```
                    ┌─────────────────────┐
                    │     feature          │
                    │   engineering        │
                    └─────────────────────┘
              ╱            │          │           ╲
        ╱                  │          │              ╲
┌──────────────┐  ┌──────────────┐ ┌──────────────┐ ┌──────────────┐
│  Feature     │  │  Feature     │ │  Feature     │ │  Feature     │
│Transformation│  │ construction │ │  Selection   │ │  extraction  │
└──────────────┘  └──────────────┘ └──────────────┘ └──────────────┘
```

Feature Transformation
→ missing value imputation
→ Handling catego — rical features
→ outlier detection
→ feature scaling

Feature construction
constructing new feature manually

Feature Selection
✓ Selecting imp feature

Feature extraction
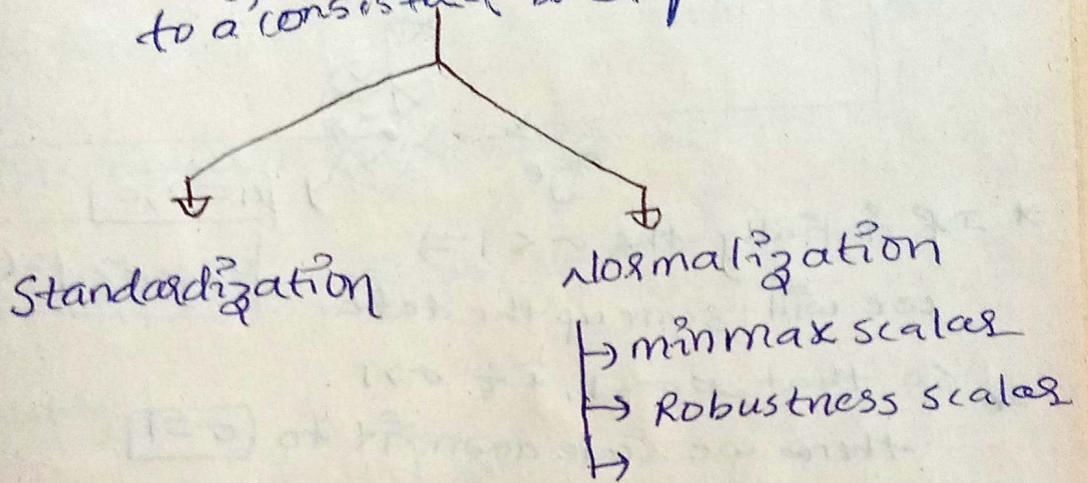extrating new feature from the given features programatically but not manually

## 1.4 feature scaling

It is technique used to transform and to standardize the independent features present in the data in a fixed range. to a consistant and predefined range



**Standardization**

**Normalization**
- min max scalar
- Robustness scalar

## Standardization :- Also called as

Z-score Normalization. No specific range *

For ex:-

| Age | Age' | |
|---|---|---|
| $x_1 = 27$ | $x_1' = \dfrac{x_1 - \bar{x}}{\sigma}$ | |
| 32 | $x_2'$ | $\bar{x}$ = mean of age column |
| 64 | $x_3'$ | $\sigma$ = standard deviation |
| 35 | 1 | of age column |
| ⋮ | ⋮ | |
| ⋮ | ⋮ | This column |
| 500 values | 500 values | $mean (\mu) = 0$ |
| | | $\sigma = 1$ |

# intution!



salary

(mean_age, mean_salary)

mean centring

age

$\mu = 0,0$

$\sigma_{age, sd} = (1,1)$

* If intially the $\sigma < 1 \Rightarrow$
  we will scale up the data
* so that $\sigma = 1$, if $\sigma > 1$
  there we scale down it to $\boxed{\sigma = 1}$

⭐ intution = mean centring + scaling by the
                              factor of $\sigma$

* when you perform standardization
  on a particular column which has
  outliers. The outliers will not
  remove on standn. we should handle
  them explictly.

when to use standardization.

