

let X be the matrix with the initial data.

let X_c be the matrix with the centered data

let X_n be the matrix with the normalized data

Let

$$\bar{T} = \left\{ \begin{array}{l} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \\ \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \end{array} \right\}$$

: non normalized
PCA

: normalized PCA

let denote $\lambda_1 > \lambda_2 > \dots > \lambda_p$

: The eigenvalues
of \bar{T} in decreasing
order

let denote a_B a unitary eigenvector
of T associated to λ_B

R₂: the fact that a_B is a unitary
vector is just $\langle a_B | a_B \rangle = 1$

Let denote V_1, \dots, V_p the P initial variables (the P columns of X_C

or x_n)

non normalized PCA

$$Z_B = \begin{cases} X_C \alpha_B & \text{the new variable} \\ \underline{x_n \alpha_k} & \text{number } k \end{cases}$$

a normalized PCA

$$\text{Var}(z_B) = \frac{1}{n} z_B^T z_B$$

If we consider the non normalized PC

$$\begin{aligned}\text{Var}(z_B) &= \frac{1}{n} e_{\alpha_B}^T X_c X_c e_{\alpha_B} \\ &= e_{\alpha_B} \left(\frac{1}{n} e_{\alpha_B}^T X_c X_c \right) e_{\alpha_B} \\ &= e_{\alpha_B} \sqrt{\text{Var}_{\alpha_B}}\end{aligned}$$

with V the covariance matrix

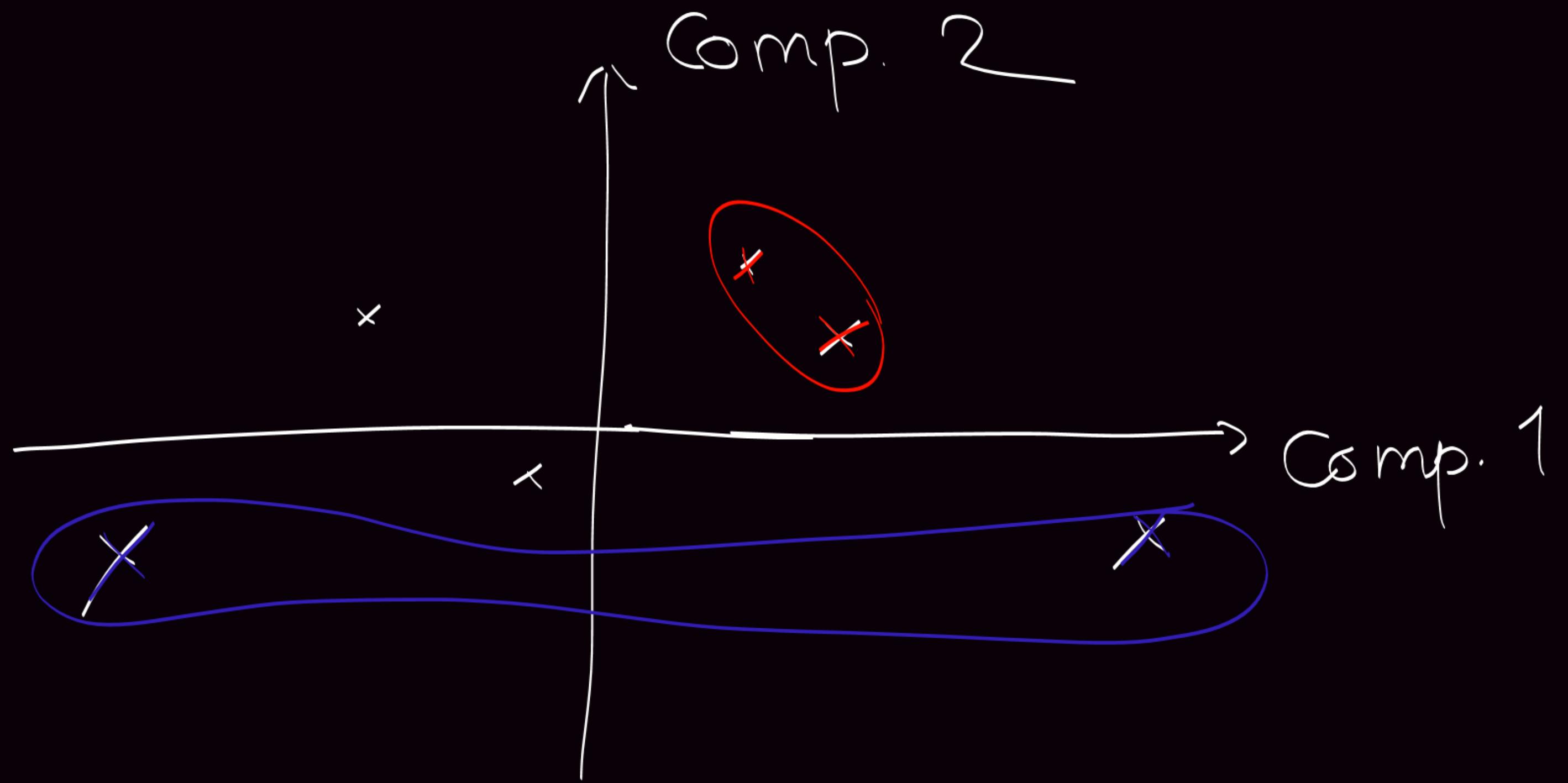
$$\text{but } \text{Var}_R = \lambda_R \cdot a_R$$

Thus

$$\begin{aligned}\text{Var}(z_R) &= t_{a_R} \lambda_R a_R \\ &= \lambda_R t_{a_R} a_R \\ &= \lambda_R\end{aligned}$$

$\sqrt{\lambda_B}$ = standard deviation
of the new variable
number B

This explains why in $P\$sdev$
(where $P = \text{princomp}(X)$) we see
the term standard deviation.



be careful:

Let u_i and u_j be two

individuals

Let h_i and h_j be their

projections

$$d^2(h_i, h_j) \leq d^2(u_i, u_j)$$

If on the first factorial
plan, two projected points
are far then, the two
individuals are far in the
global space

But if the projected points
are closed in the first factorial
plan, the associated individuals
may be far in the global space.

⇒ we need to quantify the quality
of the representation of an individual
on an axis or

on a plan.

The idea is to consider the angle formed by the vectors $\overrightarrow{Ov_i}$ and $\overrightarrow{Oh_i}$

If this angle is small \Rightarrow the projection is a good representation of the individual

For simplicity, we consider in fact the cosinus of the angle,

and better the cosinus square!

⇒ if \cos^2 near 1 \Rightarrow good representation!

Rk:

because the Factorial plan
are generated by 2 orthogonal
vectors. The \cos^2 of the angle with
the plan is equal to sum \cos^2
& the angle with each generated vector

consider the Factorial plan (β, β')

let denote $\alpha_{i, \beta, \beta'}$ the angle between

the individual u_i and its projection

on the plan (β, β')

we denote $\alpha_{i, \beta}$ the angle between u_i
and its projection on the component β

$$\cos^2 \alpha_{i,k,k'} = \cos^2 \alpha_{i,k} + \cos^2 \alpha_{i,k'}$$

comp. k'

$\rho_{i,k'}$

$\rho_{i,k,k'}$

$\rho_{i,k}$

comp. k

Rk: how to compute $\cos^2 \alpha_{i,k}$?

