

★ Feature Engineering (Handling missing value)

→ missing value occurs in dataset when some of the informations is not stored for a variables there are 3 mechanism.

① missing completely at Random (MCAR)

→ missing completely at Random is a type of missing data mechanism in which the probability of value being missing is unrelated to both the observed data and the missing data. in other words, if the data is MCAR the missing values are randomly distributed throughout the dataset, and there is no systematic reason for why they are missing.

for example:- in a survey about the prevalence of certain disease the missing data might be MCAR if the survey participants with missing value for certain questions were selected randomly and their missing responses are not related to their disease status or any other variables measured in the survey.

② missing At Random (MAR)

→ If one data is empty then another data is also empty
→ age of women X → salary of men → X

③ missing data not a random (MDAR)

→ Because of one fixed. It does not fixed another data.
→ Job Statistichication → income depended.

missing
value.
Random

① MCAR → no specific reason of missing value
→ Survey form

② MAR → data is missing due to one data is depend on another data
→ Systemic Relationship → income of men → age of women

③ MDAR → missing is not Random and is dependent on unobserved or unmeasured factors that are associate with the missing value

Import Seaborn as sns

★ checking missing value in titanic dataset

→ `df = sns.load_dataset('titanic')`

→ `df.isnull().sum()`

→ this fun is for find null value

Show All in One

True → missing value
False → no missing value

→ `sns.heatmap(df.isnull())`

using heatmap visual a missing value

missing value

★ Handling missing value By deleting row

→ `df.dropna().shape`

delete all rows having a nan value

values in tuple form

★ Handling missing value By deleting columns

→ `df.dropna(axis=1)`

→ `sns.distplot(df['age'])`

new column for mean.

→ `df['Age-mean'] = df['age'].fillna(df['age'].mean())`

→ `df['Age-median'] = df['age'].fillna(df['age'].median())`

→ `mode = df[df['age'].notnull()]['embarked'].mode()[0]`

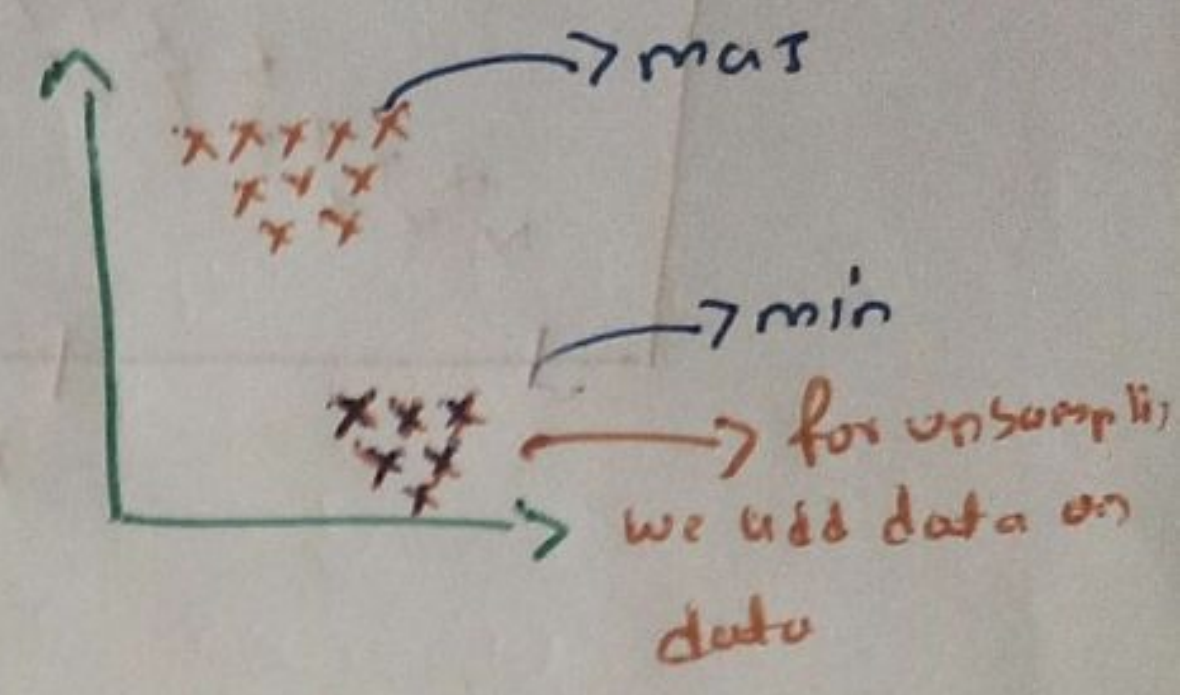
→ `df['embarked-mode'] = df['embarked'].fillna(mode)`