1. K-means                    ( " )
2. K-Nearset neighbors ( Euclidean distance)
3. principal component analysis
   ⚹ ( try to get feature with max variance)

4. ANN — because we apply gradient descent

5. Gradient descent → During back propagation we update the weights to have similar scale weights we use.

* Distribution is always retained

⇒ Eg:- By using distpbt of age



probability

0·04

0·02
0·025
0·020
0·015
0·01
0·005
0·0

10 20 30 40 50 60 70

before scaling

P

-3 -2 -1 0 1 2 3 → age

after scaling

Standarization def$^{n}$

Here all the features will be transformed in such away that it will have the properties of a standard normal distribution with $\mu = 0$ & $\sigma = 1$
(standard devic)

Normalization:- It is a technique
often applied as part of data preparation
for machine.L. The goal of normalization
is to change the values of numeric columns
to use a common scale, without distor
-ting differences in the ranges of values
or losing information.

→ Min Max scaling
→ Mean normalization
→ Max absolute
→ Robust scaling.

⇒ Normalization is generally Min Max
Scaling

* Min Max scaling !-

Age!-

$X_i = 21$
17
16
34
56
⋮
100 values

$$X_i' = \frac{X_i - X_{min}}{X_{max} - X_{min}}$$

✳ Range of Nor!- [0 to 1]

# Intuition : Normalization



* After normalization distribution can may not be retained

* we should carefully think about outliers!! ( ? )

## Mean normalization :

$$x_i' = \frac{x_i - x_{mean}}{x_{max} - x_{min}}$$

Range = [-1 to 1]

Intution : Is to converge data toward mean centre through mean which is called as Mean centring.

* Generally we don't use it because
Standardization.

## Max Abs scaling :-

$$x_i' = \frac{x_i}{|x_{max}|}$$

we have class in sklearn as
"MaxAbsScaling".

use :- when we have sparse data.

Means data having more [Zero's]

## Robust scaling :-

wt
$x_i \sim 200$

300

75

$$x_i' = \frac{x_i - x_{median}}{IQR}$$

$IQR$ = Inter Quartile range

$$= (75^{th} \text{ percentile} - 25^{th} \text{ percentile})$$

$$= (Q_3 - Q_1)$$

use :- It is Robust to outliers

* Means when we have outliers we can
use this

# Normalization Vs standardization

→ **Q1**: Is feature scaling required?

**Q2**: If yes, 90% we use standard scalar only 10% Normal scalar

\* **Tip!** - IN [NN] we deal with images which range 0-255 pixels we know the range then we use Normalization "

## 1.2 Handling categorical features:

```
                    |
         _____|_____
        |                       |
        ▽                       ▽
     nominal                  ordinal
     features                 features
```

→ No order          The feature in which
of categories        categories have order

$eg^n$:- <u>States</u>        $eg$:- <u>Ed$^n$</u>

*Telangana              PG (post graduate)

maharastra              UG (under ..)

Andhra.                 HS (High school)

$son$.                   order ⇒ PG > UG > HS

"One Hot Encoding"      $Sol^n$: 1) Ordinal
                                     Encoding

                                 2) label Encoding

# Ordinal Encoding :-

I) Label Encoding : we use this when our label/output column has categories like (yes/NO) (classification)

Ej :-

| No. of hours studied | No. of hours slept | pass/fail y |
|---|---|---|
| 5 | 8 | pass |
| 6 | 9 | fail |
| 7 | 7 | pass |
| 8 | 6 | fail |

We use label Encoding

1
0
1
0

2) ordinal Encoding :- we give order as
PG=2, UG=1, HS=0

| Education | (Edn)! |
|---|---|
| HS | 0 |
| PG | 2 |
| UG | 1 |
| PG | 2 |
| HS | 0 |
| UG | 1 |

* Both label / ordinal encodings are
same. But label is not only used on
dependent feature, ordinal used on
independent features

## One Hot Encoding

eg:- color

yellow

Blue

Red

Blue

yellow

Bule

## SK-learn Pipelines :-

* pipelines is a mechanism chains togeth
-er multiple steps so that the output of
each step is used as input to the next
step.