$$\cos^2 \alpha_{i,k} = \frac{\langle \overrightarrow{Ou_i}, \overrightarrow{Oa_k} \rangle^2}{\|\overrightarrow{Ou_i}\|^2 \|\overrightarrow{Oa_k}\|^2}$$

where u_i is the vector associated to the individual in the matrix X_C in the matrix

$$= \frac{(t_{a_k} u_{c_i})^T u_{c_i} \cdot a_k}{t_{u_{c_i}} \times u_{c_i}}$$

Contribution of an individual to

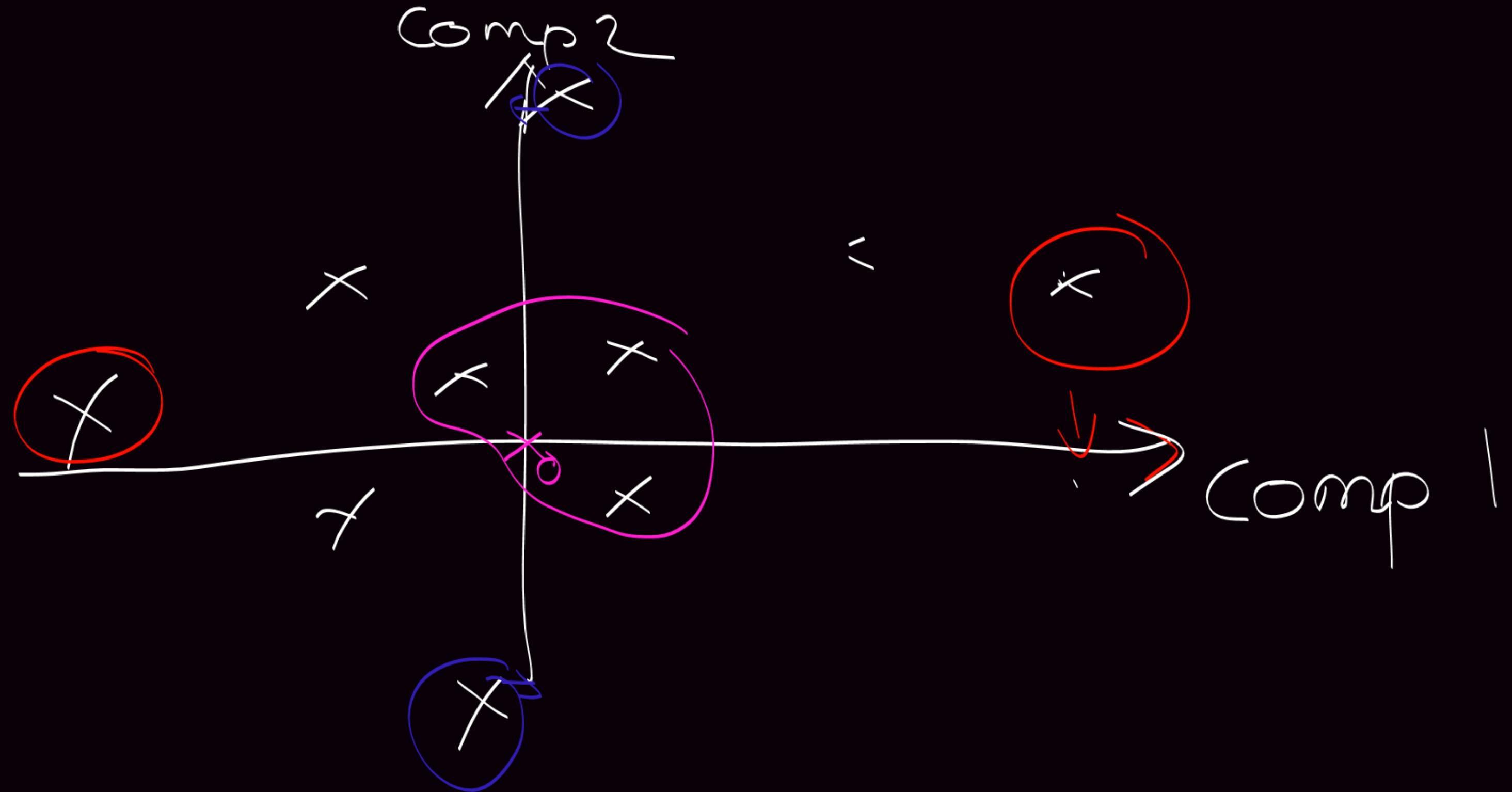
an axis Δ

The contribution of the individual on

Δ is

$$\frac{1}{n} d^2(O, R_{\Delta,i})$$

↑
the projection of individual
i and Δ .



Pierre (5.50, 6 / 14, 11.50)

annie

evelyn

(9.50, 5, 17.50, 12)

monique

(14.50 14.50
15.50 15.50)

alain

Q

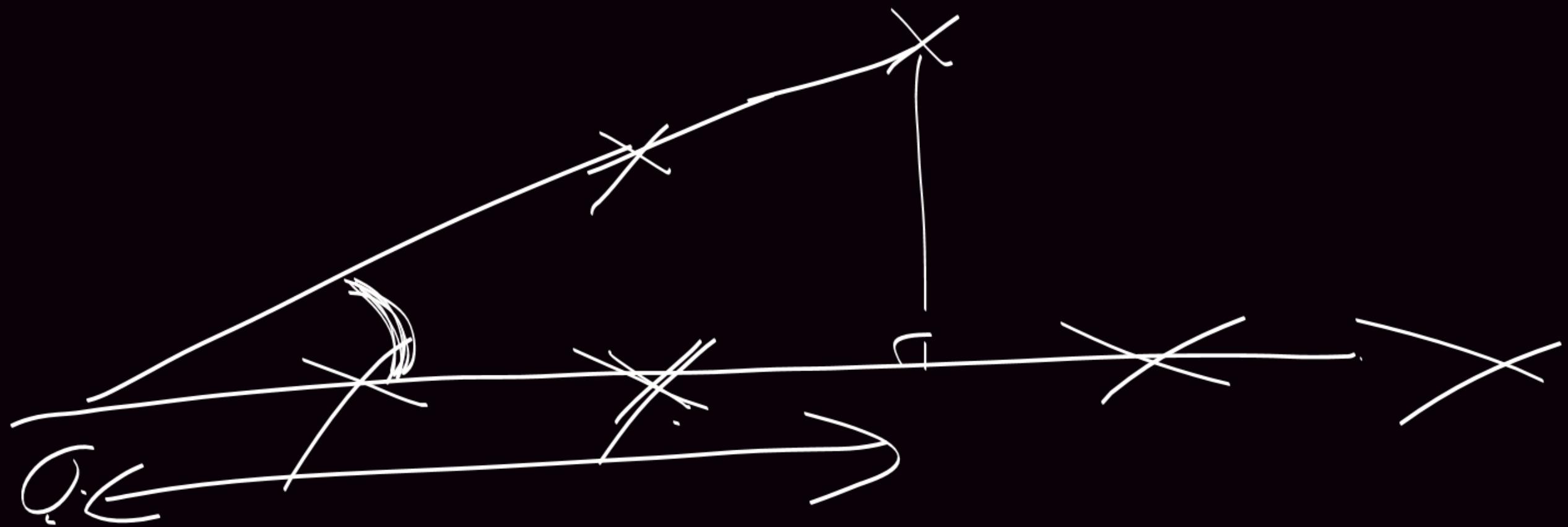
didier

brigitte (13, 12.50, 8.50, 9.50)

jean
(6, 6, 5, 5.50)

andre

(11, 10, 5.50, 7)



What happens for the variable?

$$V_{11}, V_P \rightarrow z_{11}, z_P \quad z_P = X_C Q_R$$

(centered)

$$\text{Var}(z_P) = \lambda_R$$

$$\begin{aligned} \text{cov}(V_j, z_P) &= \frac{1}{n} Q_R^T X_C \sqrt{\dots} \\ &= \frac{1}{n} t_{Q_R}^T X_C \left(\begin{array}{c} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{array} \right) \end{aligned}$$

$$\text{cov}(z_k, v_j) = t_{\text{var}} V \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

matrix of
covariance of X or matrix
of covariance de X_C

$$= \left(t V_{\text{var}} \right) \begin{pmatrix} 0 \\ i \\ 0 \\ n \\ 0 \\ i \end{pmatrix}$$

