Machine learning

Supervised                    unsupervised

→ unlabelled dat

— Labelled data            → No target

— Target is available      No prediction

y — continuous        y categorised

Regression            Classfn

__Metrics__            — MPG → HP, SP, WT, VOL        → Checked → Area, income, gender, time, usage

✓ MSE                    Cont                              Cat                        __Metrics__

✓ RMSE              ✓ Linear Rg          Logistic Regress ✓        Accuracy ⇒

✓ R²                    D tree Rg             D tree clf ⟶           Precisn ⇒

                        Random Rg            KNN clf                Recall →

                        XGBoost Rg           XGBoost clf            F1 Score →

                                                                    Auc Roc →

# Regression

Model

| $x_1$ | $x_2$ | $x_3$ | $y$ | $\hat{y}$ | $(\hat{y}-y)^2$ |
|---|---|---|---|---|---|
| | | | (35)→ | (29) | +6 |
| | | | 43 | 47 | −4 |
| | | | 76 | 79 | −3 |

actual

Train

MSE

| 65 | 61 | 4 |
|---|---|---|
| 73 | 70 | 3 |
| 57 | 54 | 3 |

→ unseen data

MSE

$$y = mx + c$$

$$\hat{y} = \beta_0 + \beta_1 x$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \quad + \beta_m x_m$$

$\beta's$ → Model Parameters

# Classfn

| $x_1$ | $x_2$ | $x_3$ | $y$ | $\hat{y}$ |
|---|---|---|---|---|
| | | | 0 | 0 |
| | | | 1 | 0 |
| | | | 1 | 1 |
| | | | 0 | 0 |
| | | | 1 | 0 |

Train

Test

→ unseen

→ ovaflthg

Acc,

| TP | FP |
|---|---|
| FN | TN |

Acc $\left(\dfrac{3}{5}\right)$ ←

Clustring →

unstandardized
↑
WCSS ✓ _ s
Silhouette scur. ✓ (−1 to +1)
Dunn Index ✗ standardized

Similarity

| | Ht | Wt | Ac | mane | mu |
|---|---|---|---|---|---|
| ✓ P₁ | — | — | | | |
| ✓ P₂ | — | — | | | |
| ✓ P₃ | — | — | | | |
| ✓ P₄ | — | — | | | |
| ✓ P₅ | — | — | | | |

$\longrightarrow$ Feature space

Wt

liquid

P₃
P₁  P₄    P₂  P₅

$\longrightarrow$ Ht
tall & thin

short

distance $\propto \dfrac{1}{\text{similarity}}$

mane

+ 0.8

$\longrightarrow$ Ac

**Fig 1**

(intra cluster distm)

1 WCSS $\rightarrow$ As small as possible

2 Intercluster distance — As large as possible

similarity

+ve 0-1

**Fig 2**

$\rightarrow$ Silhuette score $\rightarrow$ -1 to +1

zero

−ve silhouette score

→ KMeans clustering

→ DBSCAN

$n_2$

K-Means Clustering

$\hookrightarrow$ No of clusters $\rightarrow$ 2

$n_1$

$C_2$

$C_1$

Wcss — Within Cluster Sum of squared distances

Task:

Find $k$ Centroids s.t the $\boxed{\text{Wcss}}$ is minimum

$$\sum_{\wedge=1}^{n_1} \left(x_i - c_1\right)^2 + \sum_{\wedge=1}^{n_2} \left(x_i - c_2\right)^2 + \quad + k \text{ Cluster}$$

$$\sum_{j=1}^{k} \sum_{\wedge=1}^{n_j} \left(x_i - c_j\right)^2 \longrightarrow \text{Wcss}$$

NP Hard problem.

— Lloyd's approximation

Step 1    Randomly choose k-points
from the dataset as centroids

Step 2    Find the distance betn every dp
and the k centroids

Step 3    Assign points to the closest centroids

Step 4    Recalculate the centroids

KMeans ++

Problems

1  Initialisation Sensitivity
— Final clusters depend on the choice
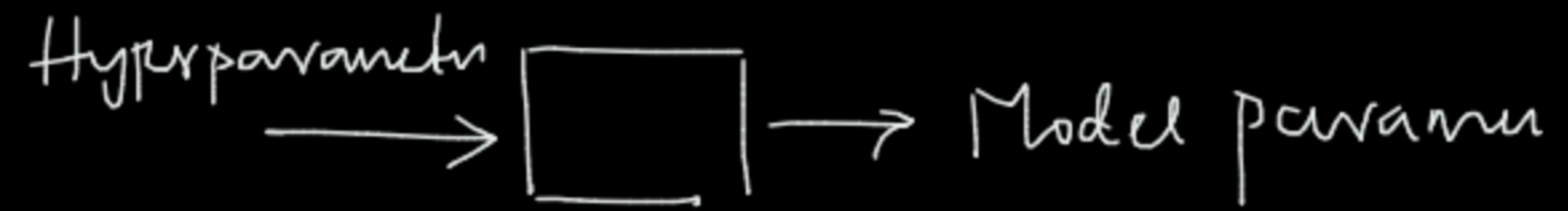of initial centroids

