# ML FUND. DAY 3

$$y = \begin{cases} A \\ B \end{cases}$$



IMBALANCED DIST.

1) UNDERSAMPLING

2) OVERSAMPLING

(RL) ——————— ML

SUPERVISED      UNSUPERVISED      SELF-SUPERVISED

$y \longrightarrow \mathbb{R}$

$\{A, \ldots, Z\}$

$\hat{f}(x) = \hat{y}$

$\cancel{y}$

$\boxed{x}$

$? \, y \, ?$

# CURSE OF DIMENSIONALITY

$$X = \begin{bmatrix} 1 & \cdots & & P \\ \vdots & & & \\ N & & & \end{bmatrix}$$

$y =$ (crossed out)

PROBLEMS →
- COMPUTATIONAL BURDEN
- MULTICOLLINEARITY
- GEOMETRY

# PCA  (SCALED $X$)

EIGENVALUES
~~EIGEN~~ EIGENVECTORS  $(\lambda_1, u_1)$

$$X_{N \times P} \longrightarrow Z_{N \times \ell}$$

$$X_{N \times P}$$

$$\vdots$$

$$\ell << P$$

$$(\lambda_P, \mu_P)$$

$\Sigma_{P \times P}$  VAR-COV
MATRIX of $X$

$$\Sigma = \begin{bmatrix} \sigma^2_1 & & \sigma_{iS} \\ & \ddots & \\ & & \sigma^2_P \end{bmatrix}$$

$$\alpha = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & \cdots & P \\ & & \end{bmatrix} \cdot \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_P \end{bmatrix} = \begin{bmatrix} \alpha_1 x_{11} + \alpha_2 x_{12} + \cdots + \alpha_P x_{1P} \\ \vdots \\ \alpha_1 x_{N1} + \cdots + \alpha_P x_{NP} \end{bmatrix}$$

