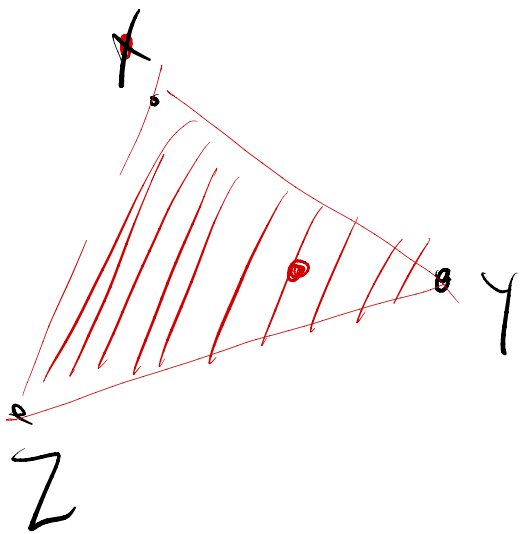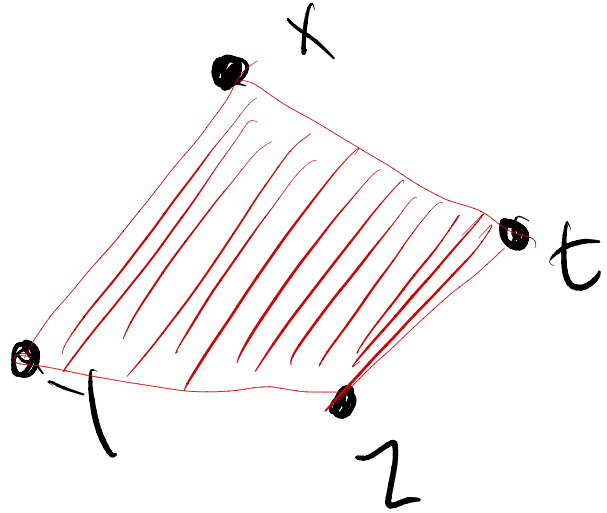# Clustering — KMEANS

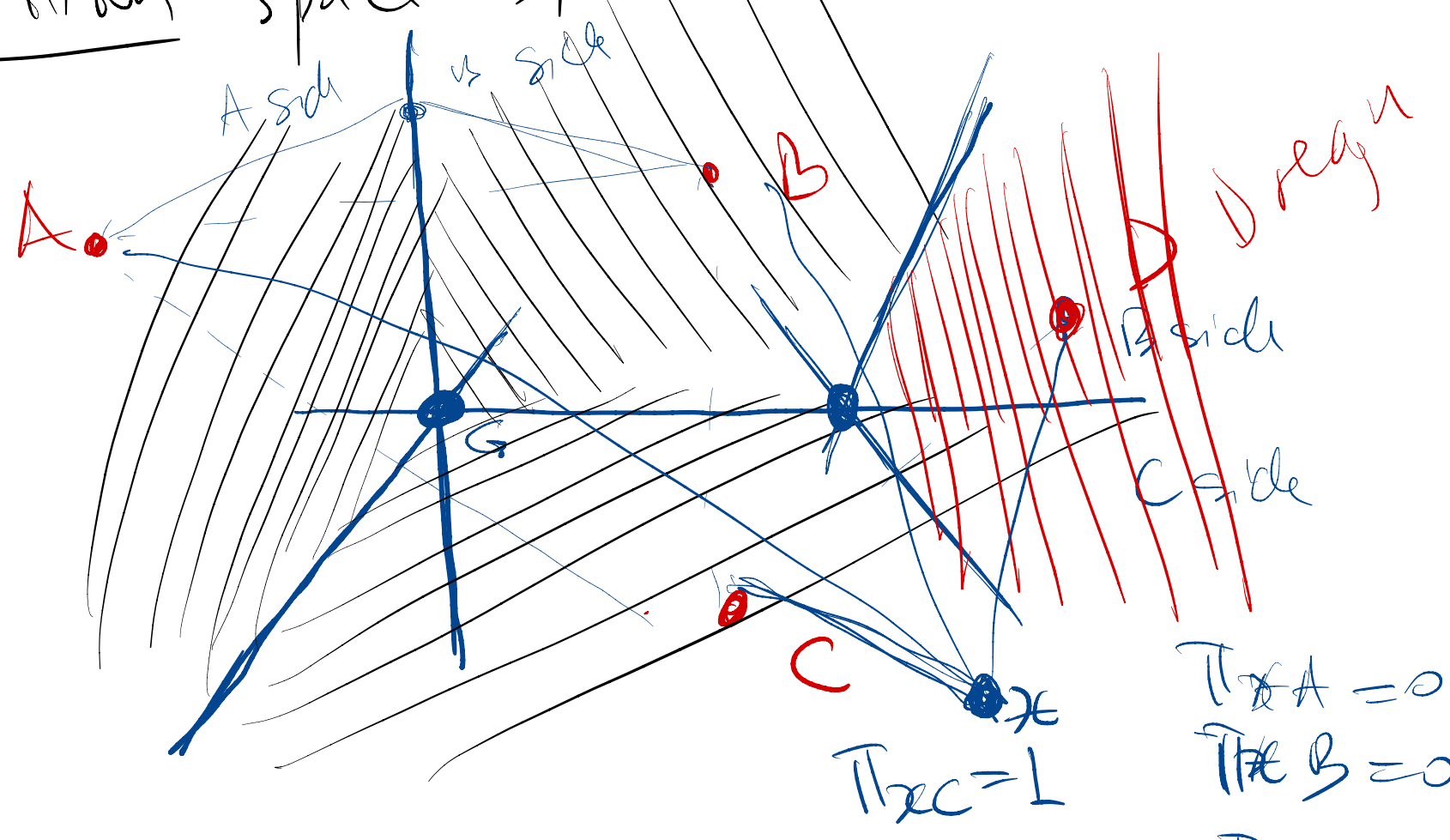convexity — closure $= \{a \mid a = c_1 x + c_2 y + c_3 z + c_4 t\}$