Imbalanced dataset Handling →

100, 900 → minority → majority

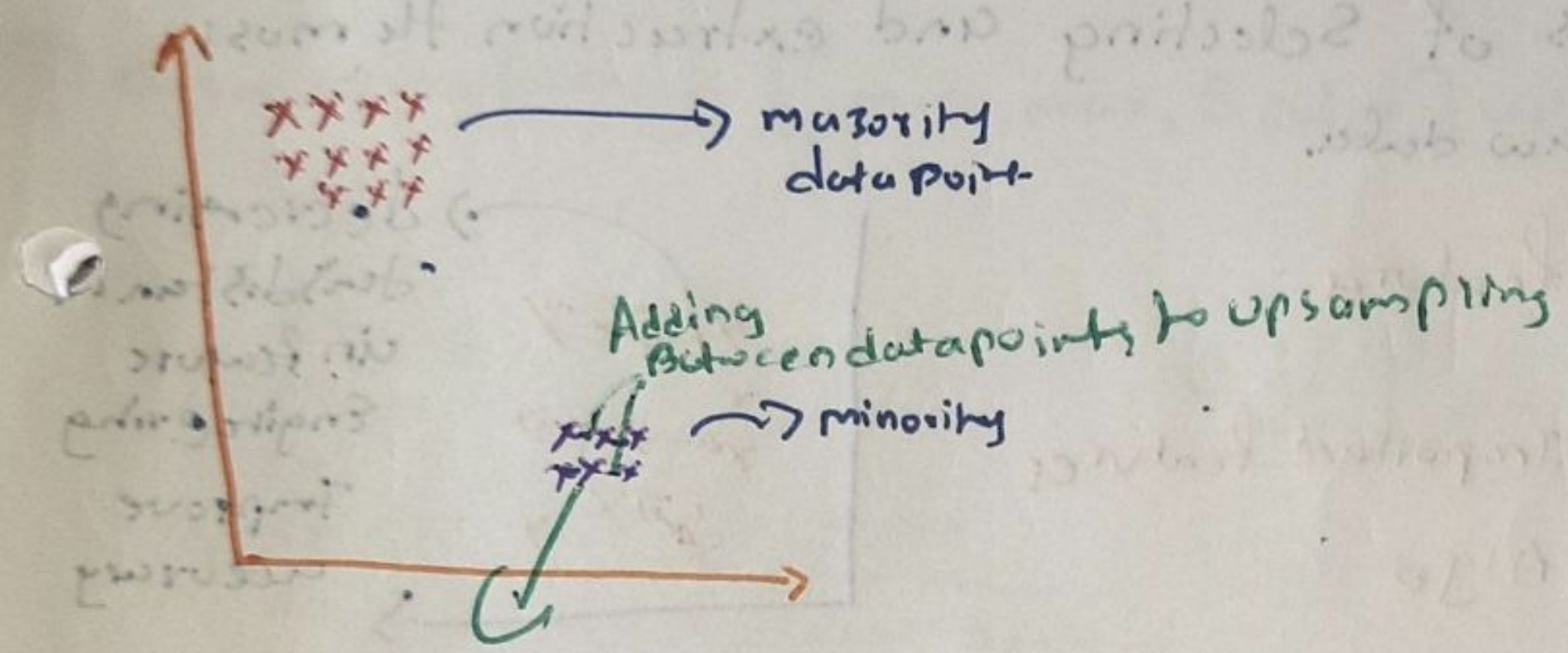
Increasing → decrease.
↓ ↓
up sampling down sampling

