Data we have

Train-test-split (

$F_1$  $F_2$  $F_3$  $y$

train 80%

| 1 |
| 2 |
| 3 |
| ⋮ |

$X_{train}$

$y_{train}$  $\hat{y}_{hat\_train}$  $\Rightarrow$

$X_{test}$

$n = 100$

$[y_{test} \Leftarrow [\hat{y}_{hat\_test}]$  ✓
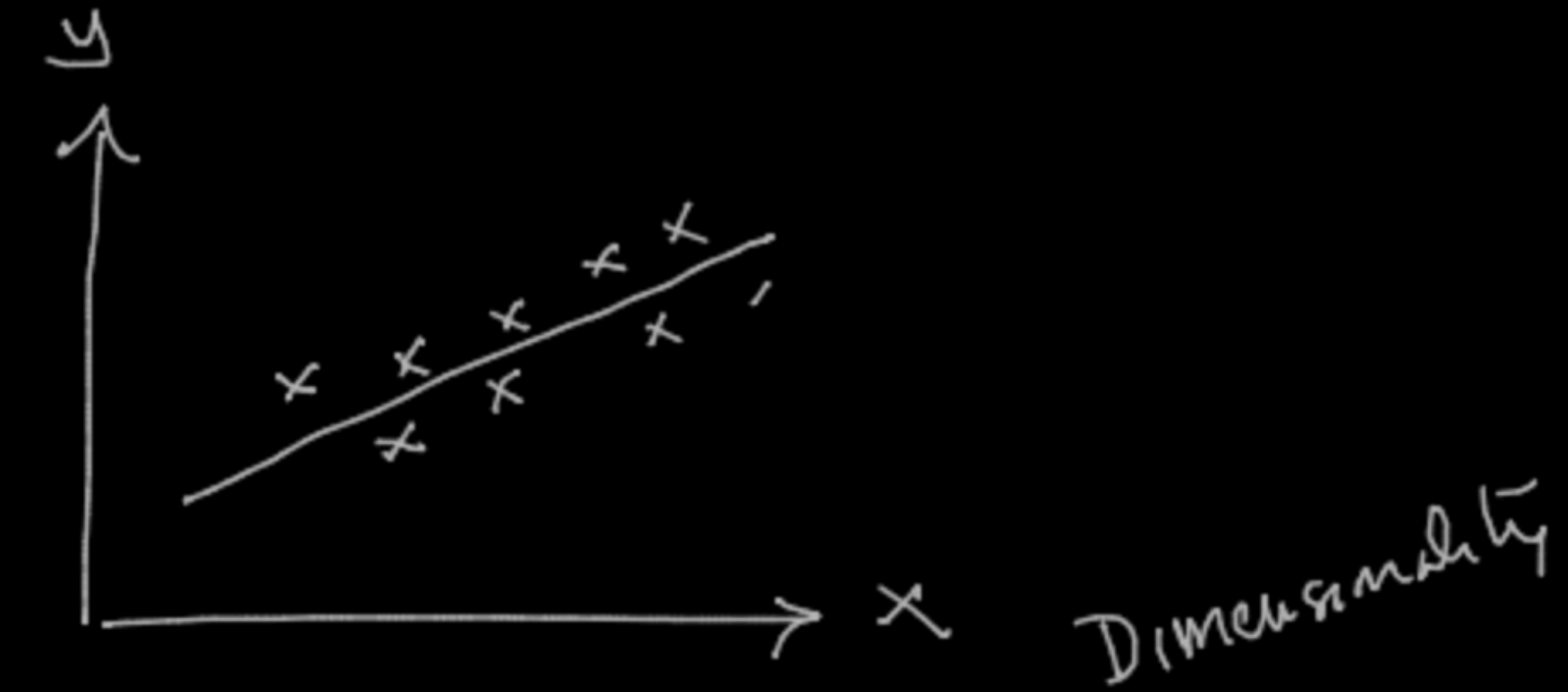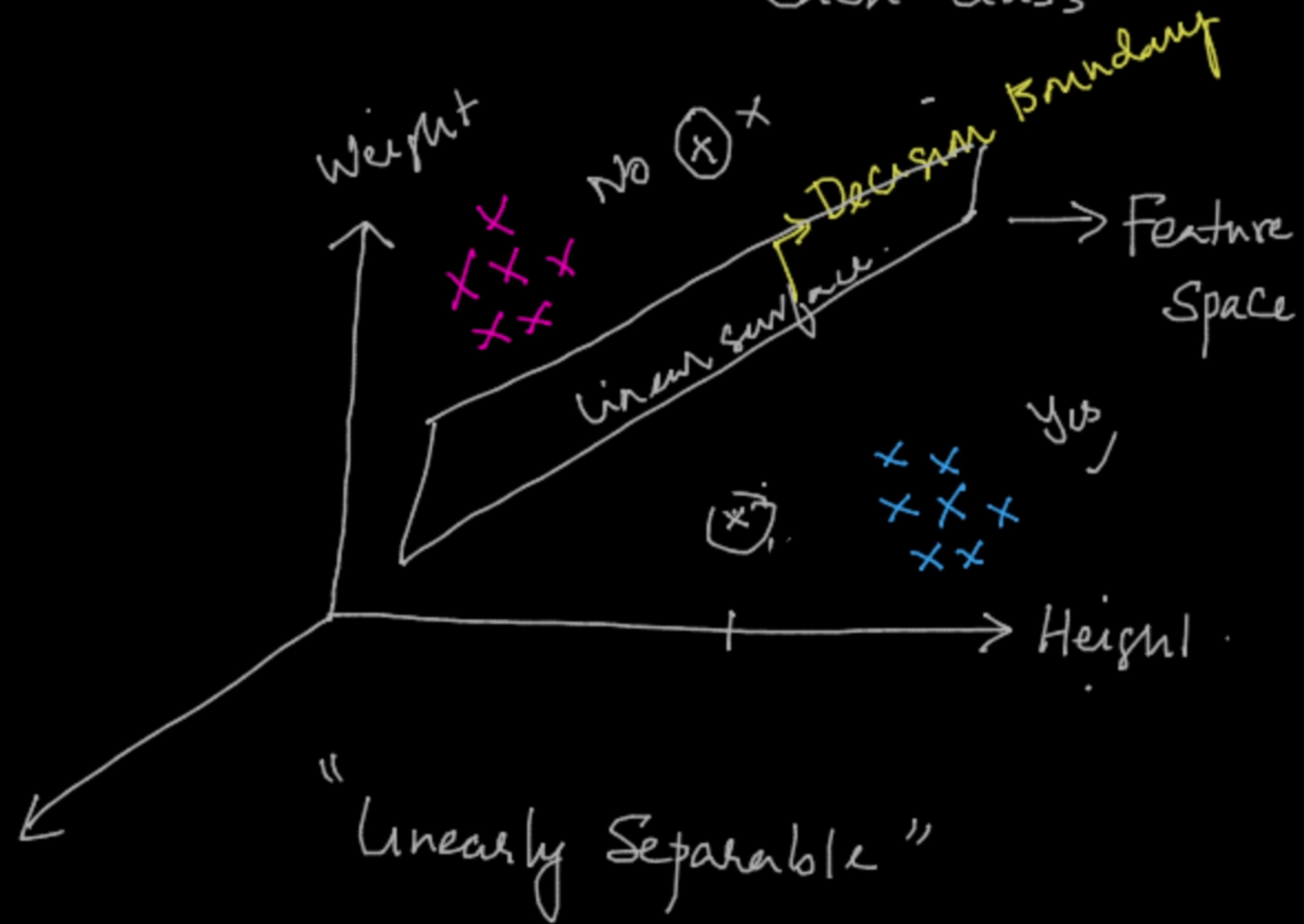
↑ Accuracy