$$z_1 = X\alpha_1 \qquad Var\left(z_1\right) = \alpha_1^T \Sigma \alpha_1$$

$$\max_{\alpha_1} \quad Var(z_1) = \alpha_1^T \Sigma \alpha_1$$

$$\text{CONST.} \quad \|\alpha_1\| = 1$$

$$p \rightarrow 1$$

$$z_2 = X\alpha_2$$

$$\max_{\alpha_2} \quad Var(z_2)$$

$$\text{CONS.} \quad - \|\alpha_2\| = 1$$
$$- \alpha_1^T \alpha_2 = 0$$

$$\left( \Sigma \right)$$

SPECTRAL DECOMPOSITION

$$X_{N \times p} \longrightarrow Z = \begin{bmatrix} \begin{pmatrix} z_1 \end{pmatrix} & \begin{pmatrix} z_2 \end{pmatrix} & \cdots & \begin{pmatrix} z_\ell \end{pmatrix} \end{bmatrix}$$
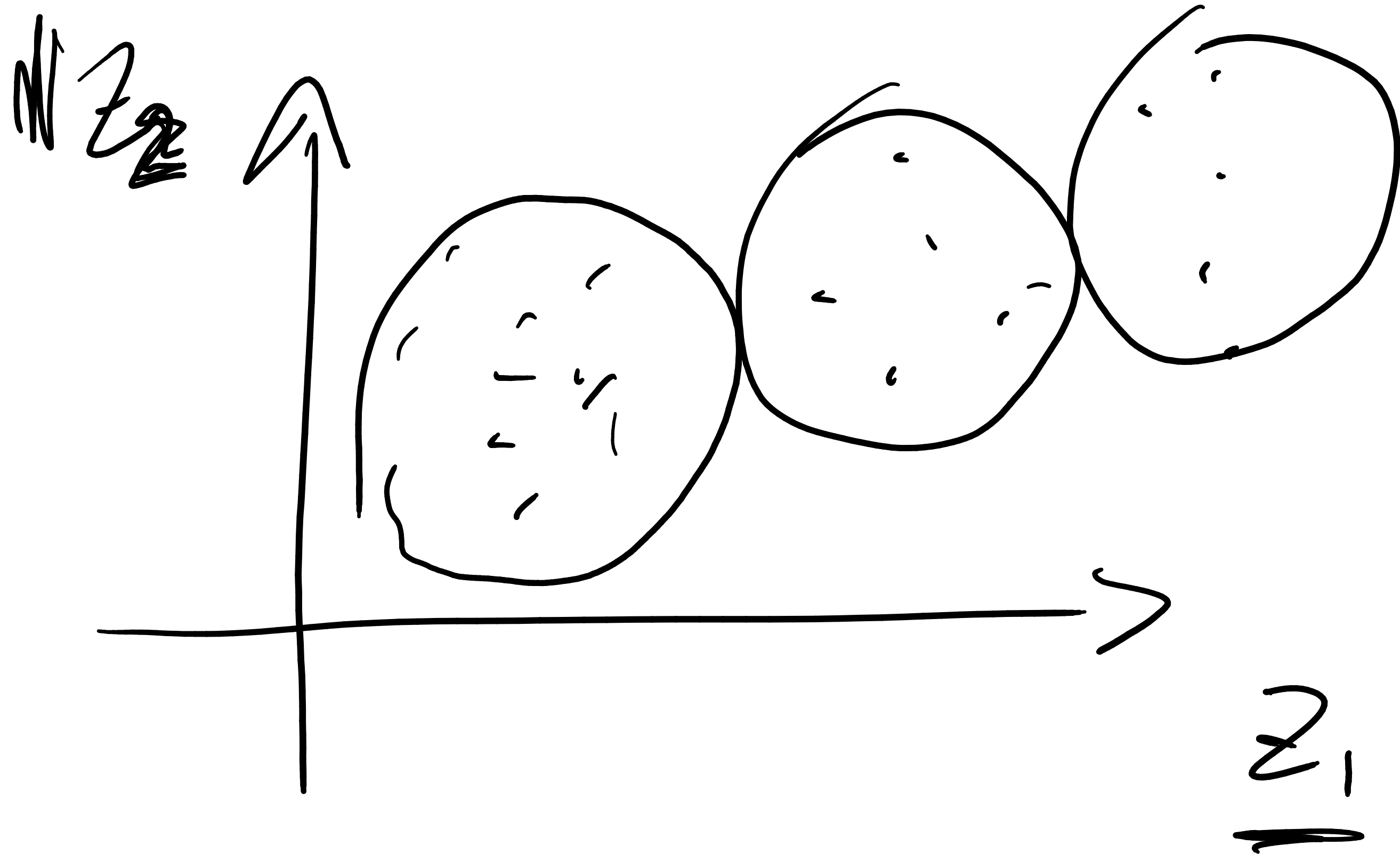
## PCA

1) FIND EIGEN(VALUES, VECTOR) FOR $\Sigma$

2) ORDER THE EIGENVECTORS by THE VALUE OF THE EIGENVALUES

$$\Sigma_Z = \begin{bmatrix} \sigma_1^2 & \cdots & 0 & 0 \\ 0 & \ddots & & 0 \\ 0 & 0 & \ddots & \sigma_\ell^2 \end{bmatrix}$$

3) USE $\ell$ pairs $(\lambda_i, \, d_i)$ TO PROJECT THE DATA ON NEW COMP. $\left( z_i = X d_i \right)$

4) USE THE NEW MATRIX $Z = \begin{bmatrix} (z_1) & \cdots & (z_\ell) \end{bmatrix}$ FOR YOUR NEEDS
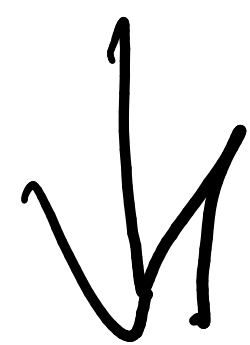
UNS.