Balanced dataset
(100, 100) (1000, 1000)



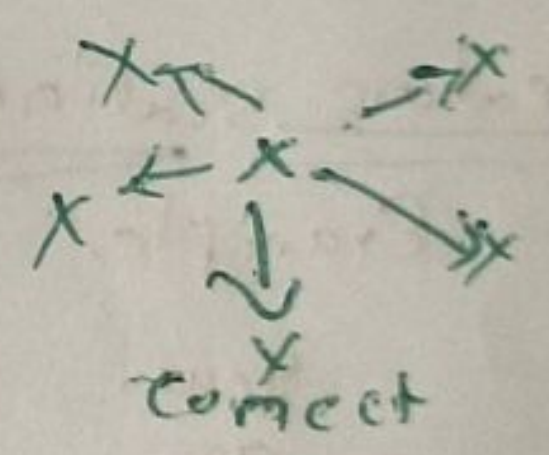
* Smote (Synthetic minority oversampling technique) → used in up sampling in minor

Adding datapoint in Between datapoints in minority



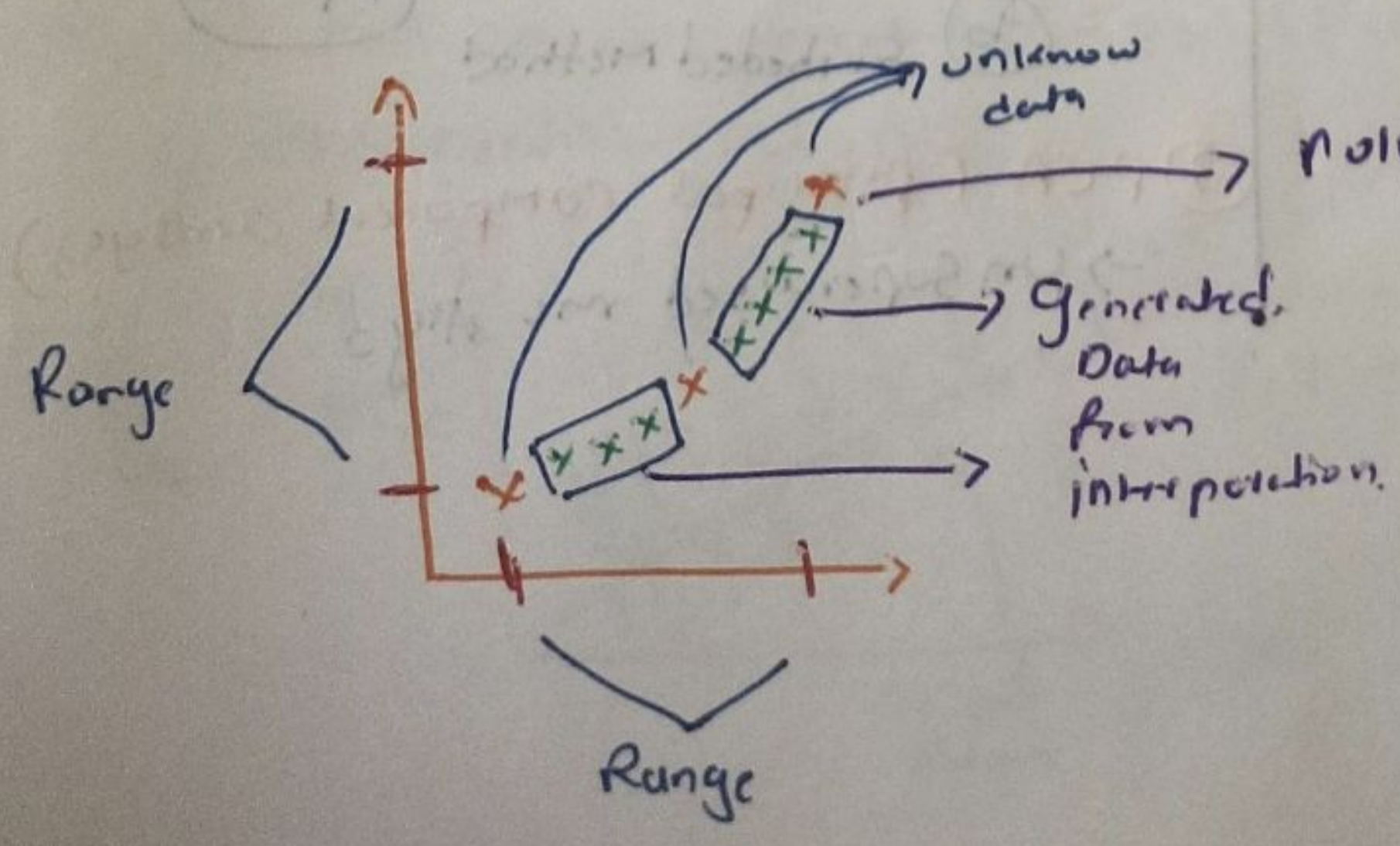
for applying Smote technique
→ pip install imblearn