model = log reg  $\leftarrow$ 80% $\rightarrow$

model . fit ($X_{train}$, $y_{train}$)

model . predict ($X_{train}$) $\Rightarrow$

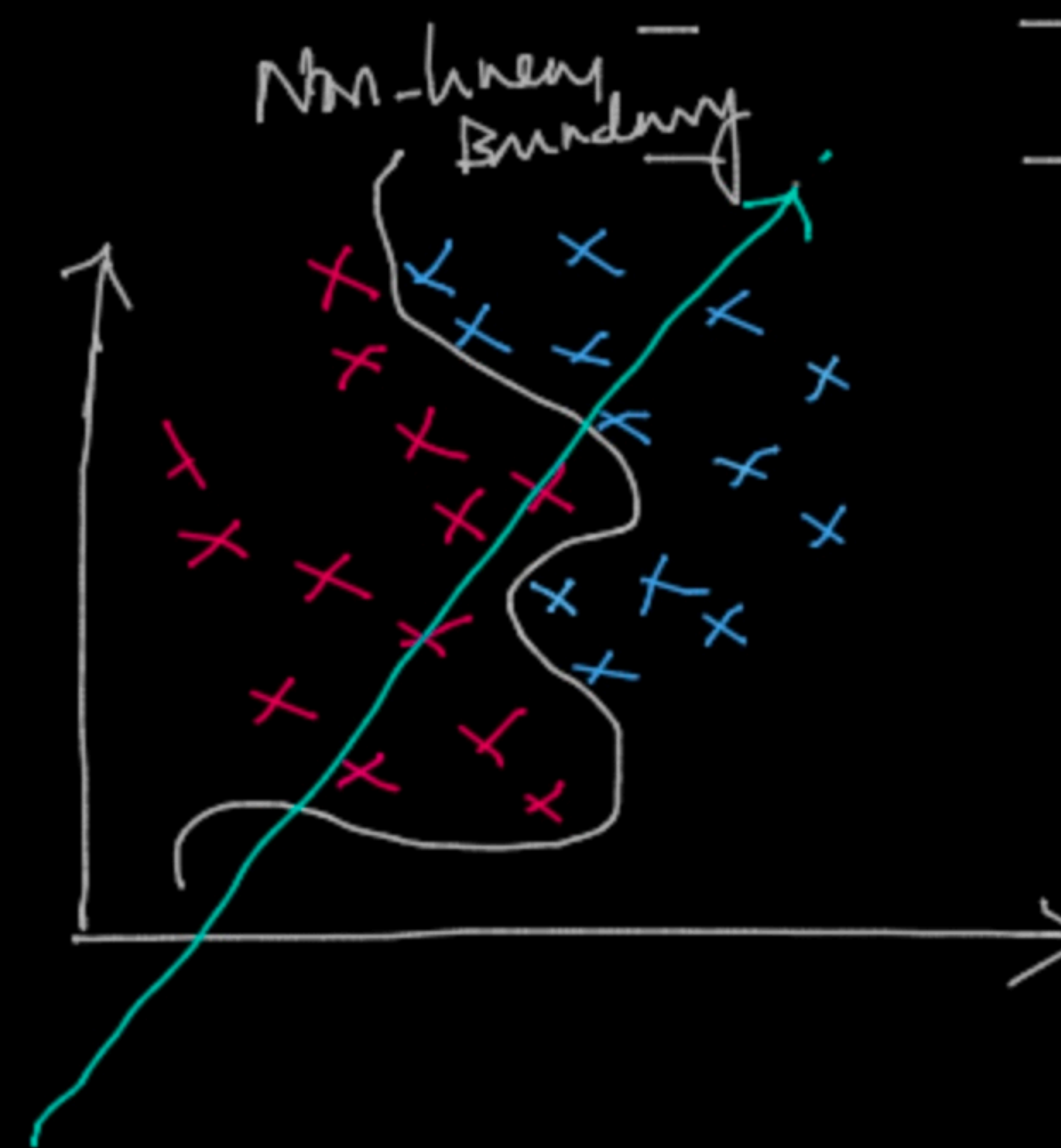model . predict ($X_{test}$) $\Rightarrow$

Unseen ✓

80% ── 20%

Talking: Geethika

# logistic Regression — Linear Model

Classification → " Partitioning the feature space into pure regions assigned to each class "



Weight

No $\otimes$ $\times$

→ Decision Boundary

$\times \times$
$\times \times \times$
$\times \times$

Linear surface

→ Feature Space

Yes

$\otimes$

$\times \times$
$\times \times \times$
$\times \times$

→ Height

" Linearly Separable "

y

$\times \times \times \times$
$\times \times \times$
$\times$

→ X

Dimensionality

| Chd | Wc | BMI | flight $x_1$ | Weight $x_2$ | Select |
|-----|-----|-----|-----|-----|-----|
|  |  |  | — | — | ✓ |
|  |  |  | — | — | ✗ |

Non-linear Boundary

# Clustering

— unsupervised ✓

— No `y` Value .

