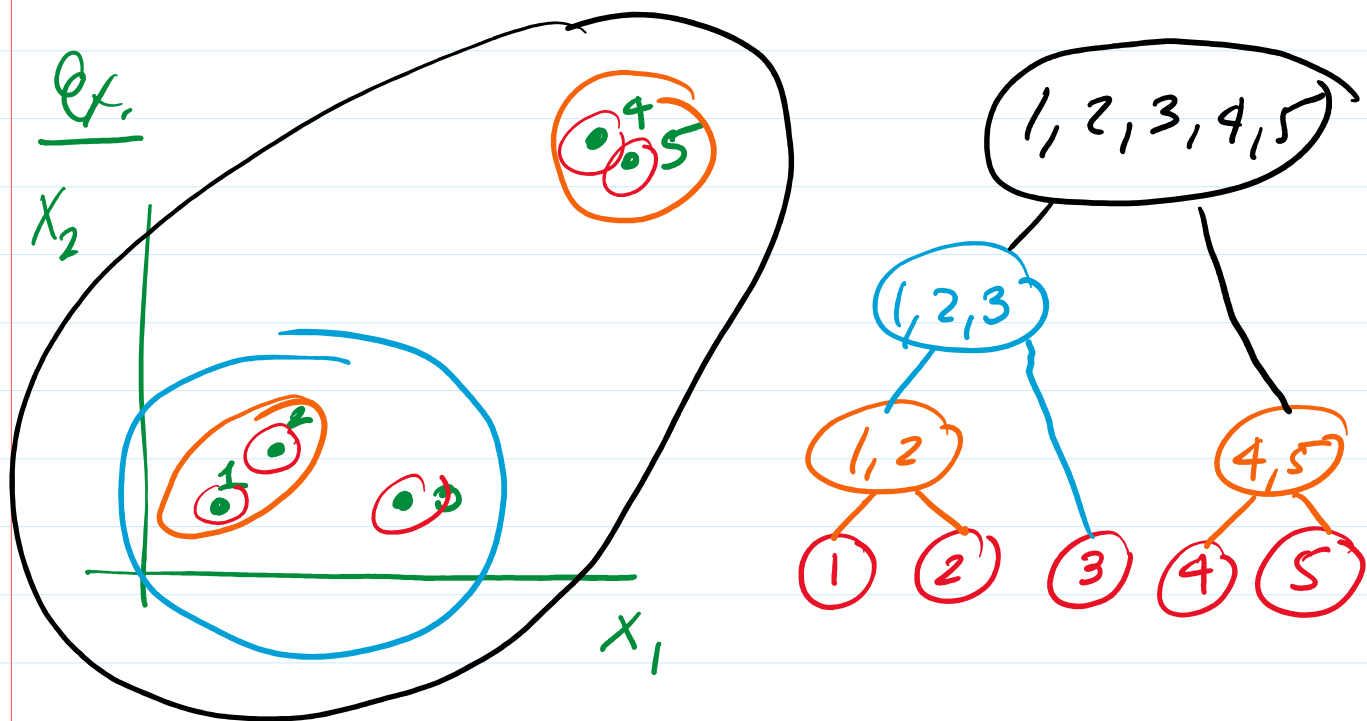$$i_k^* = \underset{i \in G_k}{\text{argmin}} \sum_{i' \in G_k} D_{ii'} .$$

## Hierarchical Clustering

Build up a collection (hierarchy) of nested clusters.

### Agglomerative clustering : bottom-up

①  start w/ each pt as individual cluster

②  merge clusters that are close

③  recursively do ② until everything is in a single cluster

To do this, need some measure of "closeness" of clusters

## Many ways to do:

Single-linkage: dist. btwn G and H is the min dissim btwn pts

$$d_{SL}(G,H) = \min_{\substack{i \in G \\ i' \in H}} D_{ii'}$$
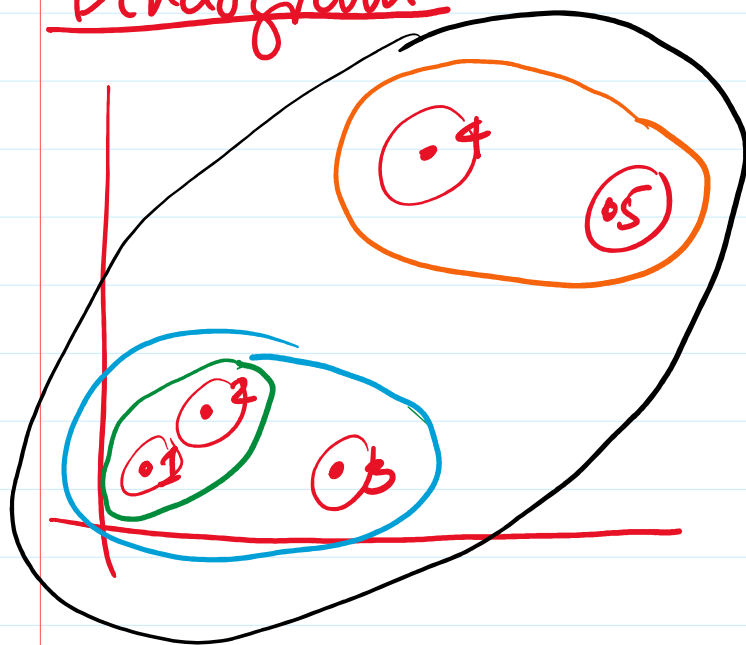
Complete-linkage: dist is max dissim

# Complete-linkage : dist is max dissim

$$d_{CL}(G, H) = \max_{\substack{i \in G \\ i' \in H}} D_{ii'}$$

# Average-Linkage

$$d_{avg}(G, H) = \frac{1}{N_G N_H} \sum_{i \in G, i' \in H} D_{ii'}$$

## Dendogram



dist btwn clusters when merged $\rightarrow$ Dendo

1   2   3   4   5