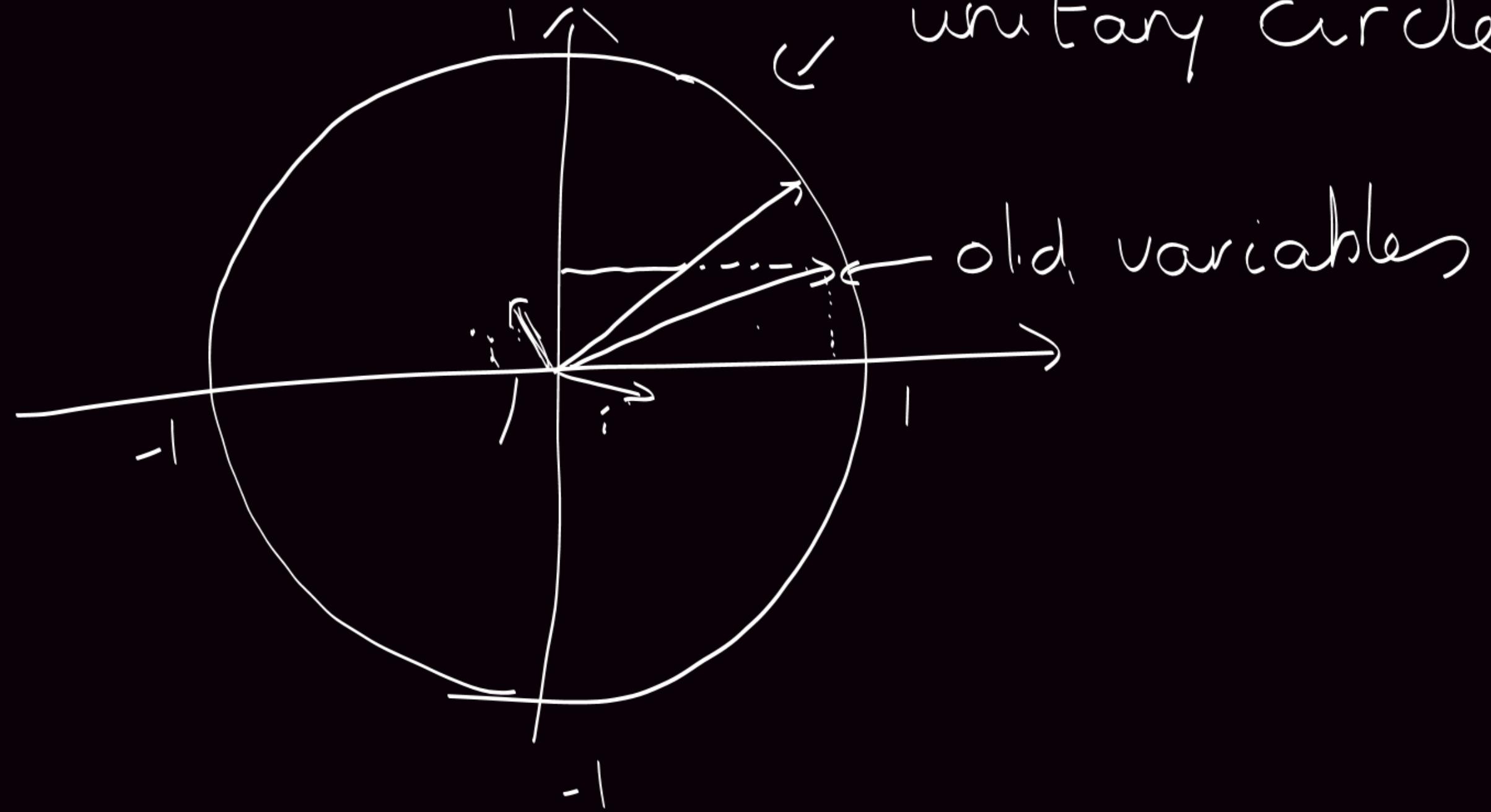
$$\begin{aligned}
 &= t(\sqrt{\alpha_R}) \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \\
 &= t(\lambda_R \alpha_R) \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \\
 &= \lambda_R t_{\alpha_R} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} = \underbrace{\lambda_R \alpha_R}_j \text{ element } j \text{ of } \alpha_R
 \end{aligned}$$

Thus

$$\text{Cor}(z_k, v_j) = \frac{\lambda_k \alpha_{kj}}{\sqrt{\lambda_k} \sqrt{\text{Var } v_j}} = (\lambda_k)^{\frac{\alpha_{kj}}{\sqrt{\text{Var } v_j}}}$$

$$\text{Var}(v_j) = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_{\cdot j})^2 \rightarrow \text{square of a norm}$$

correlation : scalar product



unitary circle.

old variables

$$\underline{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} x_{11} & \cdots & \cdots & x_{1p} \\ \vdots & & & \vdots \\ x_{n1} & & & x_{np} \end{pmatrix}$$

$P \gg n$

regression on
PCA

$$\begin{pmatrix} z_{11} & \cdots & z_{1q} \\ \vdots & & \vdots \\ z_{n1} & & z_{nq} \end{pmatrix} \quad q < n$$

J PCA

Trade off : Partial least
Square

Clustering

Aim: We have n individuals
can we determine several
groups among the n individuals?

Several methods

↳ all are based on the notion

of distance

↳ several distances!

P

let X be the space of the observations.

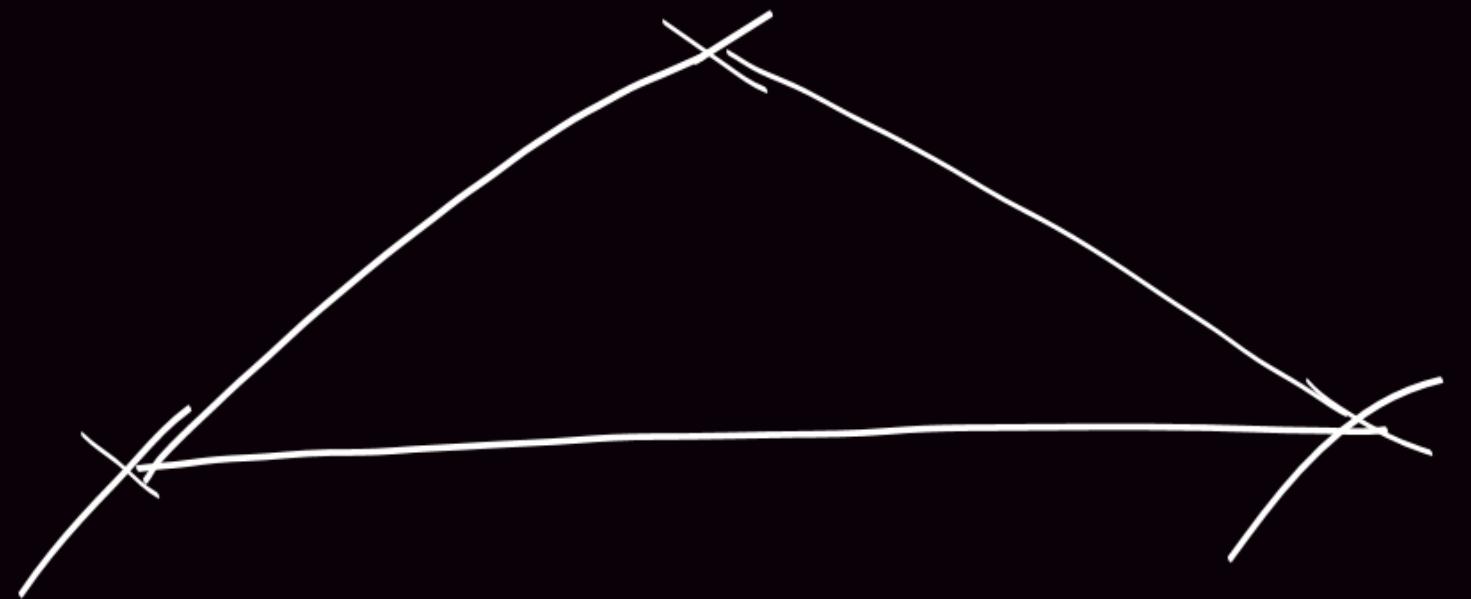
A distance d on X is a application

from $X \times X$ to \mathbb{R}^+ such that:

$$-\forall x, y \in X \quad d(x, y) = d(y, x)$$

$$-\forall x, y \in X, \quad d(x, y) = 0 \Leftrightarrow x = y$$

$$-\forall x, y, z \in X, \quad d(x, y) \leq d(x, z) + d(z, y)$$



$$X = \mathbb{R}^P$$

euclidian distance:

$$\forall x, y \in \mathbb{R}^P, d_2(x, y) = \sum_{i=1}^P (x_i - y_i)^2$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_P \end{pmatrix}, y = \begin{pmatrix} y_1 \\ \vdots \\ y_P \end{pmatrix}$$

ℓ_1 -distance

$$d_1(x, y) = \sum_{i=1}^p |x_i - y_i|$$

ℓ_∞ -distance

$$d_\infty(x, y) = \max_{i \in \{1, \dots, p\}} \{|x_i - y_i|\}$$

Exercise:

We take $P = 2$

Define $B(0, 1) := \{x(x_1, x_2) \in \mathbb{R}^2 / d(0, x) \leq 1\}$

$$\text{for } d = \begin{cases} d_1 \\ d_2 \\ d_\omega \end{cases}$$

d_2 :

$\mathcal{B}(O; 1)$

center

radius

$$d_2(O, x) \leq 1 \Leftrightarrow$$

$$\Leftrightarrow (x_1 - 0)^2 + (x_2 - 0)^2 \leq 1$$

$$\Leftrightarrow x_1^2 + x_2^2 \leq 1$$

$$x_1^2 + x_2^2 = 1$$

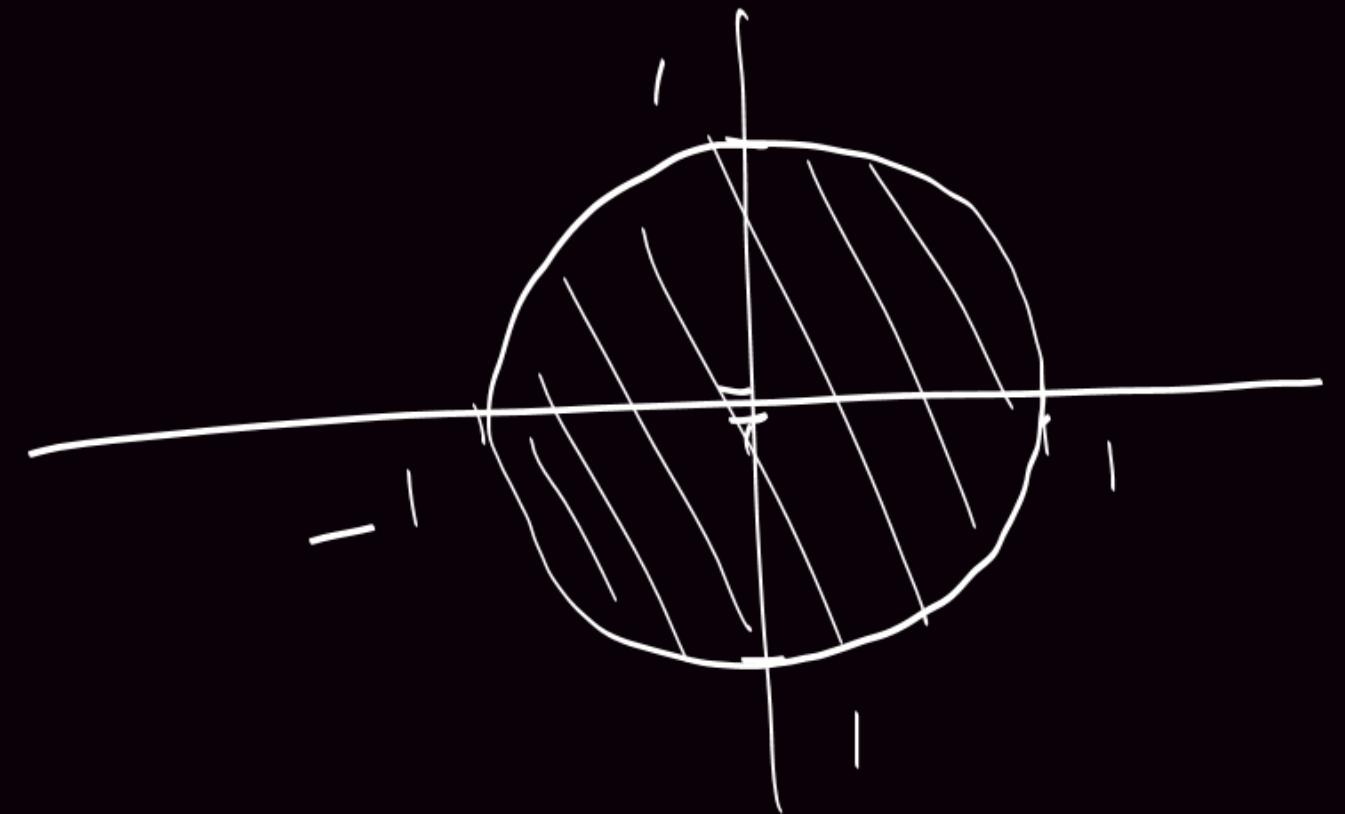
↳ equation of the circle

whose center is O and

radius 1

$$(x_1 - 2)^2 + (x_2 - 4)^2 = 2$$

Rk: $(x_1 - 2)^2 + (x_2 - 4)^2 = 2$
→ circle whose center is $(2, 4)$
and radius $\sqrt{2}$

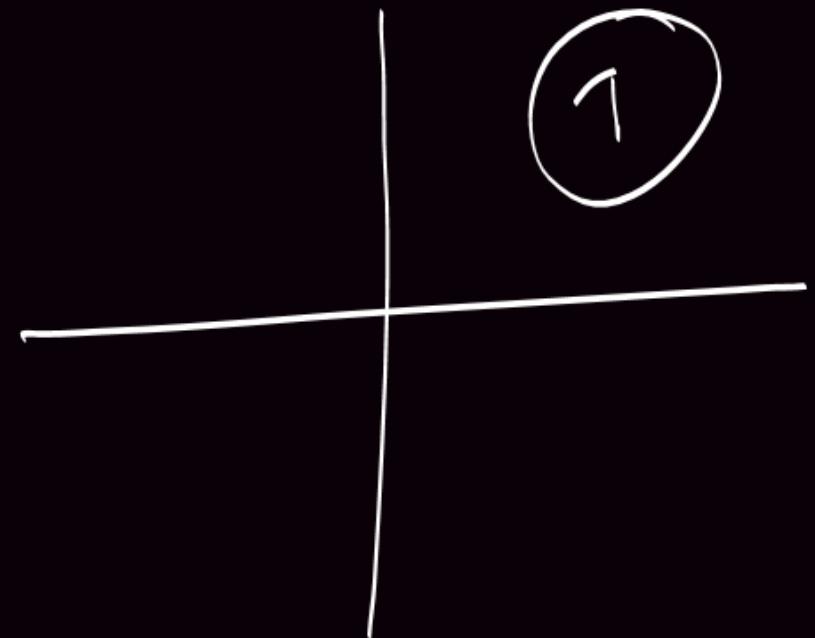


euclidian distance

$\frac{d_1}{x} \cdot$
 $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ such that $|x_1 - 0| + |x_2 - 0| < 1$
 $\Rightarrow |x_1| + |x_2| < 1$