— Clusters of similar Data.

$$y - \hat{y} = Rgre$$

$$y \; \hat{y} \Rightarrow Accuracy.$$

Good cluster?



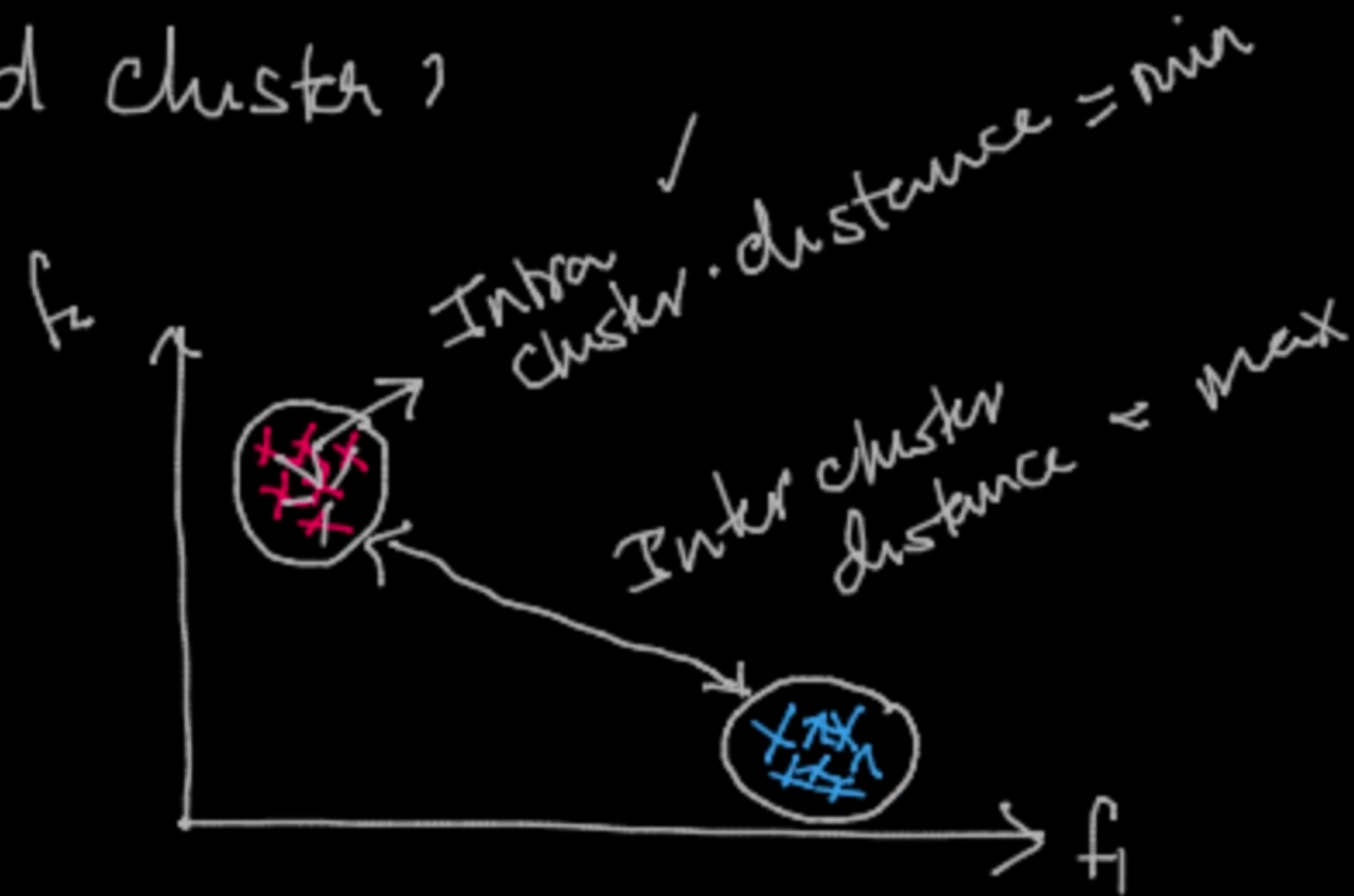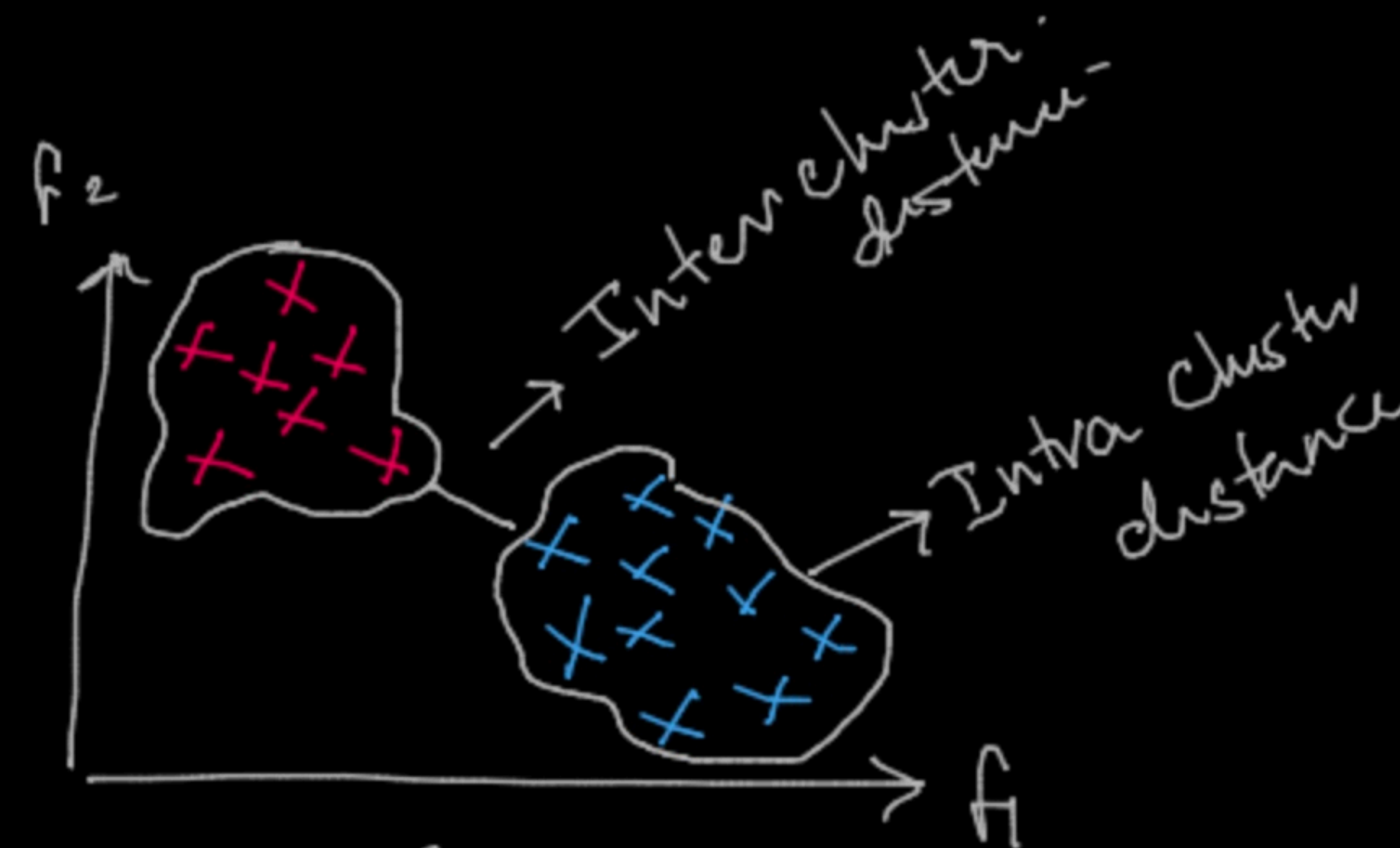Intra cluster distance = min
Inter cluster distance = max

fig 1

fig 2

Inter cluster distance
Intra cluster distance

Supervised
Regression
Classification

unsupervised
✓ Clustering

$$Dunn\ Index \Rightarrow \left[\dfrac{Inter}{Inter}\right]$$

✓ 1 Minimum Intra cluster distance (WCSS)

— Within Cluster sum of squared distances.