from imblearn. over sampling
import smote



* Data interpolation

→ Data interpolation is the process of Estimating unknown values within a dataset Based on the known values. In python there are various libraries available that can be used for data interpolation, such as numpy, scipy, pandas, Here is an example of how to perform data interpolation using numpy library

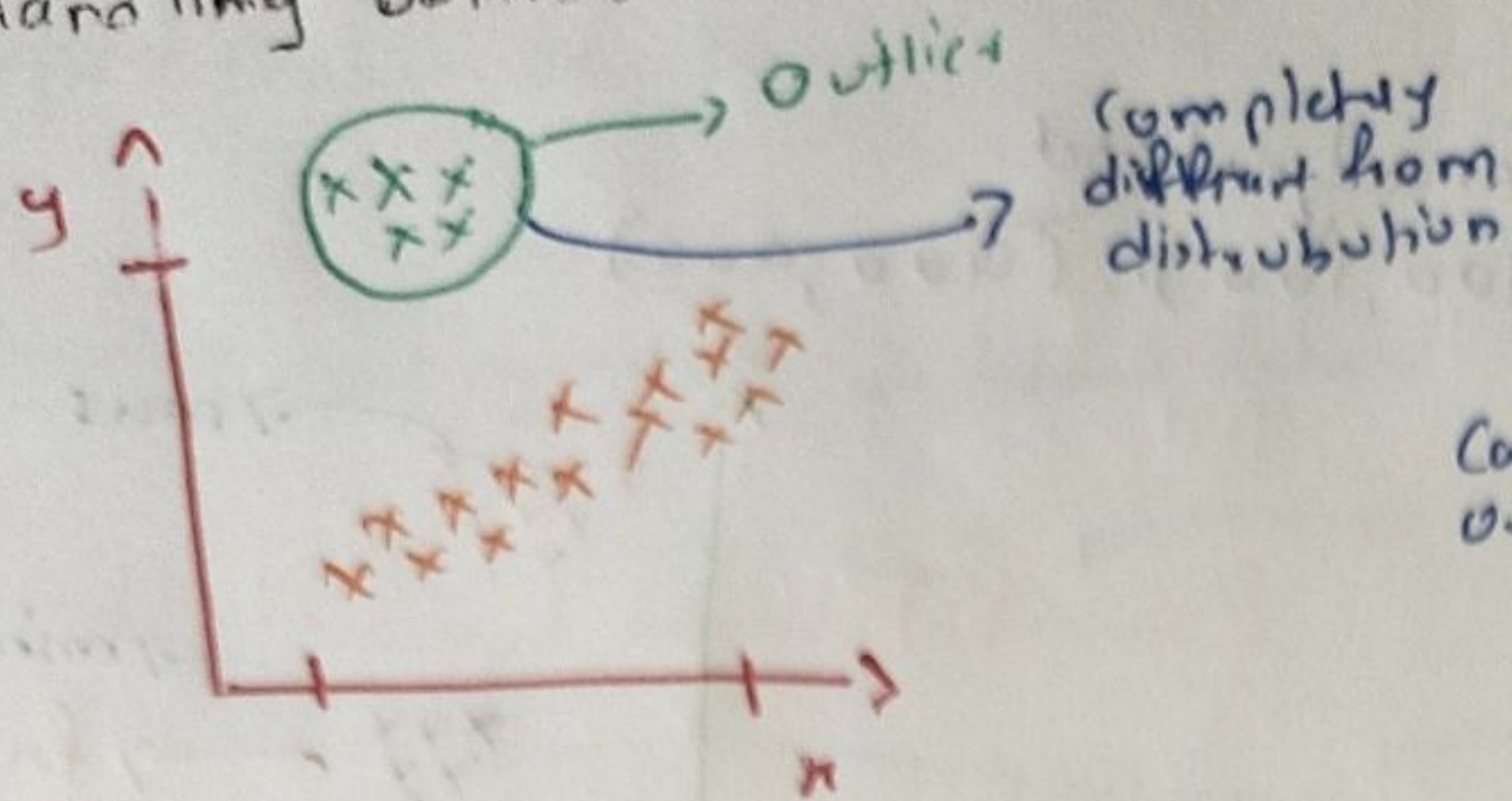


* linear interpolation

* cubic interpolation

* polynomial

★ Handling outliers



★ 5-number summary

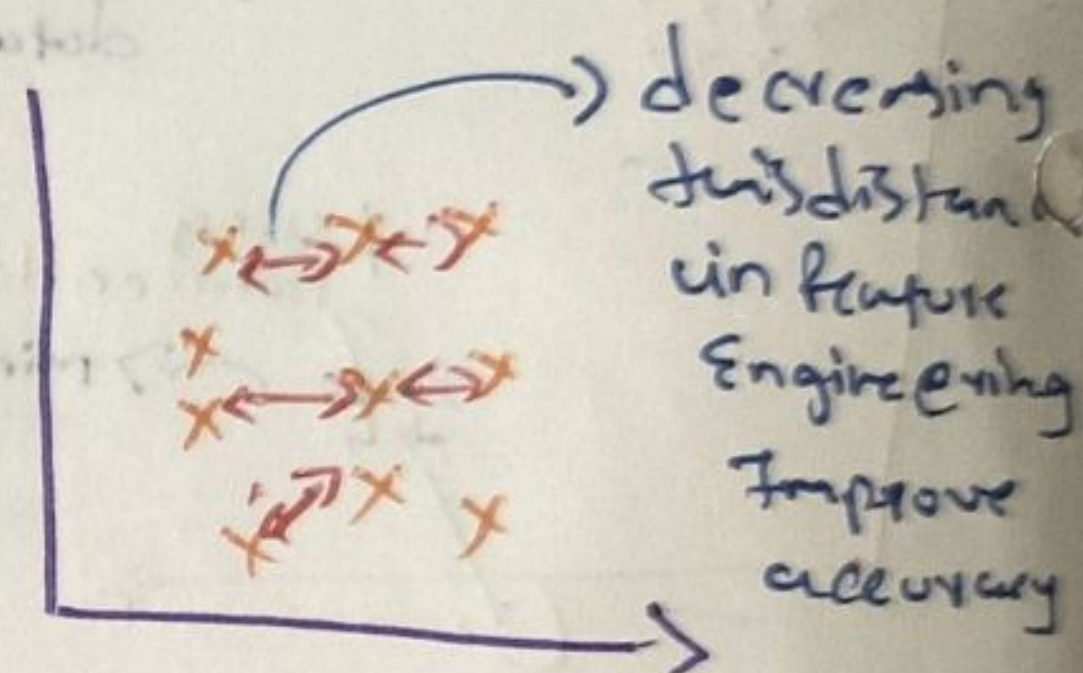
1. minimum value
 2. Q1 - 25 percentile
 3. median
 4. Q3 - 75 percentile
 5. maximum
- Calculate outliers

★ Feature Extraction → data scale down.

→ Feature extraction is process of selecting and extraction the most important features from raw data.