1st case: $x_1 \geq 0, x_2 \geq 0$

$$\Rightarrow x_1 + x_2 < 1$$



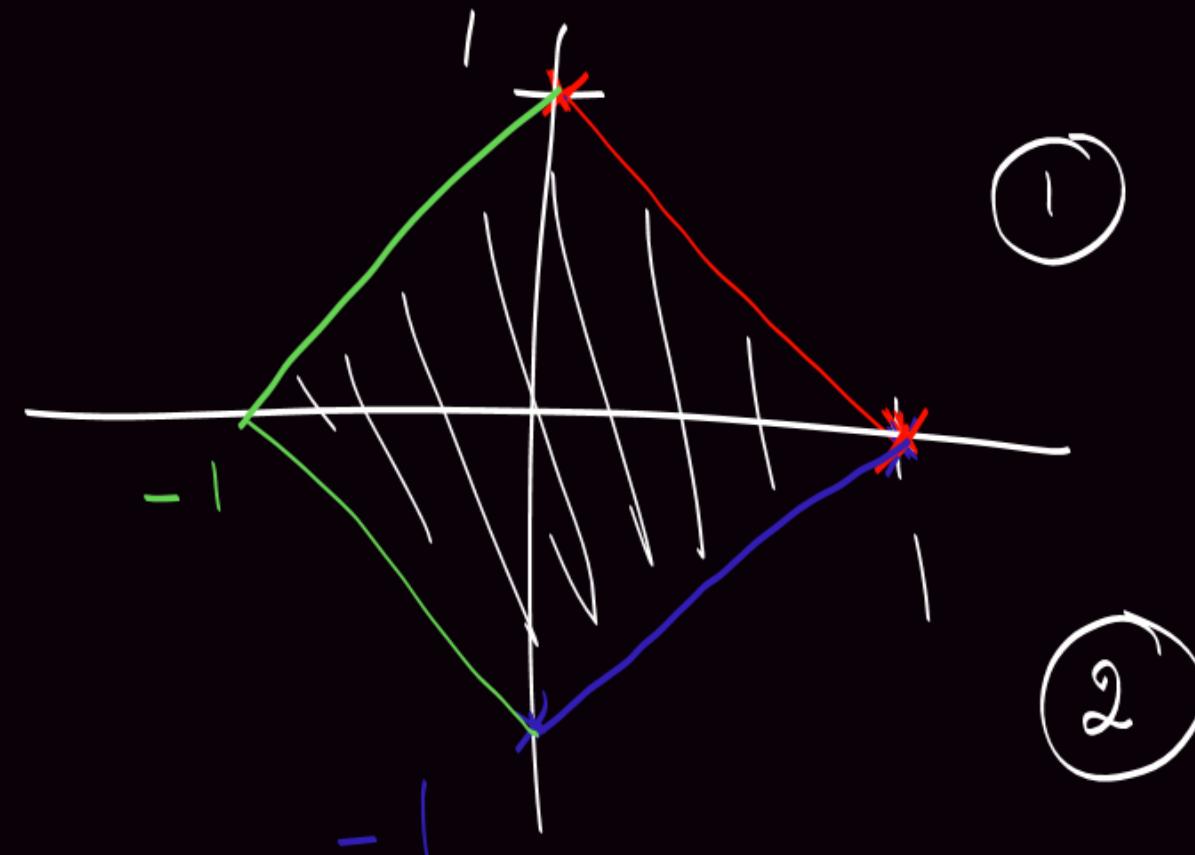
We solve $x_1 + x_2 = 1$

$$\Rightarrow x_2 = 1 - x_1$$

2nd case: $x_1 \geq 0$ and $x_2 \leq 0$

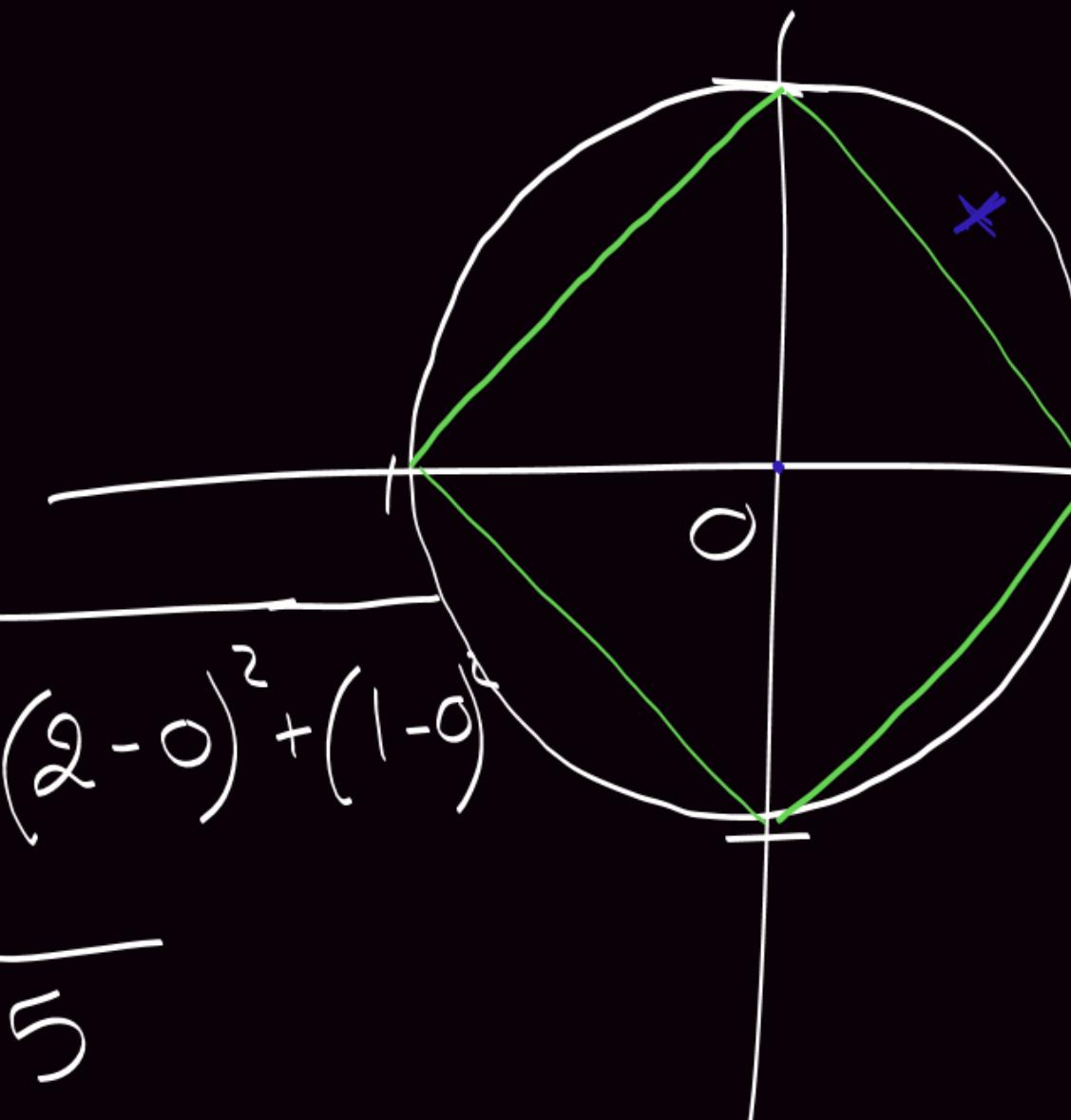
we solve: $|x_1| + |x_2| = 1$

$$x_1 + (-x_2) = 1$$



Rb:

$$d_2(O, x) = \sqrt{(2-0)^2 + (1-0)^2} \\ = \sqrt{5}$$



$$x = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

$$d_1(O, x) = \sqrt{|2-0| + |1-0|} \\ = 3$$

$$\frac{d}{d_\infty}(o, x) = \max \{ |x_1 - o|, |x_2 - o| \}$$
$$= \max \{ |x_1|, |x_2| \}$$
$$\leq 1$$