✓ 2 Maximum Inter cluster distance

Weight

$P_3$ x $(Ht, Wt, BMI)$

$P_2$
x

$P_1$

BMI

→ Height

Distance ⟺ Similarity

Distance $\propto \dfrac{1}{\text{Similarity}}$

|       | $f_1$ | $f_2$ | $f_3$ | $f_4$ |       |
|-------|-------|-------|-------|-------|-------|
| $P_1$ | Nr | Cat | Nr | Nv | } Mixed Data |
| $P_2$ | N₂ | Cat₁ | Nr | N₂ | |
|       | Nr | Nr | Nr | Nv | } Numeric |
|       | Wr | Nr | N₂ | N₂ | |
|       | Cat | Cat | Cat | Cat | } only categorical |
|       | Cat | Cat | Cat | Cat | |

Numeric

1. Euclidean Distance ✓
2. Manhattan distance ✓
3. Minkowski distance ✓
4. Mahalanobis dist

Categorical

1. Binary Euclidean
2. Simple Matching Co-efficient
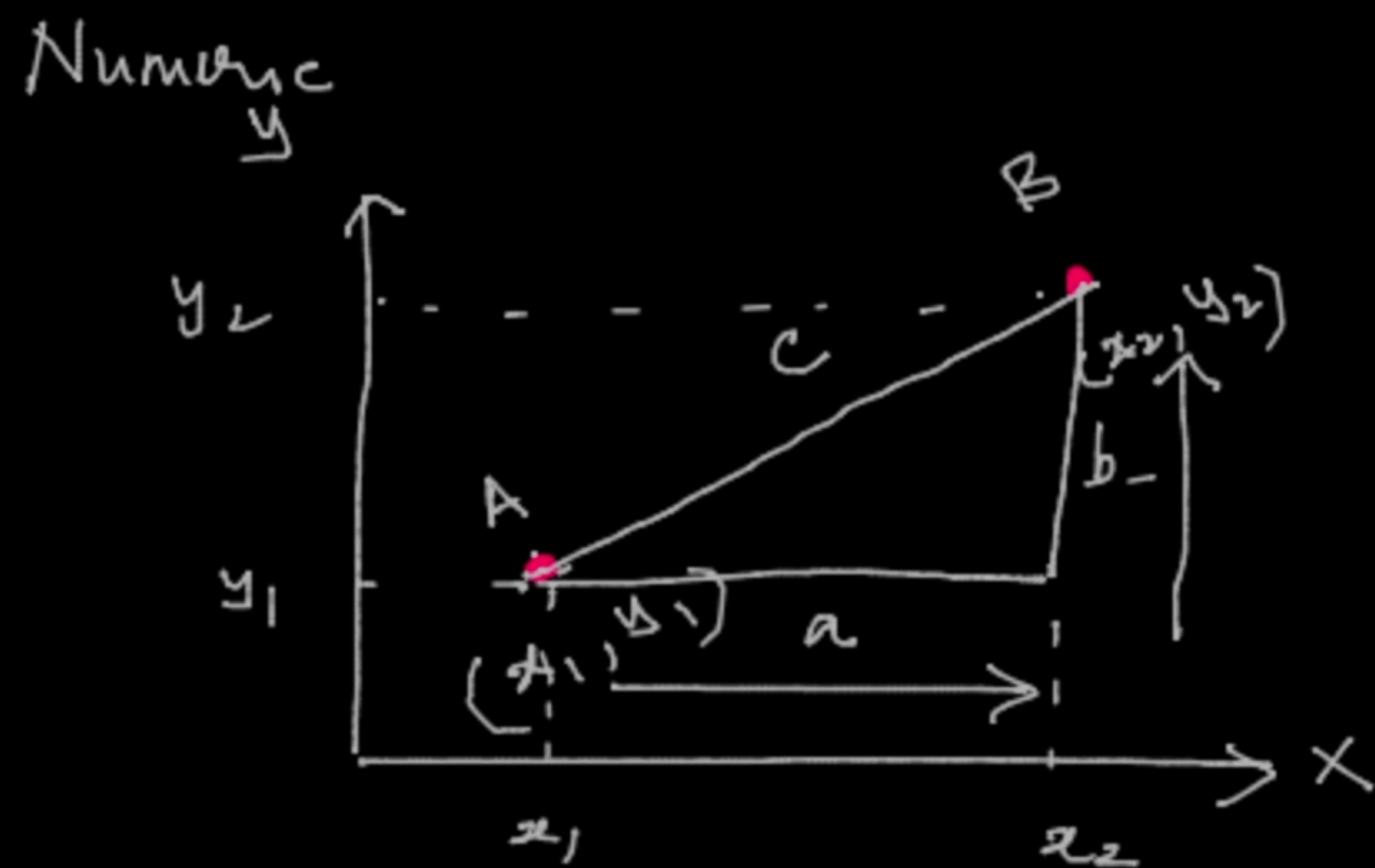3. Jacquard's dist

Mixed

1. Grower's dissimilarity Index

Numeric



**1. Euclidean Distance**

$$c = \sqrt{a^2 + b^2 + c^2 + \cdots}$$

$$\sqrt{(x_2-x_1)^2 + (y_2-y_1)^2 + (z_2-z_1)^2 + \cdots}$$

'As the crow flies'

**2. Manhattan distance**

$$c = a + b$$

"Taxi Distance"

**3. Minkowski distance**

$$c = \left[ (x_2-x_1)^p + (y_2-y_1)^p + \cdots \right]^{1/p}$$

$P = 1$ ; Manhattan
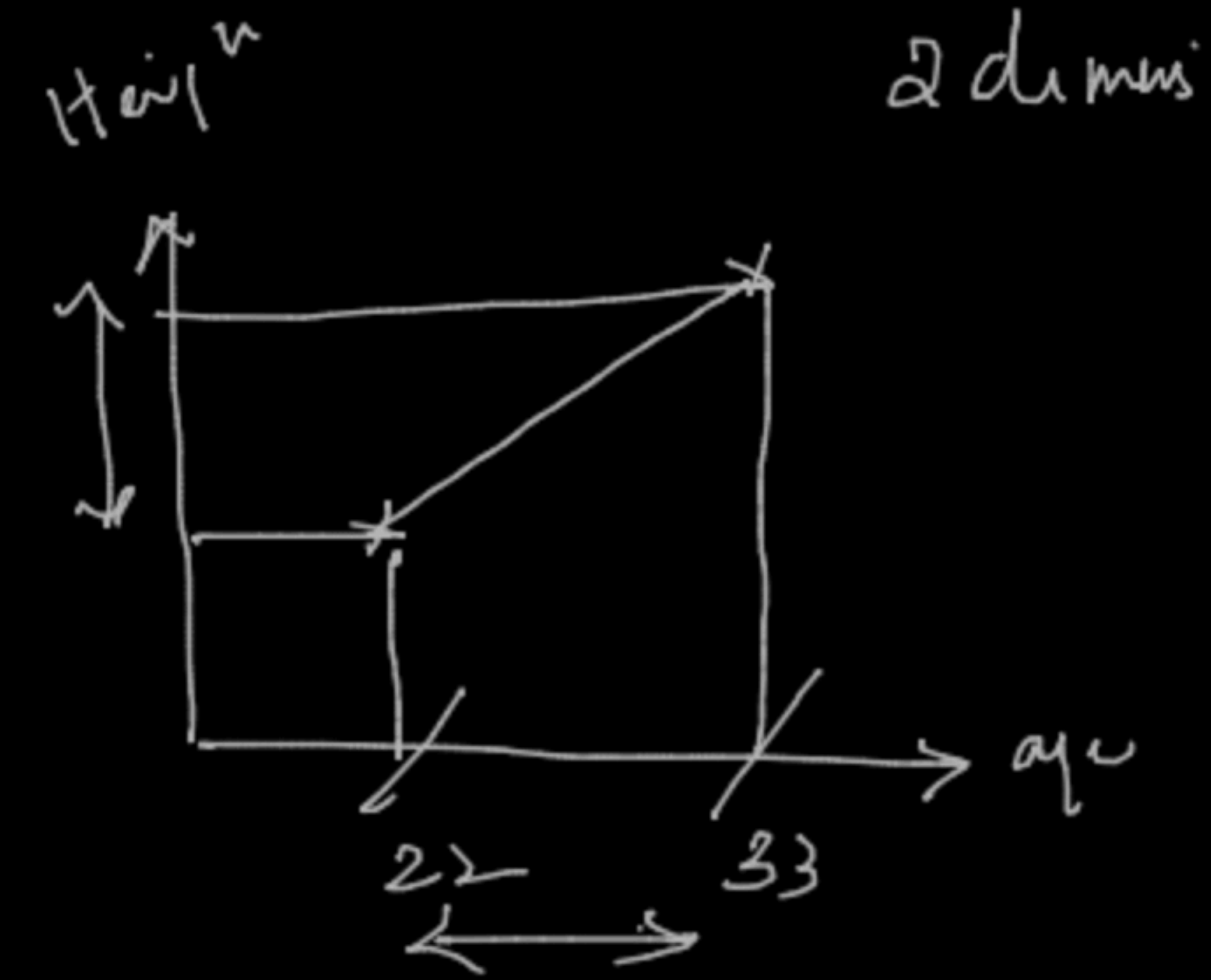
$$c = (x_2-x_1) + (y_2-y_1)^1$$
$$a + b$$