* pipeline makes it easy to apply the
same preprocessing to train & test

# Function Transformers :-

<u>Goal</u> :- Is to make ~~non~~ normal distribution

<u>normal distribution</u> :

probability → normal distribution

probability

probability
distribution
~~curve~~ Function

↑ → normal distribution



20  30  4  50 age

(PDF)

P ↑

→ Not
normal
distribution

→ age column

use of normal distribution :- In ML
we have Statistical algorithms
there are linear regression, logistic
regression :- etc. so that we can get
good accuracy.

# Sklearn Transformers

Function Transformers
- → log trans
- → Reciprocal trans
- → sqrt trans
- → custom (we can make)

power trans former
- → Box-Cox
- → Yeo-Johnson

Quantil trans

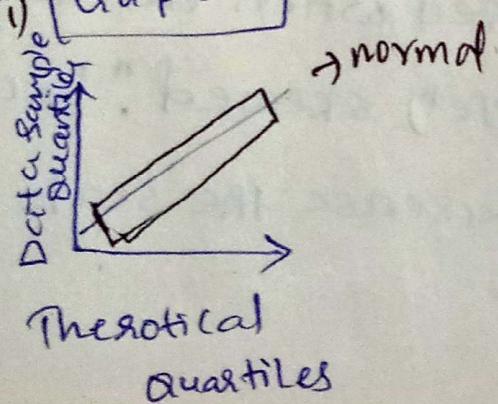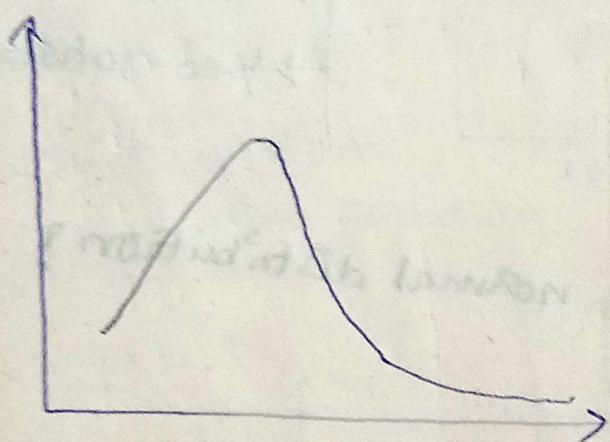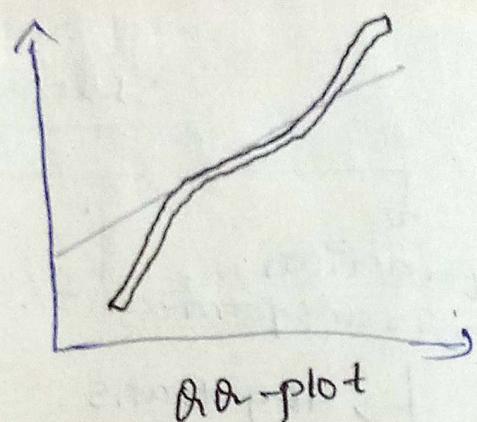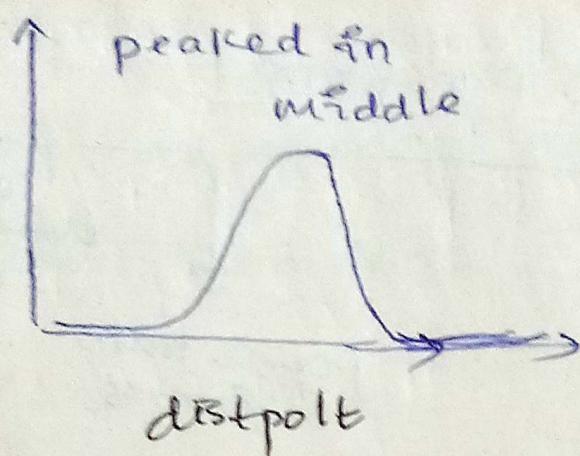1) How you find the normal distribution?