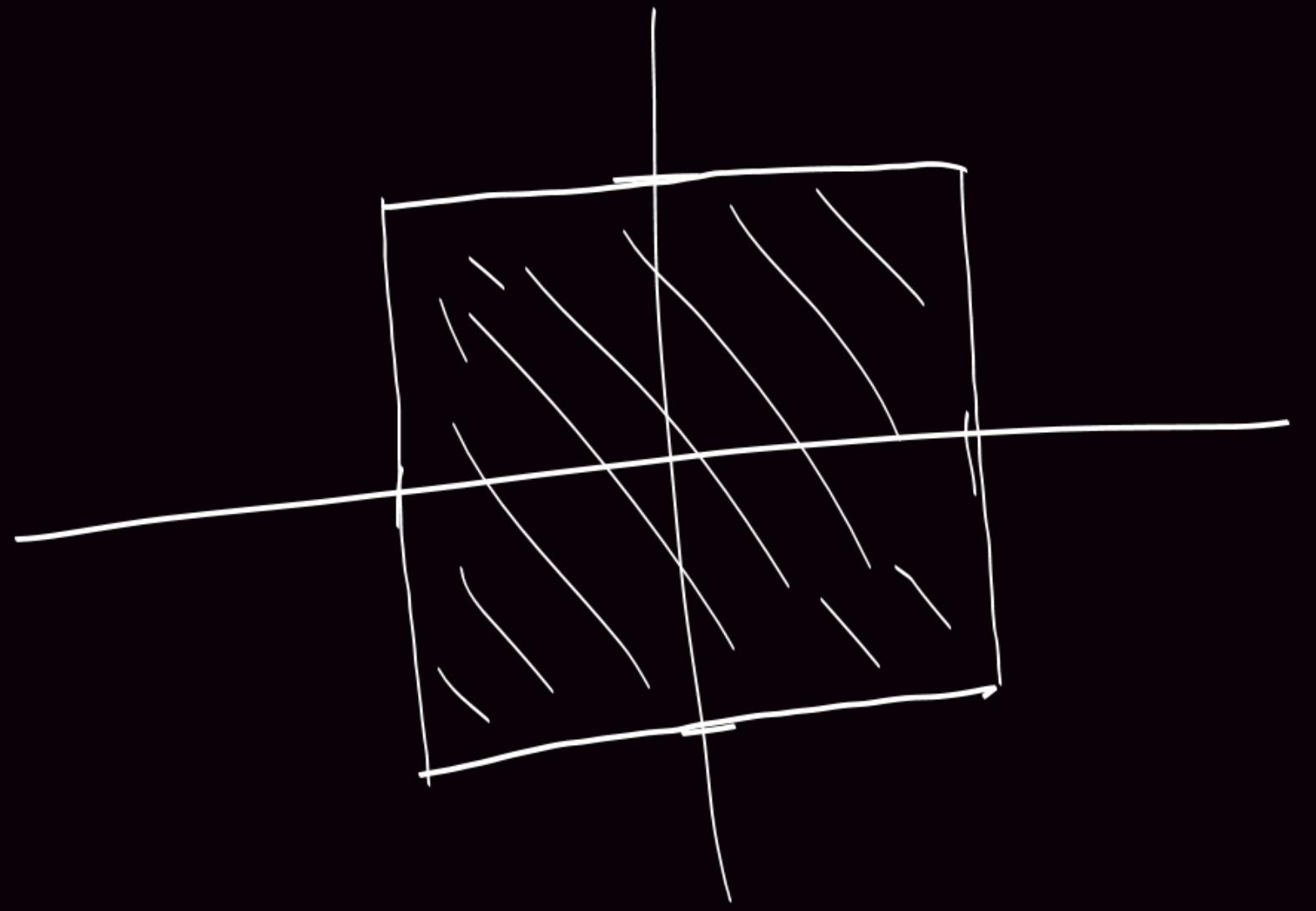
$$\max\{|x_1|, |x_2|\} = 1$$

$$\Leftrightarrow |x_1| = 1 \text{ and } |x_2| \leq 1$$

or

$$|x_1| \leq 1 \text{ and } |x_2| = 1$$

$$\Leftrightarrow \begin{cases} (x_1 = 1 \text{ or } x_1 = -1) \text{ and } (-1 \leq x_2 \leq 1) \\ (-1 < x_1 \leq 1) \text{ and } (x_2 = -1 \text{ or } x_2 = 1) \end{cases}$$



- × Ascending Hierarchical Classification
- × K-means

Ascending Classification

This is a recursive method
that depends on two distances:
- individual distance
- group distance

→ individual distance is used
to compute distance between
elementary individuals.

But as we are making groups
we need to compute distances
between groups.
→ group distance

- * group distance
- * Ward : let A and B be two groups

let n_A : the number of elementary individuals in A
and n_B : the same for B

The G_A the mean point of A

$$d_W(A, B) = \frac{n_A \times n_B}{n_A + n_B} \times d(G_A, G_B)$$

individual
distance

* minimal Link:

$$d_{mi}(A, B) = \min_{\substack{x \in A \\ y \in B}} d(x, y)$$

* maximal Link:

$$d_{ma}(A, B) = \max_{\substack{x \in A \\ y \in B}} d(x, y)$$

✓ average Link

$$d_{av}(A, B) = d(G_A, G_B)$$

description of the method

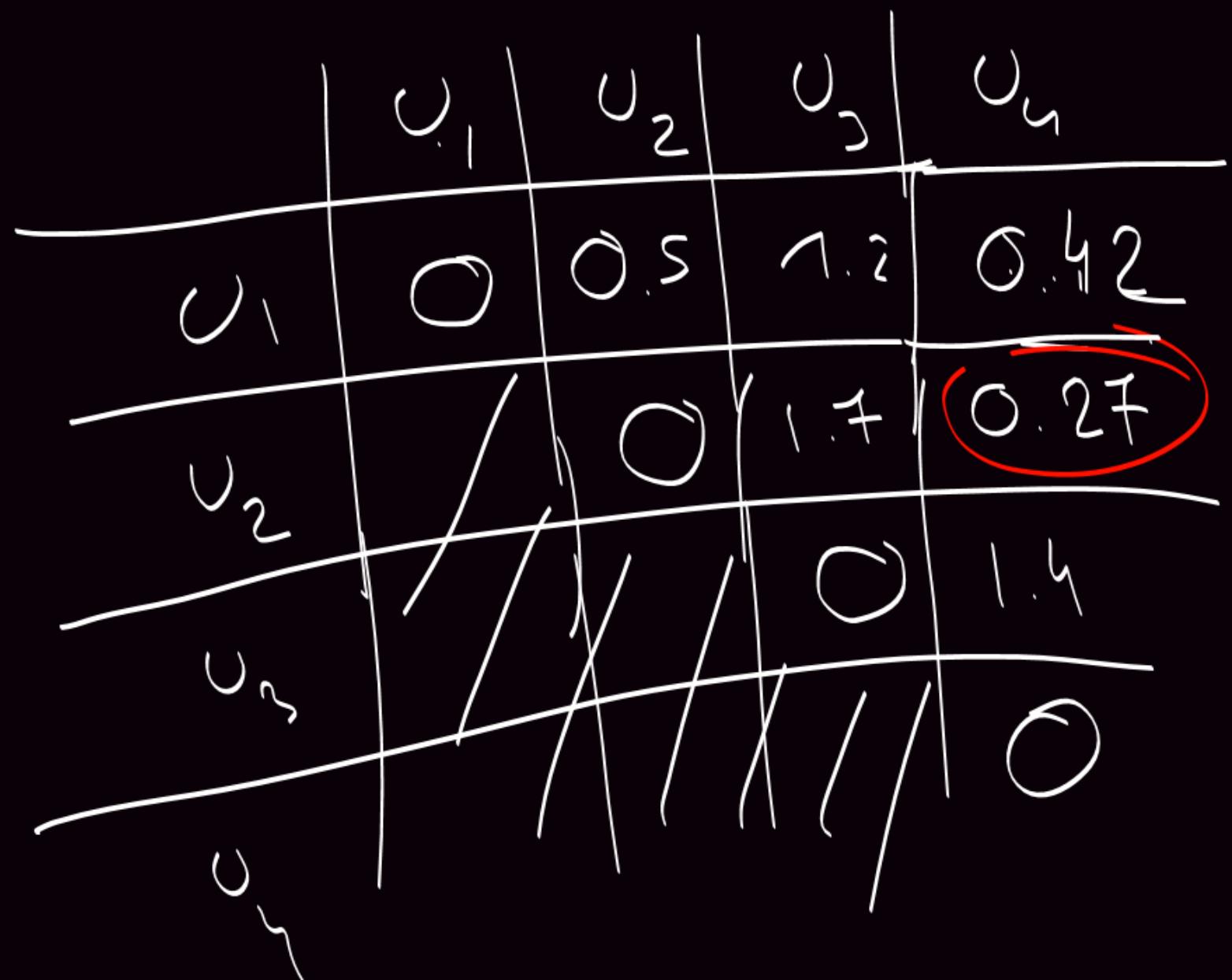
Step 1: we have n groups, each group contains 1 individual.

We compute all the possible distance between two individuals

We look at the two individuals associated to the smallest computed distance, not located on the diagonal!

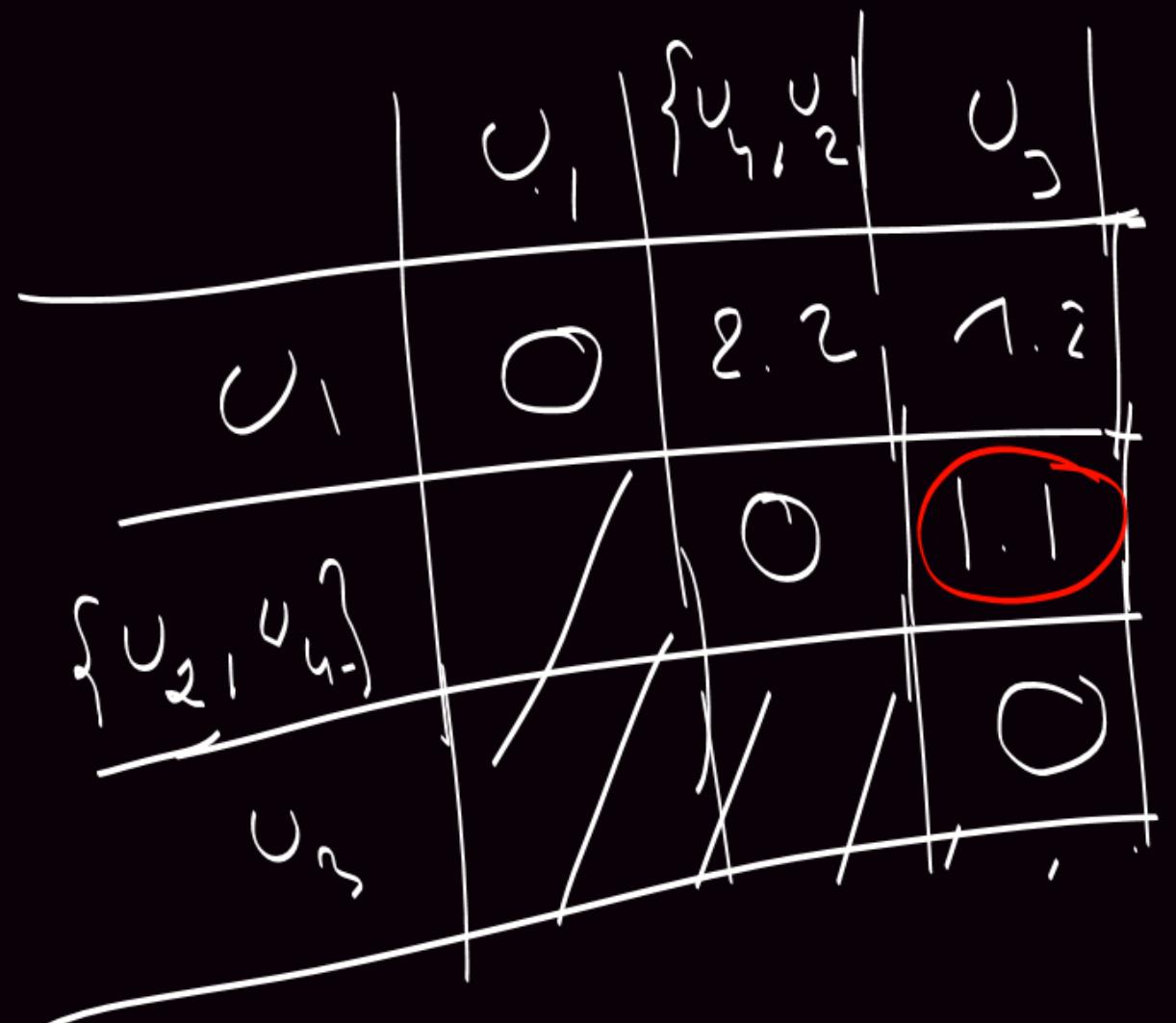
→ We put those two individuals in a group
→ at the end of step 1, $(n-2)$ groups with 1 individual, and 1 group with 2.

4 individuals.



→ we group U_4 and
 U_2

4 individuals



\Rightarrow we group
 $\{U_2, U_4\}$ with U_3

Step 2:

We compute the distance
between the $(n-1)$ groups.

→ smallest distance

→ group

R_k:

Total variance with n points:

$$\frac{1}{n} \sum_{i=1}^n d^2(e_i, G)$$

||

$$\sqrt{\sum}$$

e_i : individual i

G : mean point

$$= \begin{pmatrix} x_{1,1} \\ \vdots \\ x_{1,p} \end{pmatrix}$$

Suppose that we have K groups
 G_i : the mean points of group i

Variance between groups: V_B

$$V_B = \frac{1}{n} \sum_{i=1}^K d^2(G, G_i) n_i$$

n individuals
inside group i

Variance inside groups : V_w

$$V_w = \sum_{i=1}^K \frac{n_i}{n} \cdot \frac{1}{n_i} \sum_{j=1}^{n_i} d^2(e_{ij}, g_i)$$