ML Application → 1000 features

↓
most important features
↓
ML Algo



① Feature scaling → [Scale] → [1 to 0]

Example:

Year age	weight	Height	BMI
32	70	140	—
28	75	160	—
35	80	155	—

Normalize and standardize.

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

↓
Standardization

★ normalization

[0, 1]

[min, max, scalar]

unit vector

② Feature Selection:-

→ we just pick the most important feature

(a) filter method

(b) Embedded Method

top 10

③ PCA (Principal Component Analysis)

→ Unsupervised ML Algo

Feature Scaling

① Standardization $\rightarrow [0, 1]$

\rightarrow Z-Score

Feature
(Age)

24
25
26
30
35
36

$$Z\text{-score} = \frac{x_i - \bar{x}}{\sigma}$$

$\mu = 0$
 $\sigma = 1$

Age'
 $\mu = 0$ $\sigma = 1$

② Normalization $\rightarrow [min, max, scaler] \rightarrow 0$ to 1

Age

24
25
26
28
29
30

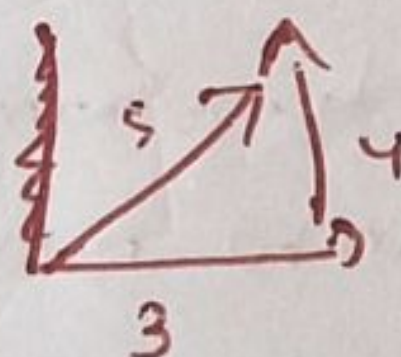
transform

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

$Age' [0-1]$

③ Unit vectors \rightarrow magnitude $\rightarrow 1$

$$\hat{a} = \left(\frac{3}{\sqrt{3^2 + 4^2}}, \frac{4}{\sqrt{3^2 + 4^2}} \right)$$