$P = 2$ ; Euclidean

$$c = \left[ (x_2-x_1)^2 + (y_2-y_1)^2 \right]^{1/2}$$

$$=$$

Height$^n$                    2 dimns

Talking: Geethika

age

22        32

$\longleftrightarrow$

11

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

22          33

$\longleftrightarrow$

age

| age | Hight |
|-----|-------|

P$_1$    22                    P$_1 \rightarrow$ 22    150

P$_2$    33                    P$_2 \rightarrow$ 33    167

11 yrs

$\longleftarrow$ Numeric $\longrightarrow$

Age   Income   height   Wt

$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix}$

Euclidean

Manhattan

$\longleftarrow$ Categorical $\longrightarrow$

Married      Manager  Smoker  gender

$\begin{bmatrix} P_1 \\ P_2 \end{bmatrix}$

$P_3$

Binary Euclidean

Single matching coeff

Jacquard'

| Id |
|----|
| Age |

$\begin{array}{c|c} P_1 & - \\ P_2 & - \\ P_3 & - \end{array}$

$\longleftarrow$ Mixed  Data $\longrightarrow$

Age  Income   Gender  Married  height   Wt·   Smoker

$P_1$

$P_2$

$P_3$

Gower's dissimilarity Index

# K-Means Algorithm

$\hookrightarrow$ No. of clusters.

— "Hyperparameter"

We decide what
Value to specify

$f_2$

$n_1$ $n_2$



| Hght | Wt |
|------|-----|
| $f_1$ | $f_2$ |

$x_1(B) \rightarrow \quad \left( \dfrac{22}{-} \quad \dfrac{150}{-} \right) \uparrow$

$x_2(0) \rightarrow \quad \left( \begin{array}{cc} - & - \\ 3 \end{array} \right)$

$x_3(0) \rightarrow \quad \begin{array}{cc} - & - \\ - & - \end{array}$

$x_0 \cdot \quad \left( \begin{array}{cc} - & - \\ - & - \end{array} \right)$

$\underline{\quad} \text{ data} \rightarrow \boxed{\text{Model}}$

$x_{1i}$

$\rightarrow x_{11} \quad (\underbrace{\longrightarrow} )$

$x_{12}$

$k$ - clusters

$k \rightarrow$ No. of clusters

Task $\rightarrow$ Find 'k' centroids s.t the WCSS is minimum

$$WCSS = \sum_{i=1}^{n_1} (x_{1i} - C_1)^2 + \sum_{i=1}^{n_2} (x_{2i} - C_2)^2 + \sum_{i=1}^{n_3} (x_{3i} - C_3)^2 + \cdots \sum_{i=1}^{n_k} (x_{ki} - C_k)^2$$

$\Downarrow$ $WCSS_1$ $\Downarrow$

$$WCSS \downarrow = \sum_{j=1}^{k} \sum_{i=1}^{n_j} (x_{ji} - C_j)^2 \Rightarrow$$

$\Rightarrow$ 'NP Hard Problems'

Approximation $\rightarrow$ Lloyd's Approximation ✓

$= WCSS -$

Talking: Geethika