ⅰ) sns.distplot

ⅱ) pd·skew ()

- ↓ -1 negatively skewed
- ↓ 0 normal
- ↓ +1 positively skewed



ⅲ) QQ plot

Data Sample Quantiles

→ normal

Therotical Quartiles

peaked in middle

distpolt

QQ-plot



+ve skewed

QQplot

## log transformation

eg: Age

$x_1 = 21$

$$\boxed{x_i' = \log(x_i)}$$

35

40

* The log transform is

45

applied when our data

47

is "+vely skewed". Because

50

log decrease the scale.

Eg:-



20  30  40  70  90

20  3  4  5  6  7  8



10  20  30    100    100

$\sharp \log$

0    $\log(10)$    $\log(100)$    $\log(000)$    10
     10            10            10

→ range is decreased
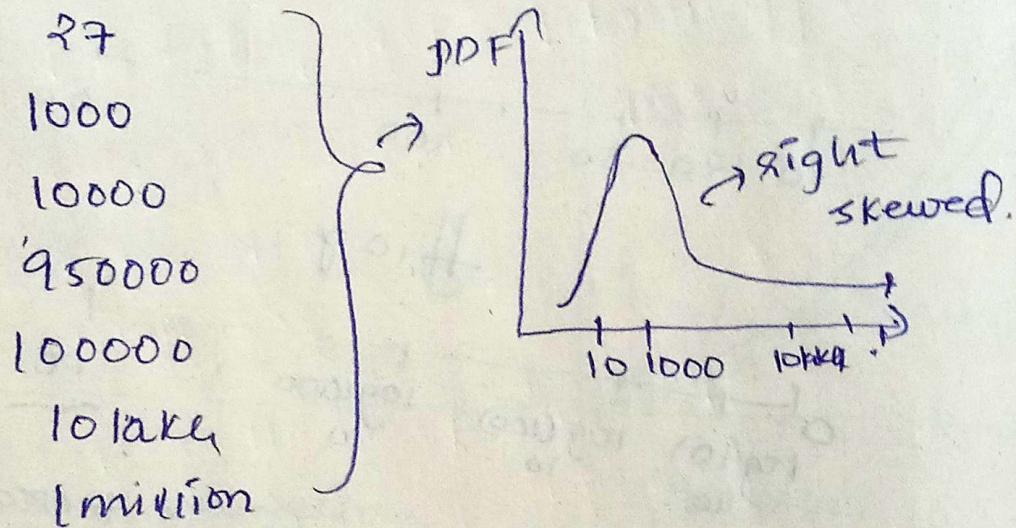
<u>Reciprocal</u> $(\frac{1}{x})$ : square $(x^2)$ ↑ : sqrt $(\sqrt{x})$
              : used for :
              : left skewed :
              : data :

# Encoding numerical features :-

Eg: In playstore if you see the count of
downloads for different apps, they
will be like this :

## downloads

27
1000
10000
950000
100000
10 lakh
1 million



PDF

→ right
    skewed.

10 1000    10kk

But when you convert them into certain
ranges (or) bins like this:

10 - 1000
1K - 10K
10K - 100K
100K - 1000K
1000K - 10000K



* what we have done above ?

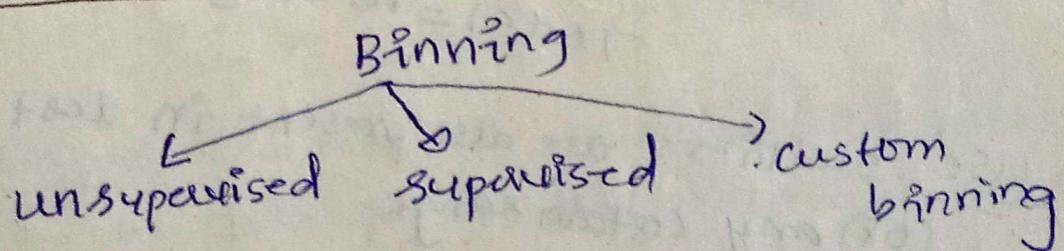we have converted the "numerical data" to "categorical data".