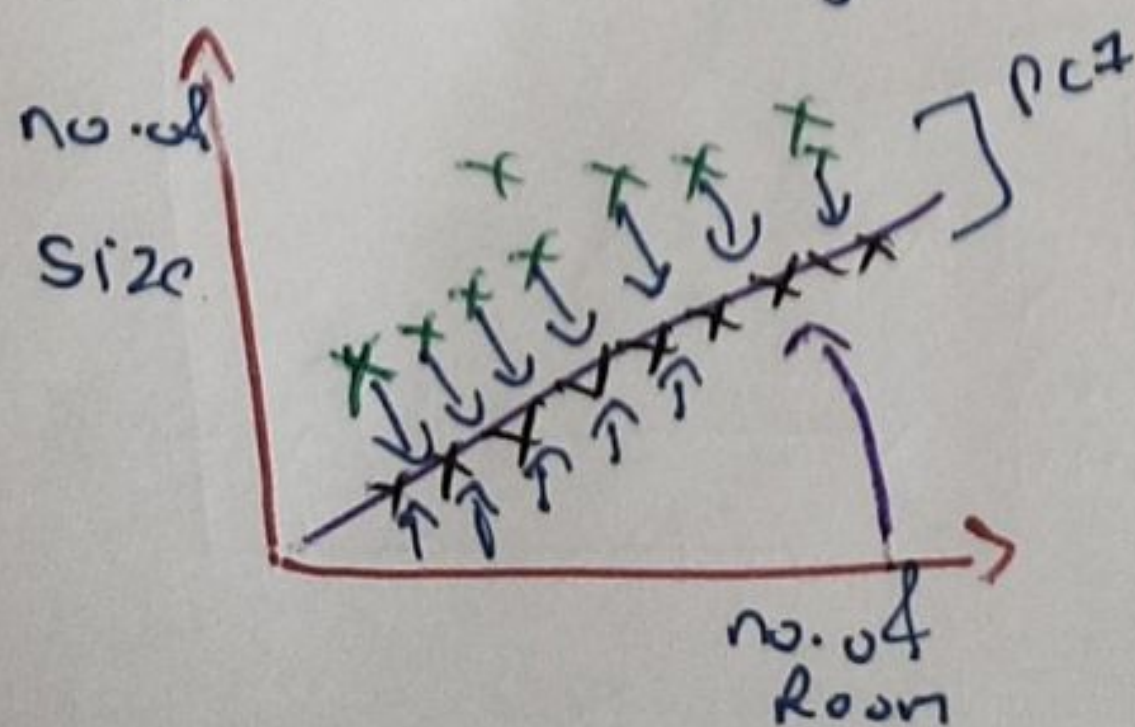


$$\|\hat{a}\| = \sqrt{\left(\frac{3}{5}\right)^2 + \left(\frac{4}{5}\right)^2}$$

$$= \sqrt{\frac{9+16}{25}}$$

$$= \sqrt{\frac{25}{25}} = 1$$

★ Principal component analysis (PCA) \rightarrow 2D \rightarrow 1D also, 3D \rightarrow 1D
 \rightarrow Unsupervised ML Algo