where $e_{i,1}, \dots, e_{i,n_i}$
denote individuals
inside group i

$$\sqrt{T} = \sqrt{B} + \frac{\sqrt{W}}{\gg 0}$$

\Downarrow

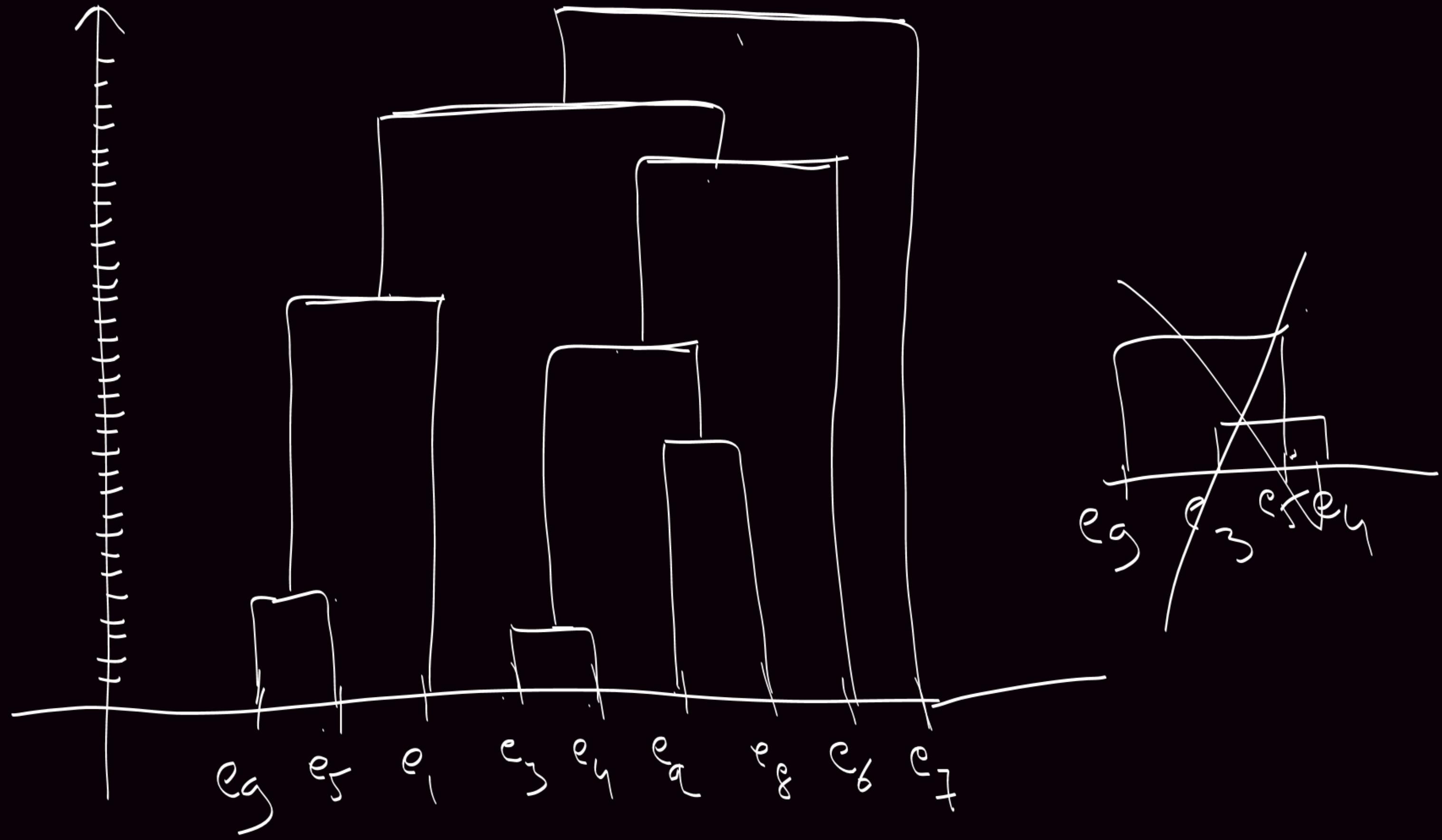
$$\sqrt{T} = \sqrt{B} + \frac{\sqrt{W}}{0}$$

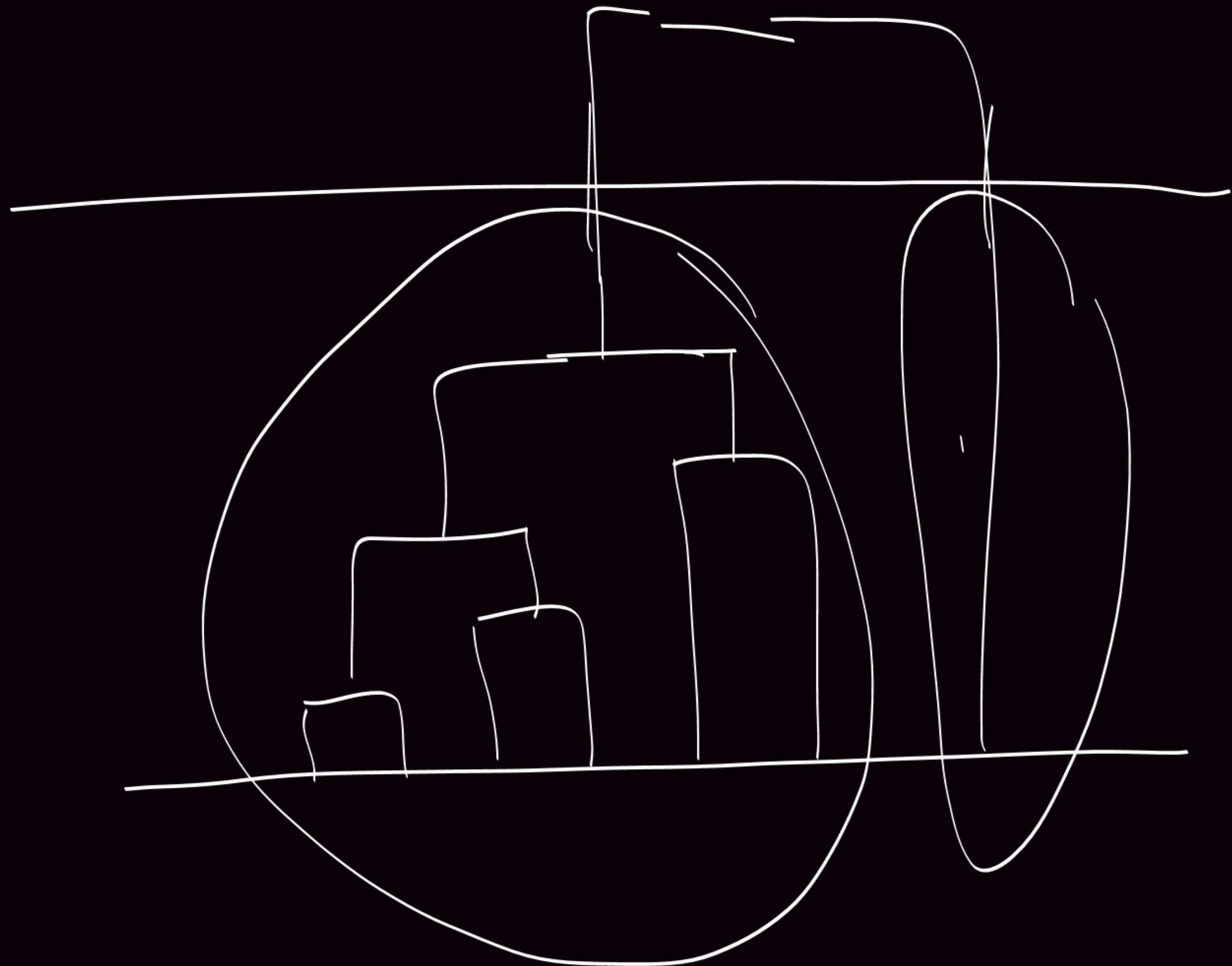
\Downarrow

$$\sqrt{T} = \sqrt{B} + \sqrt{W}$$

We stop the ascending method
when all the individuals are
together!

drawing of the ascending method:
cladogram





WE ARE
BACK