$c_1, c_2, c_3, c_4 \in \mathbb{R}^+$

$c_1 + c_2 + c_3 + c_4 = 1$

Space Partition by distance-idea

A B C centroids

partition space by num-dist to $A, B, C$.



A side   B side

A

B

D region

D

B side

C side

C

$\Pi_{xC} = 1$

$\Pi_{xA} = 0$

$\Pi_{xB} = 0$

$\Pi_{xD} = 0$

$x_i$ = datapoint $i$ $\quad$ $i = 1:N$ $\quad$ $K = $ # of clusters

$\mu_K$ = centroid for k-th group/cluster
$$k = 1:K$$

$\pi_{ik}$ = membership indicator $= \begin{cases} 1 & \text{if } x_i \rightarrow \text{cluster } k \\ 0 & \text{otherwise} \end{cases}$
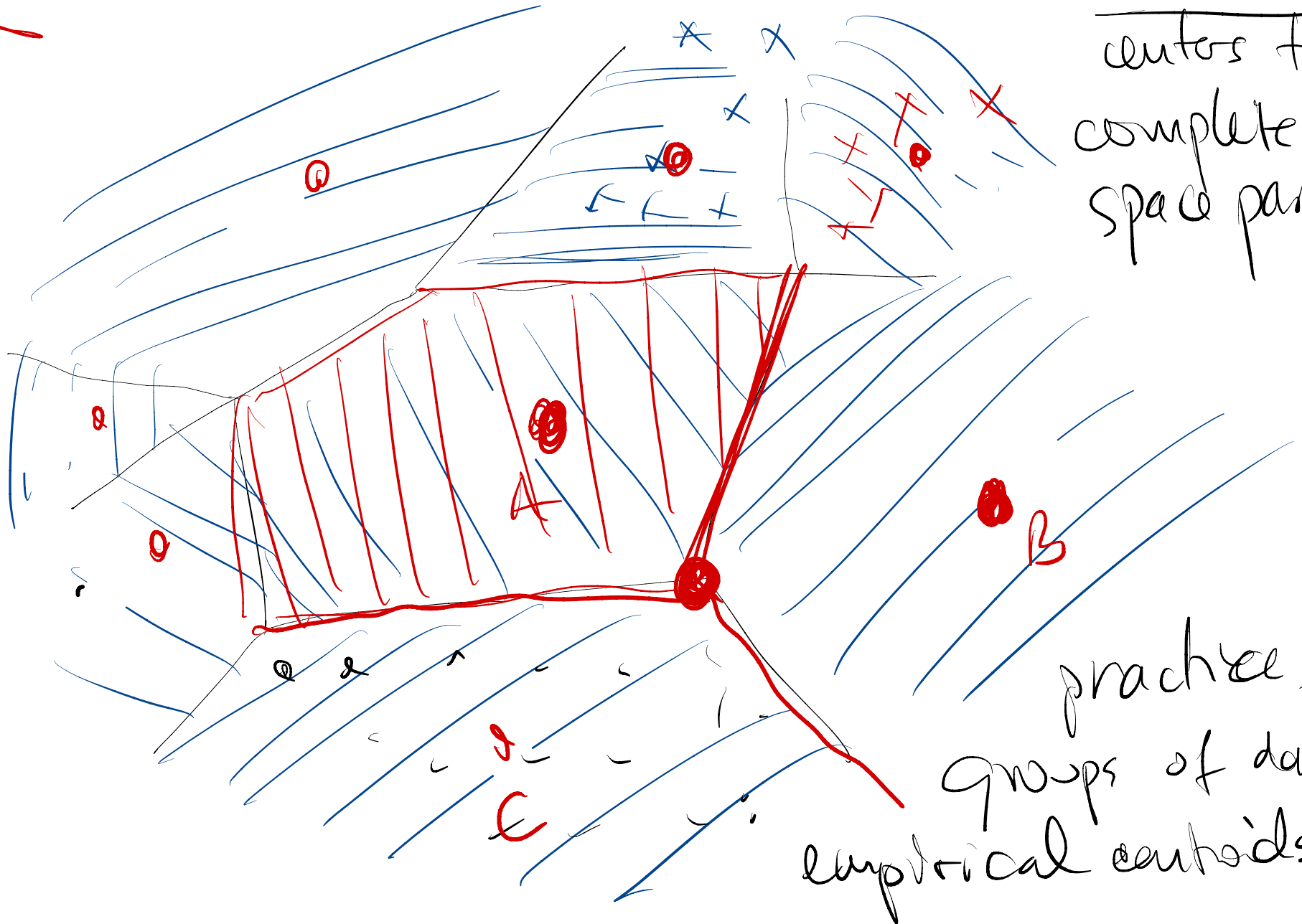
---