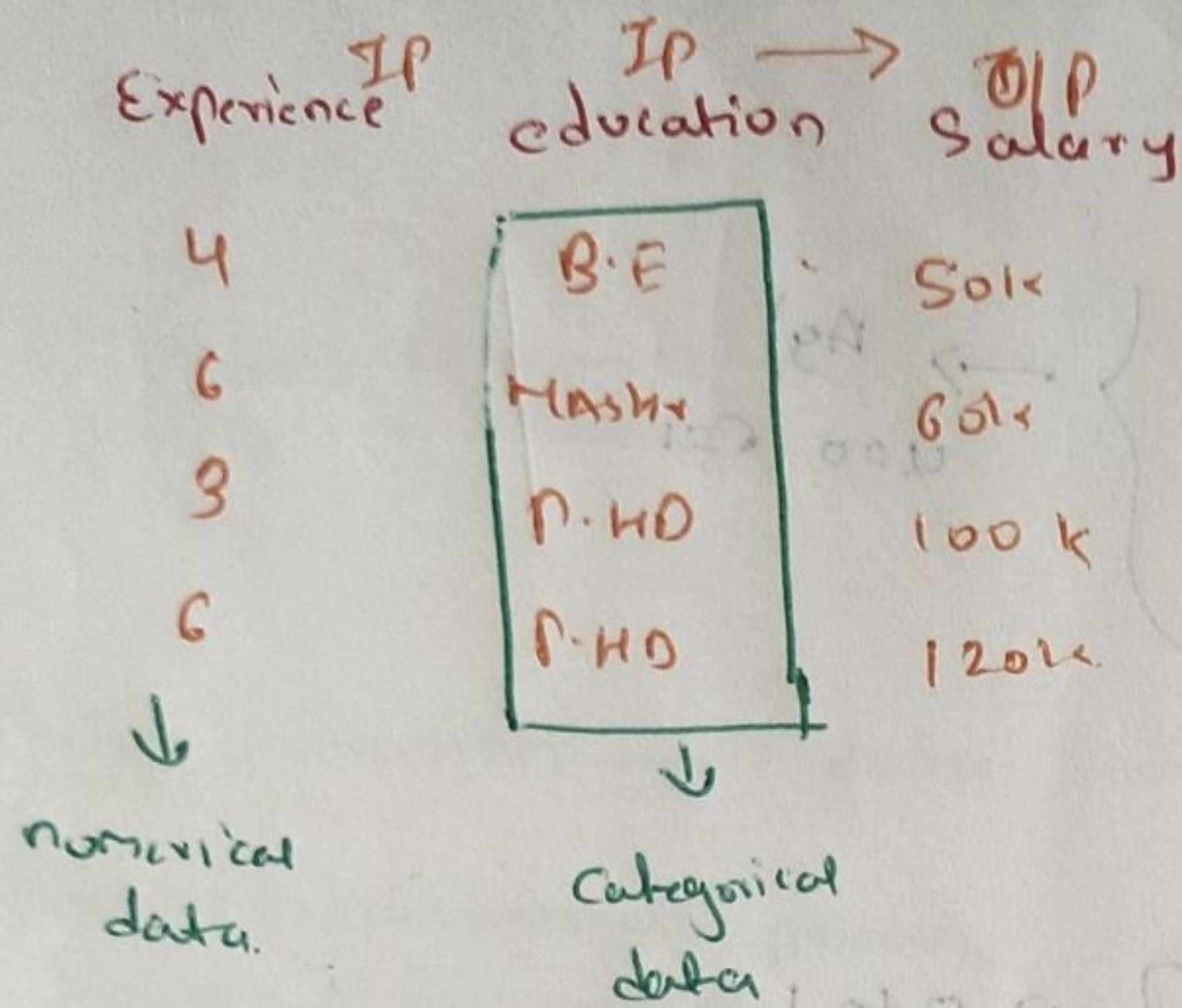


2D \rightarrow 1D
 $x, y \rightarrow$ mid

$\rightarrow [pci] \rightarrow [price]$

Train mode.

★ Data encoding → categorical data → numerical data



Data Encoding

Aim → categorical feature →

numerical



ML model



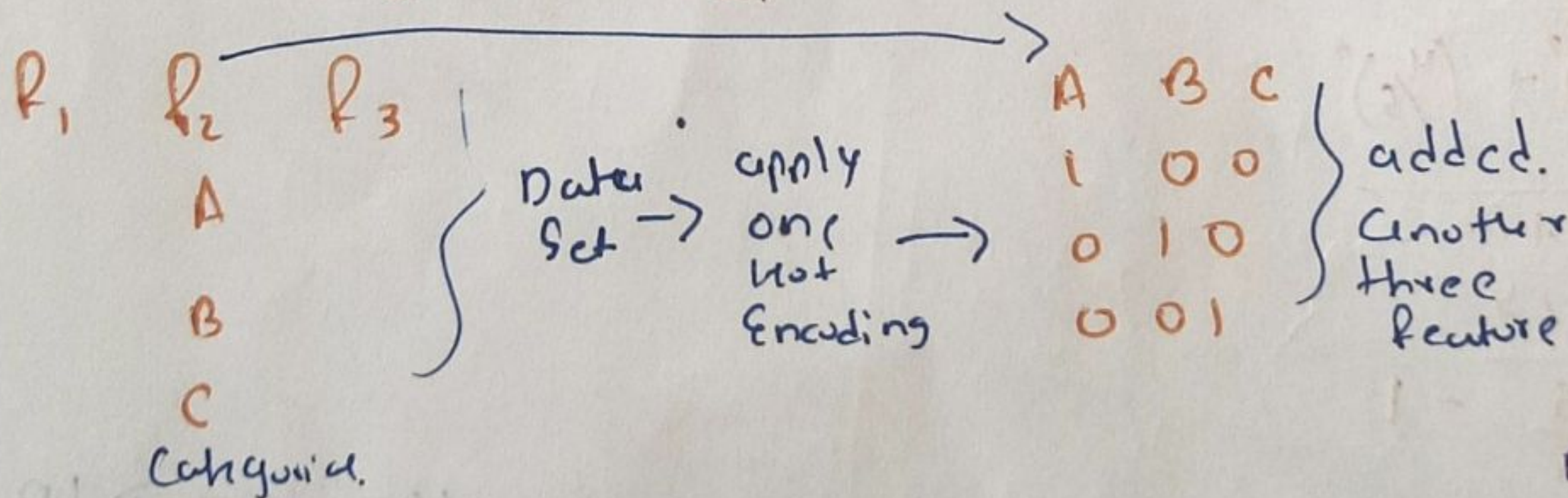
Train

★ type of Data encoding

- (a) nominal / one (one hot encoding)
- (b) ORDINAL and LABEL encoding
- (c) Target Guided Ordinal encoding

categorical → numerical.

(a) nominal encoding is a technique used to transform categorical variables that have no intrinsic ordering into numerical value that can be used in ML models



Note! - don't when there are many feature.

- Disadvantage
- (1) Sparse matrix. - [overfitting]
 - (2) 1000 feature. - 1000 matrix

③ Target Guided ordinal encoding

→ targeted variable in a mean replaced and make
new dataset for categorized data

④