$Z_2$

CLUSTERS ON PCA

$Z_1$

SUP.

$$X \rightarrow Z \qquad \hat{y} = f(Z)$$
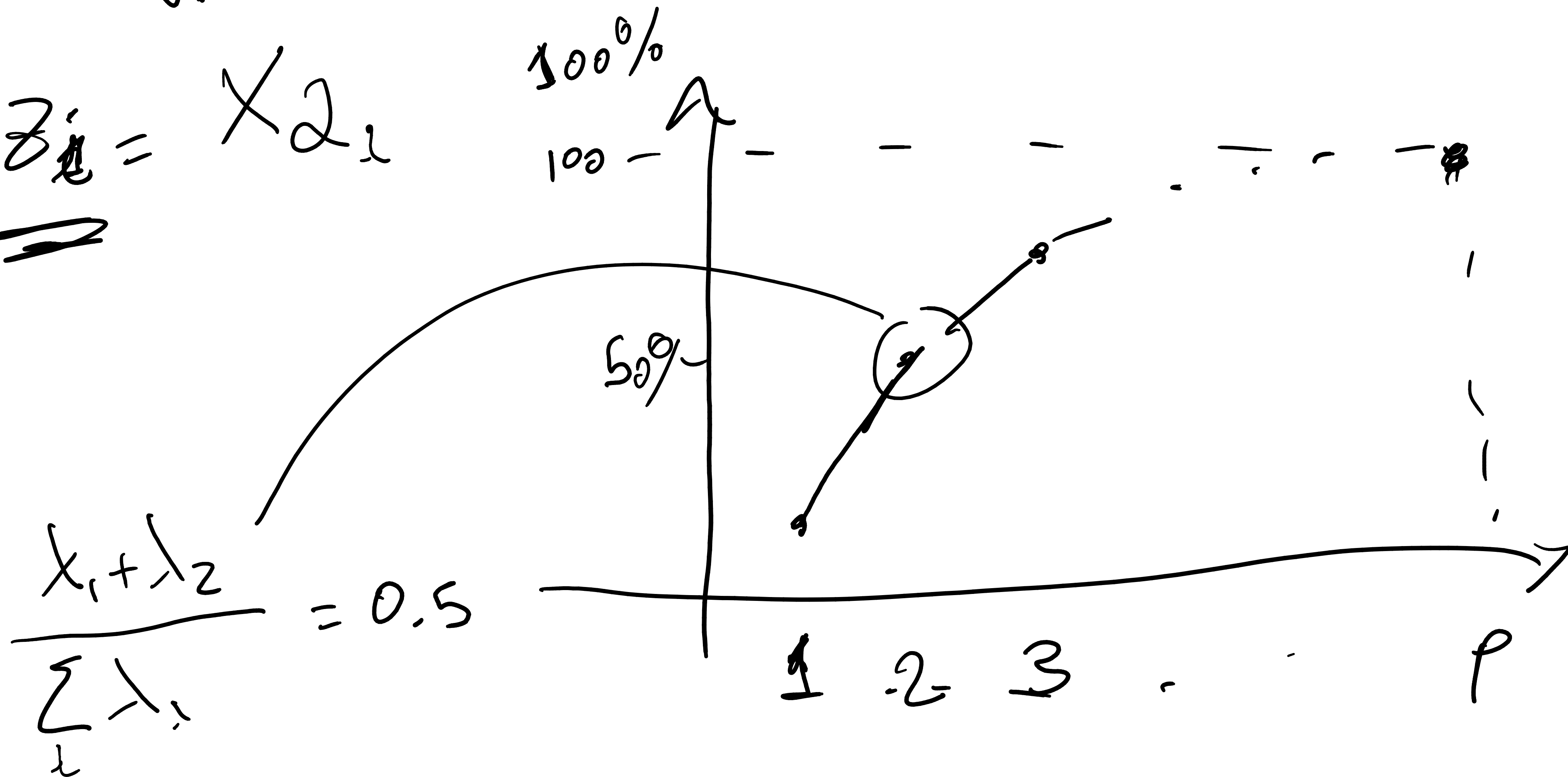
Pipeline (StandScaler(), PCA(), RF())