Decide Membership / Given centroids $\mu_K$

$\boxed{\text{for } x_i}$

<span style="color:red">E step</span>

$$\text{cluster} = \arg\min_k \{ \text{dist}(x_i, \mu_k) \}$$

$$\pi_{ik} = 1 \text{ for that argmin, } 0 \text{ for the others}$$

# Decide Centroids / Given Membership Trick

M-step

theoretical
centers for
complete
space partition

practice:
groups of datapoints
empirical centroids

Q Q Q Q Q

A B C

$$\mu_k = \text{avg of datapoints in cluster } k$$

$$k\text{-fixed} = \frac{\sum_{x \in \text{cluster } k} x}{\# x \text{ in cluster } k} \qquad \rightarrow \text{sum of vectors}$$

$$(\text{goal}) = \frac{\sum_{i=1}^{N} \pi_{ik} \cdot x_i}{\sum_{i=1}^{N} \pi_{ik}}$$

$$\text{optimal for}$$
$$\text{dist}(x, \mu) = \|x - \mu\|^2$$

Medium

K Means: improve iteratively E step / M step

until convergence

$\rightarrow \text{dist}_2$  Euc

$\text{SSE} = \text{obs} \left[ \text{MIN} \left\{ \sum_i \pi_{ik} \cdot \text{dist}(x_i, \mu_k) \right\} + \right]$  fix k

proof · easy : E step $\Rightarrow$ optimal $\pi_{ik}$ | given $\mu_k$

· easy M step $\Rightarrow$ optimal $\mu_k$ | given $\pi_{ik}$

· not easy global objective

$$\text{MIN} \left\{ \sum_{i,j,k} (\pi_{ik} \cdot \pi_{jk}) \, dist(x_i, x_j) \right.$$
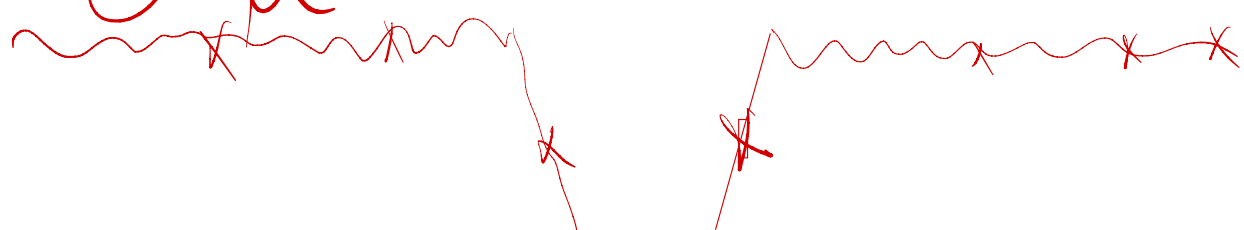
points in same cluster (k)
have small dist on avg

OBS :


global min    local min

$$\text{OBJ SSE} \quad \sum_i \sum_k \pi_{ik} \|x_i - \mu_k\|^2$$

$$= \sum_k \left\{ \sum_{\substack{i \\ \pi_{ik}=1}} \|x_i - \mu_k\|^2 \right\}$$

$$\text{Min} = \sum_k \left\{ \sum_i \pi_{ik} \|x_i - \mu_k\|^2 \right\}$$

$$\frac{\partial \text{OBJ}}{\partial \mu} = 0 \quad \overset{\text{exercise}}{\Longrightarrow} \quad \mu_k = \frac{\sum \pi_{ik} \, x_i}{\sum \pi_{ik}}$$

Evaluation Clustering    idea ① All pairs    $\binom{W}{2}$

All pairs
same cluster

$$\sum_{i,j} sim(X_i, X_j)$$

$\pi_{ik} = \pi_{jk} \ \forall k$

big

All pairs
diff cluster

$$\sum_{i,j} sim(i,j)$$

$\pi_{ik} \neq \pi_{jk} \ \exists k$

small

idea 2 — supervised (labels/tags correlated with desired clusters)

BETTER EVAL

Same label ⟹ same cluster

Same cluster ⟹ same/similar labels

⟹ small label variance

* deal with granularity/specificity
— clusters
— tags/labels