* The technique we use to convert numeri—cal to categorical is "Discretization".

1) **Discretization** :- It is a process of transforming continuous variables into discrete variables by creating a set of contiguous intervals that span the range of variable's values. Discretization is also called **binning**, where bin is an alternative name for interval.

Why use Discretization :-

1. To handle outliers
2. To improve the value spread.

Types of Discretization:—

Binning

unsupervised        supervised        custom binning

unsupervised          supervised          custom binning

├→ Equal width        ├→ Decision
│   (uniform)         │    Tree
├→ Equal frequency    │   Binning
│   (quantile)
└→ k means
    binning

## Equal width / uniform Binning :-

Age :- 27, 34, 84, 26, ........160

$$max_{age} = 160$$
$$min_{age} = 0$$

we select the no. of bins = 16

$$A \; Range = \frac{max\_age - min\_age}{10 \; bins}$$
$$= \frac{160 - 0}{160} = 10$$

Now bins range will be

$(0-10), (10-20), (20-30), \; ......... \; (150-160)$

Total = 16 bins

use :- outliers are also present in last
       (or) any certain range.

ii) No change spread (or) distribution.

# Equal Frequency/Quantile Binning :-
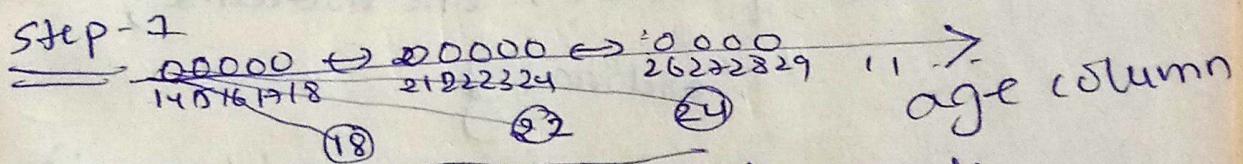
⇒ we should decide no. of intervals = 10

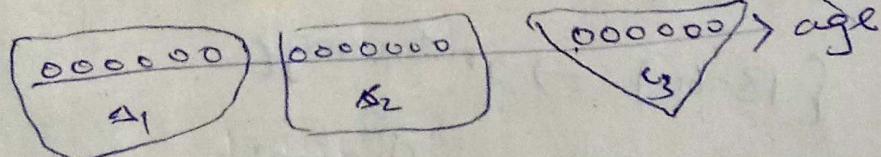Each interval contains 10% of total observations.

0-16; 16-20; 20-22; 22-25; ----

{ 50-74

Note! Here bin width is not constant

## K-Means Binning :-

Used! when our column is in the form of clusters (or) groups

Step-1

00000 ⟺ 00000 ⟺ 0000  11 ⟶
14151617 18    21 22 23 24    26 27 28 29      age column
      ⑱              ㉓         ㉔

randomly intilized centriods

Step-2 Now distance b/w one point and every centroid is calculated. If the distance is less then that the point assigned to a particular centriod.

( 0000000 )  | 0000000 |  ( 0000000 ) ⟶ age
    △₁            △₂            △₃

Step-3 Now mean of each point is selected in cluster is taken then centriod position is changed within the cluster.

Step 4 :- process repeated when position of
           of centriod is changed.

* finally the range of each bin is

$$0 - c_1, \quad c_1 - c_2, \quad c_2 - c_3 \ldots$$

Code :

from sklearn. preprocessing import KBinsDiscre
                                        tizer

parameters :-

   n_bins = our wish (any int)

   encode = {'onehot', 'one-hot-dense,
               'ordinal'}

   strategy : {'uniform', 'quantile', 'kmeans'}

Custom/Domain_Based_Binning :-

* we create own range like

   [0-18] → kids

   {18-60] → 2nd range

   {60-80] → 3rd range

* This class is not present in
  sklearn. we should use

pandas.