Pipeline.fit(X_train, y_train)

$$(\lambda_2, \quad \alpha_2)$$

$\uparrow$ VALUE $\qquad \uparrow$ VECTOR

$$\frac{\lambda_2}{\sum\limits_{i=1}^{P} \lambda_i} = \text{FRACT. of EXP. VARIANCE}$$

$$\Downarrow$$

$$\vec{z}_i = X\alpha_i$$



100%

100 ---

50%

1   2   3   .   .   P

$$\frac{\lambda_1 + \lambda_2}{\sum\limits_{i} \lambda_i} = 0.5$$

# CLUSTERING

$$X = \begin{bmatrix} \\ \\ \vdots \\ \\ \end{bmatrix} \begin{matrix} 1 \\ \\ \\ \\ N \end{matrix}$$

$$\begin{matrix} 1 & & P \end{matrix}$$

$X_1 \in \mathbb{R}$

$$X = \begin{bmatrix} \\ \\ \\ \vdots \\ \end{bmatrix}$$

DISTANCE

# CLUSTERING

## (A) NON HIERARCHICAL

- \# CLUSTERS

- LOOK FOR GROUPS

## HIERARCHICAL (B)

- EACH ROW IS A CLUSTER, AGGREGATE
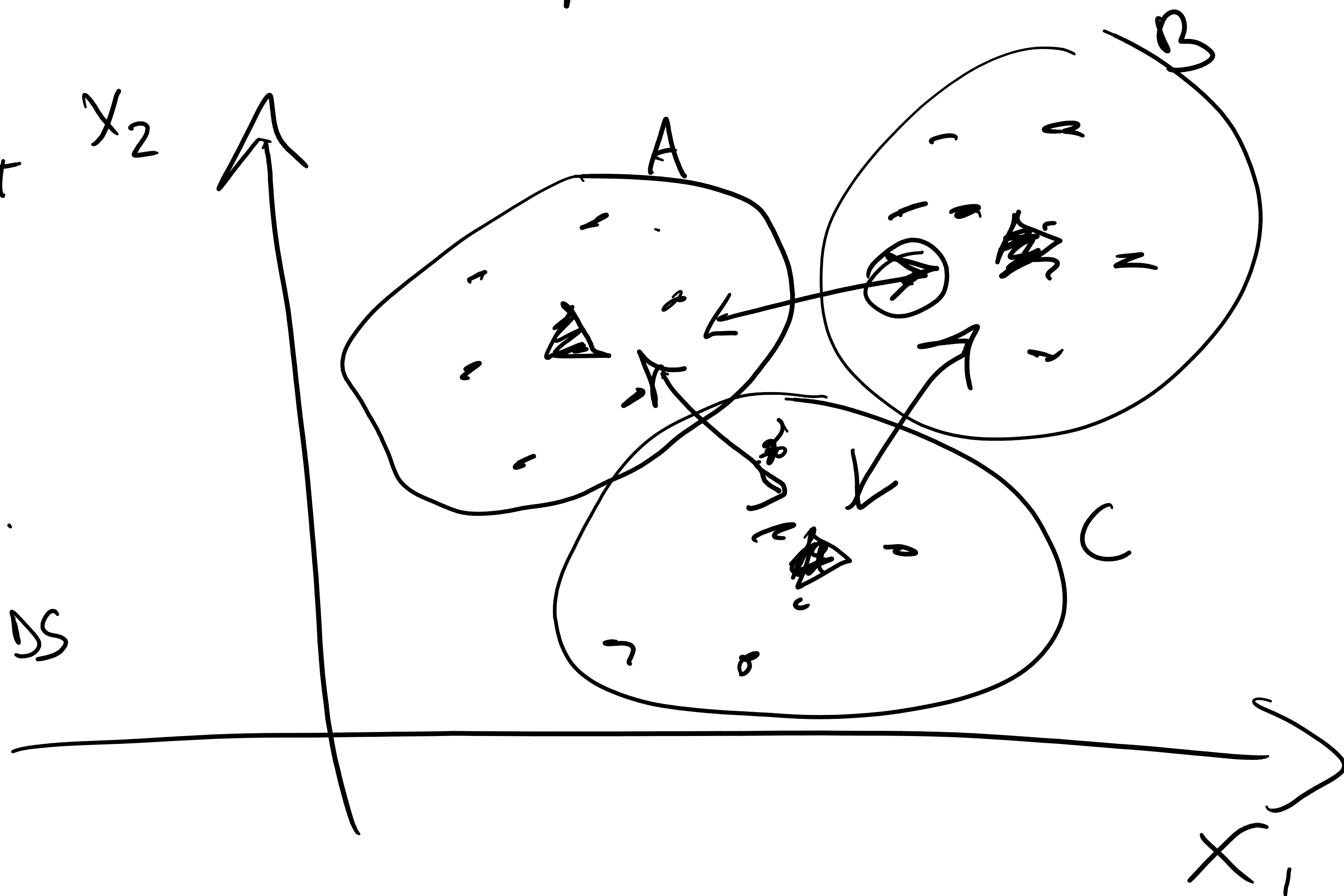
- SINGLE CLUSTER, DIVIDE

# K - MEANS
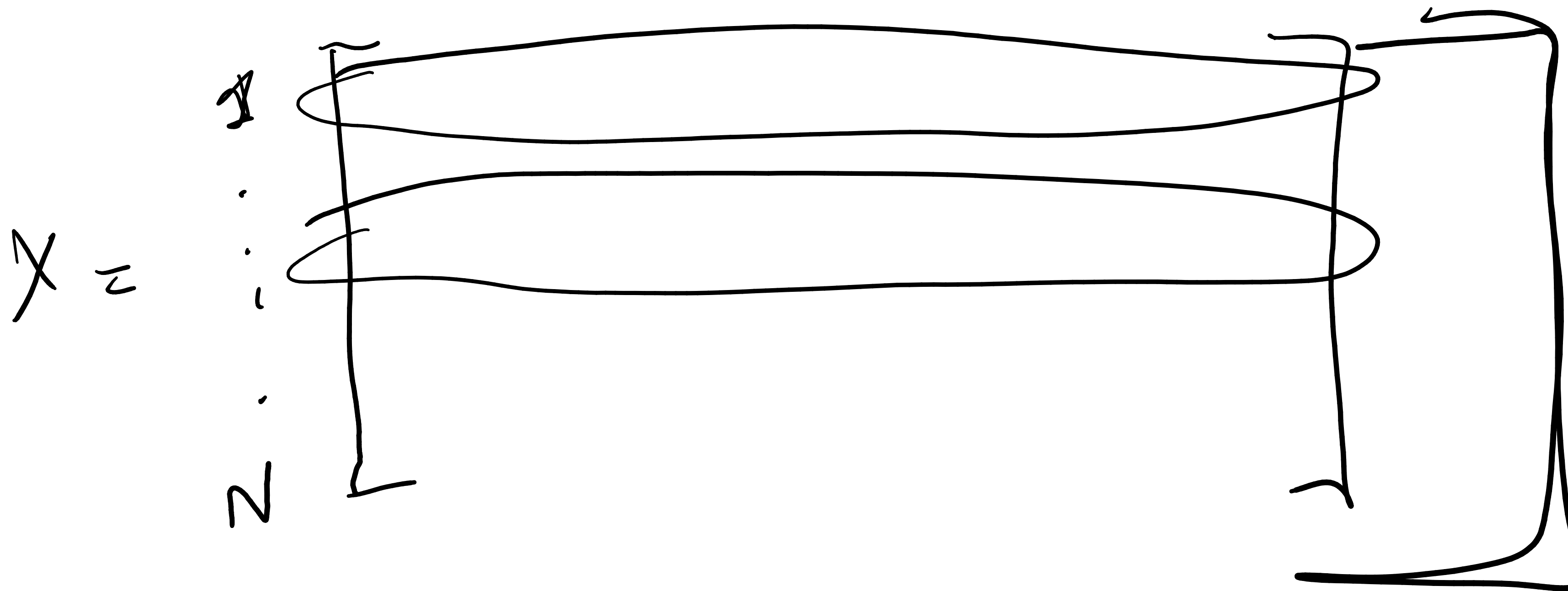
← SET A k VALUE

VARIANCE WITHIN ↓
VARIANCE BETWEEN ↑

→ CENTROIDS

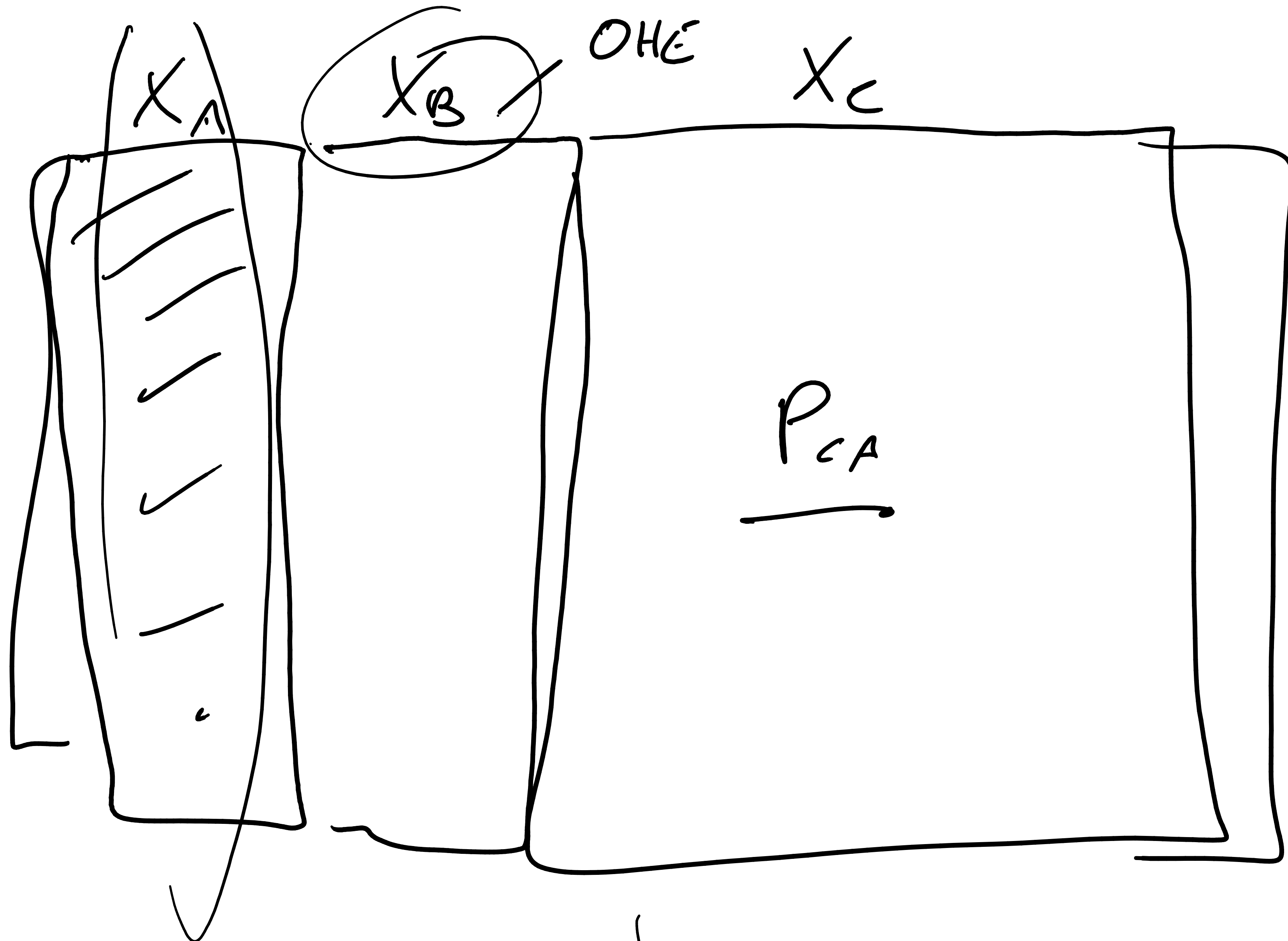— ASSIGN EACH POINT TO ONE CLUSTER
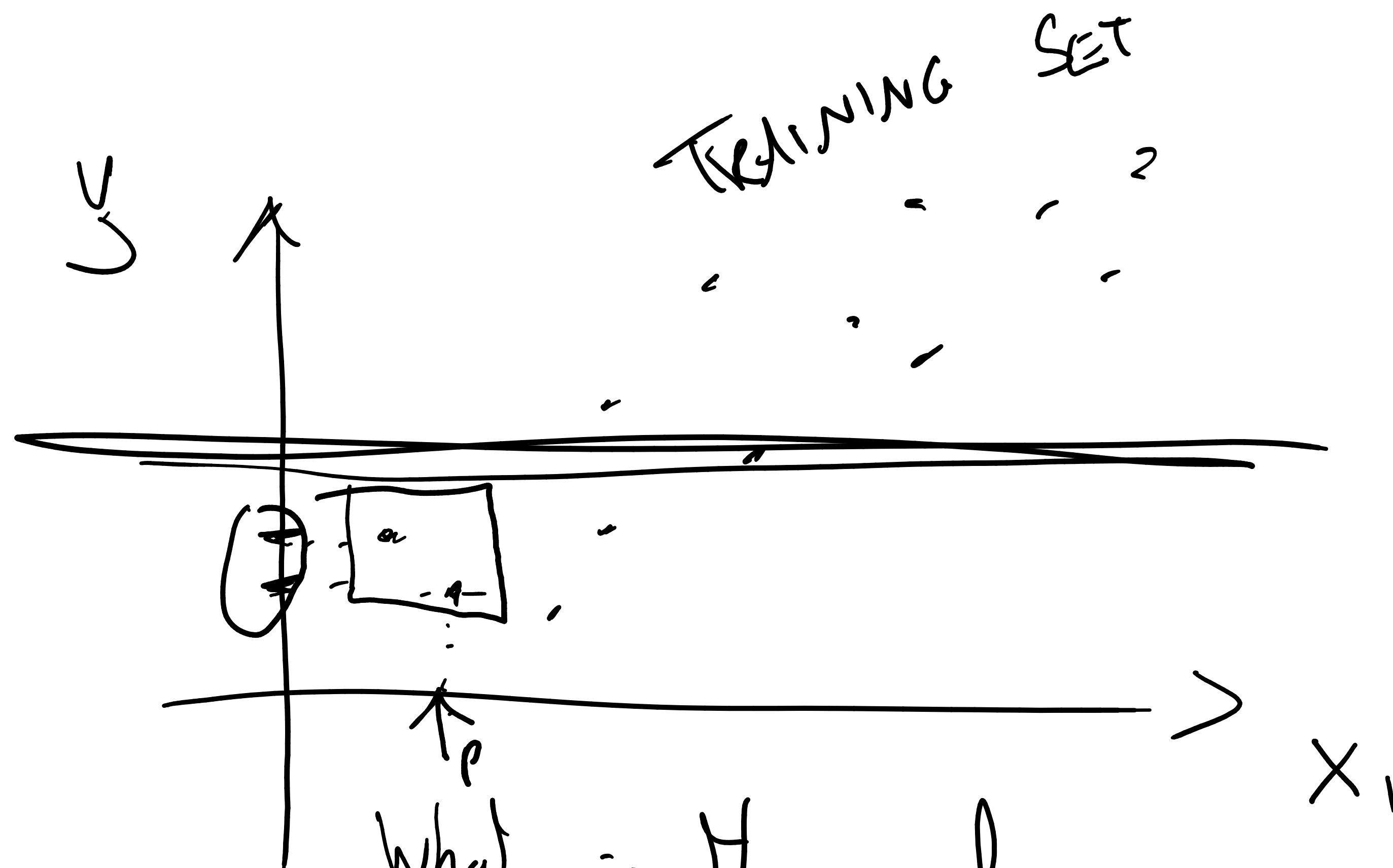
$$K = 3$$

→ COMPUTE AGAIN THE CENTROIDS

$$X = \begin{matrix} 1 \\ \vdots \\ i \\ \vdots \\ N \end{matrix} \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

- AGGREGATIVE

- DIVISIVE

$X =$ 

$X_A$  $X_B$  OHE  $X_C$

$P_{CA}$

F. SELECT.

# KNN ( REGR. CLASS. )

TRAINING SET

$y$

$x_1$

$p$

$K = 2$

$K = N$

What is the value of $y$ for the $k$ nearest points?

$$\hat{y}_p = \frac{\tilde{y}_1 + \tilde{y}_2}{2}$$

SSE