2  outliers are not handled properly

$(x_2, y_2)$

$(x_1, y_1)$

$(x_3, y_3)$

$C\left(\dfrac{x_1+x_2+x_3}{3}, \dfrac{y_1+y_2+y_3}{3}\right)$

Hyper parameters

— We supply to the model

Model parameters

— estimated by the model from the data
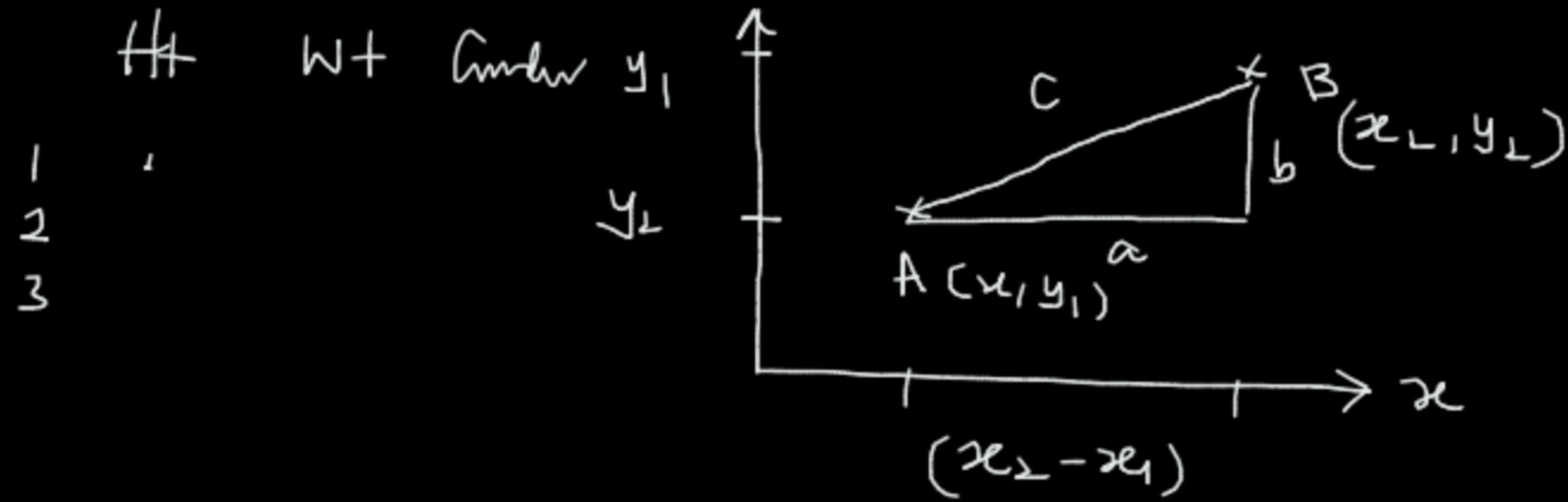
Hyperparameter $\longrightarrow$ ☐ $\longrightarrow$ Model parameter

$K \rightarrow$ Hyperparameter

$\rightarrow$ No. of clusters

$\longrightarrow$ Elbow Method

Athul  Ansu  Acsip  Not shqn



WCSS

$\nearrow$ Elbow pt

1    2    3    4    5    6

# Euclidean distance

Ht   Wt   Cmdw   $y_1$

1
2
3

$y_2$



C

$B$
$b$ $(x_2, y_2)$

$A (x_1, y_1)$   $a$

$(x_2 - x_1)$

$x$

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

$$= \int \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2} \, t$$

`As the crows flies`

## Manhattan Distance — taxi distance

$p = 2$

$p = 1$



B

$b$

A

$a$

$c = a + b$

$$= (x_2 - x_1) + (y_2 - y_1)$$

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

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

$p = 1$       $(x_2 - x_1) + (y_2 - y_1)$

## All Features are Nr

— Euclidean

— Manhattan

— Minkowski

— Mahalanobis

#### Ht  Wt  age

$P_1$

$P_2$

$P_3$

## Categorical

Binary Euclidean

Simple Matching Co-efficient

Jacquard's dist

#### Gndr  Job  Manner  Similar

$P_1$

$P_2$

## Mix of Nr & cal

Gower's dissimilarity index

#### Gndr  Age  income  % Married

$P_1$

$P_2$

$P_3$

|      | Age ✓ | Income ✓ |
|------|-------|----------|
| $P_1$ | $28x$ | $15L$ → |
| $P_2$ | $32x$ | $25L$ → |

$$\sqrt{(28-32)^2 + (15L - 25L)^2}$$

$$= \underline{16} + 10,00,000\ 000\ 000$$

Scaling

Standardize → $Z_x = \dfrac{x-\mu}{\sigma}$ → $-3$ & $+3$

Normalizing → $N_x = \dfrac{x - x_{min}}{x_{max} - x_{min}}$ → $0$ & $1$

| $\overline{Age}$ | $\overline{N_{Age}}$ |
|------|------|
| 27 |   |
| 32 |   |
| 45 |   |
| 20 → | 0 |
| 54 | 1 |

$$\dfrac{27-20}{54-20} = \dfrac{7}{34} \to 0$$

$$\dfrac{54-20}{54-20} = 1$$

$$\dfrac{20-20}{34} \to 0$$