## 2. Binarization!

For ex:- 

| age | age' |
|-----|------|
| 16  | 0    |
| 17  | 0    |
| 12  | 0    |
| 22  | 0    |
| 35  | 0    |
| 45  | 0    |
| 77  | 1    |
| 65  | 1    |
| 62  | 1    |
| 54  | 1    |
| 80  | 1    |

I set a threshold = 45

If many age ≤ 45 ﹢ch

If $age_i \leq 45$   $age_i' = 0$

$age_i > 45$   $age_i' = 1$

## correlation vs co-variance

| correlation ↓ | co-variance ↓ |
|---------------|---------------|
| Direction ✓   | Direction ✓   |
| Strength ✓    | Strength ✗    |
| Range = -1 to 1 | Range:- $-\infty$ to $+\infty$ |

# Mixed data:

For ex:

| cabin | categorical | Numerical |
|---|---|---|
| B5 | B | 5 |
| D4 | D | 4 |
| S2 | S | 2 |
| H3 | H | 3 |
| C23 | C | 23 |
| D41 | D | 41 |

Type-1:-

Type-2:-

| column | cat | num |
|---|---|---|
| 7 | 7 | NaN(null) |
| 3 | 3 | NaN |
| 1 | 1 | NaN |
| A | NaN | A |
| C | NaN | C |
| 4 | 4 | NaN |
| 3 | 3 | NaN |

* ## MISING Values

Remove the rows

Impute
  - univariate
    - numerical
      - → mean/median
      - → Any random value
      - → end of distribution
    - categorical
      - → most frequent value
      - → missing
  - multivariate
    - KNN Imputer
    - Iterat Imput

# Complete case Analysis (CCA):-
### (or)

* Complete case analysis (CCA), is also called "list-wise deletion" of cases, consists in discarding observations where values in any of the variables are missing.

* Complete case analysis means literally analyzing only those observations for which there is information in all of the variables in the dataset

## Assumption For CCA:-

we apply the CCA when our data is randomly missing

* we don't apply CCA:-

i) Let us assume in one of the column if 1st 50 rows values or last 50 rows values are missing then we don't apply it

* we follow [MCAR] = missing complete at random.

* df. isnull().mean() * 100

we go percentage of missivalues in each column

Numerical data:-

1)* Mean/Median is discussed before.

2)* Arbitary value constant

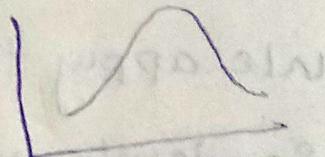Simple Imputer (strategy = 'constant', fill-value = any value)

when to use? When daf is not missing at random.

3) End of Distribution Imputation:-

when to use :- When data is not missing at random.

Case-1

*If normality distributed



the we fill with $(mean + 3\sigma)$ ✓
(or)
$mean - 3\sigma$ ✓

&
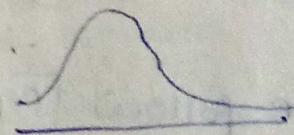
Case-2

If +vely or -vely skewed then we use



* we fill with $Q_1 - 1.5 IQR$

$Q_3 + 1.5 IQR$

$Q_1 = 25^{th}$ percentile value

$Q_3 = 75^{th}$ percentile value

$$\boxed{IQR = Q_3 - Q_1}$$

## Disadvantges

→ Outlies

→ Distribution (pdf)

→ covariance change

→ variance per (SD) des largely

## Univariate-categorical

* most frequent :

    codn! If data is MCAR

    $\overline{ii}$) $5 < \%$

* "Missing" :- Here we will with the word 'Missing' in place of Nan values

    codnt NOT MCAR

    $\overline{i}$) $\frac{5}{10} > \%$

# Random Imputation :-

* It is applied on both numerical and categorical data.

* It doesn't matter
  
  i) MCAR
  
  ii) <5%

* Eg
  
  | Sex | Age |
  |-----|-----|
  | M | 26 |
  | F | 32 |
  | M | 54 |
  | F | 29 |
  | | 12 |
  | M | |
  | F | 52 |

  randomly chosen

  These missing values will be replaced by the random number in that column

* **Imp!**
  * preserves the variance of the variable
  * Memory heavy for deployment, as we need to store the original training set to extract values from and replace the NA in coming observations
  
  * Well suited for linear model as its does not distrot the distribution, regardless of the % of NA

# KNN Imputer

|  | col1 | col2 | col3 | col4 |
|---|---|---|---|---|
| R1 → | 12 | 14 | — | 13 | ← focused |
| R2 → | A1 | 15 | 16 | 12 |
| R3 → | 14 | 12 | 17 | 25 |

$(12, 14/-, 13) = $ point 1

point 2 $= (-, 15, 16, 17)$

* The distance b/w $(1,2,3) \times (7, 8, 9)$ is calculated through E.q $\overset{c}{\text{Euclidean}}$ distance.

* But we have missing values so we use "nan_euclidean"

For ex: nan_Euclidean b/w P1 & P2

$$= \sqrt{\text{weight} \times \left[ (x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2 \cdots \right]}$$

$$\text{weigh} = \frac{\text{Total no. of actual pairs}}{\text{no. of pairs present}}$$

R1 & R2
dist A1 $= \sqrt{\text{weigh} \times \left[ (14-15)^2 + (13-A1)^2 \right]}$

R1 & R2 $= \frac{3}{2} \times \left[ ( \quad ) + ( \quad ) \right] = 10$

$$R_2 \ast R_3 = \sqrt{\frac{3}{3} \times \left[ (\quad)^2 + (\quad)^2 + (\quad)^2 \right]} = 50$$

\* If no. of neighbors $= 2 = K$

The $A_1 = \frac{12 + 14}{2}$

weights $=$ uniform $\Rightarrow$ we have take the mean of the nearest Rows.

If weights $=$ distance

then $A_1 = \dfrac{\frac{1}{d_1}(x_1) + \frac{1}{d_2}(x_2)}{2}$

$$= \dfrac{\frac{1}{10}(12) + \frac{1}{50}(14)}{2}$$

$$= \quad --- \,,$$

Advantages:

1) most accurate values than mean/ median/ random

Disadv:

→ If the dataset is large then it takes so much time to fit. means be - cause it should calculate more distances.

ii) During the deployment time we should upload our floating dataset. because if any input value of the user is missing we should calculate the dB 7omces.

Iterative Imputer / MICE :-

MICE = Multivariate jomputation by Chained Equations.

## Assumptions

MCAR = Missing completely at random
MAR = " at random
MNAR = " Not " "

\* we get Good results of mice when [MAR] cod$^n$ is present.

Adv ∧ Dis adv
↓                    ↓
most              7) on large dataset its
accurate           working is very slow
                   ii) memory = we should
                   upload the data set on the
                   server.!

# Outliers:- The data that behave differently from others.

→ Big Question!- In some situations outlier can also be use & dangerous, when should then.

Effect of outliers on ML algorithms:-

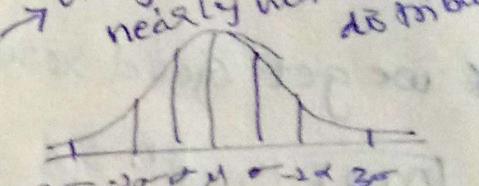* The algorithms in which the weights are assigned are generally effected. They are:

     i) Linear Regression

     ii) Logistic Regression

     iii) Adaboost

     iv) Deep learning

## How to [defect] outliers:

1. Normal distribution → only if our data is nearly normal distributed *

For ex! age col$^n$

* If data in age col is m

     $> \mu + 3\sigma$ OR $< \mu - 3\sigma$ then it is outlier.

2. Skewed Distribution:

IQR = Interquartile Range

$(Q1 - 1.5 \times IQR)$   $Q_1$   median   $Q_3$     $(Q_3 + 1.5 \times IQR)$

$(25^{th}$ parcntile$)$   $(q^{th}$ parcntile$)$