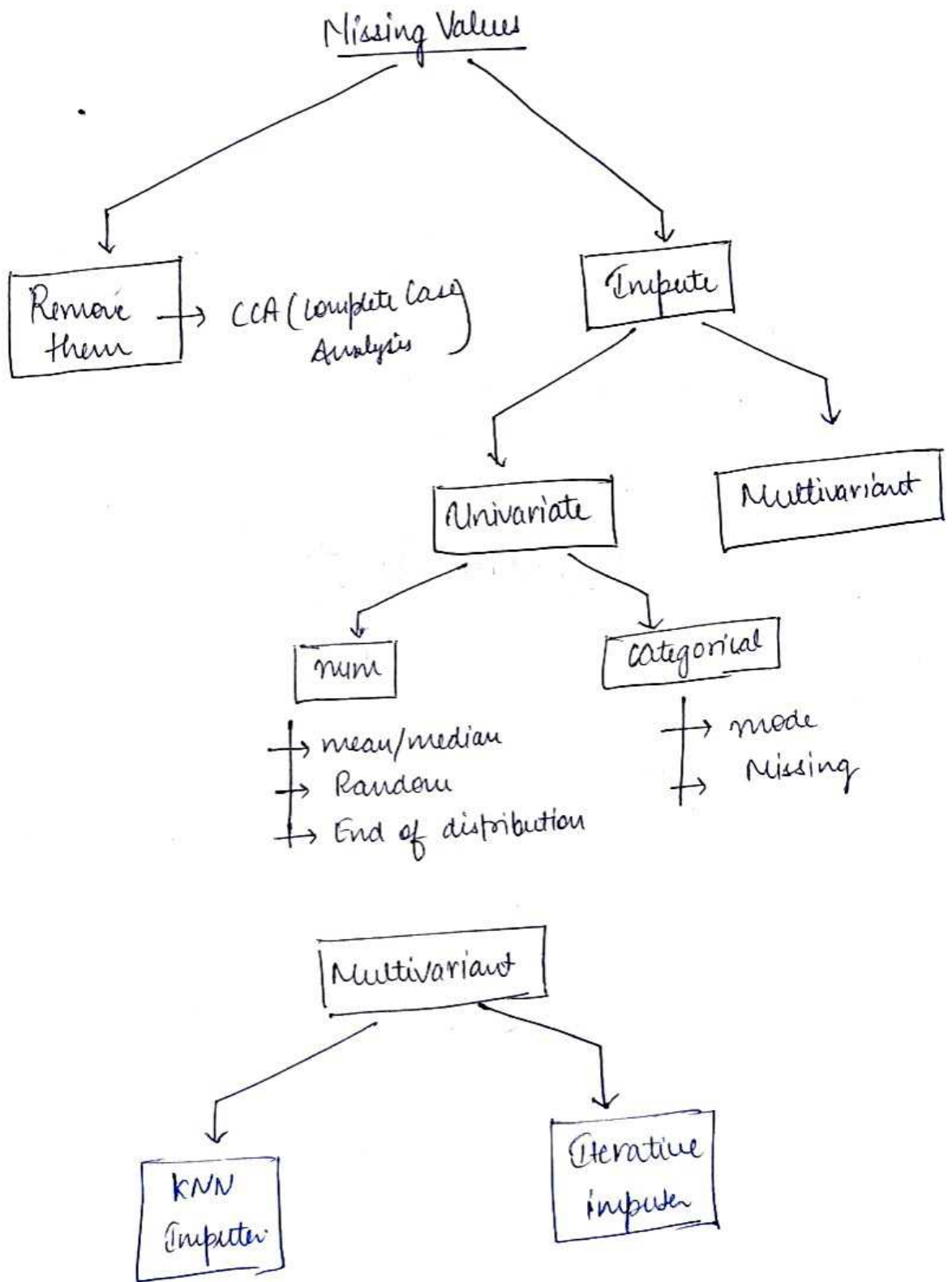


# Handling Missing Data

①



## Complete Case Analysis

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

Complete Case Analysis means literally analyzing only those observations for which there is information in all of the variables in the datasets.

### Assumption for CCA

1) Missing Completely at Random

1000, 47

Age → 50 → missing value (Randomly)

↳ 950, 4

### Advantage / Disadvantage

#### Advantage

1. Easy to implement as no data manipulation required.
2. Preserves variable distribution (if data is MCAR, then the distribution of the variables of the reduced dataset should match the distribution in the original dataset).

## Disadvantages

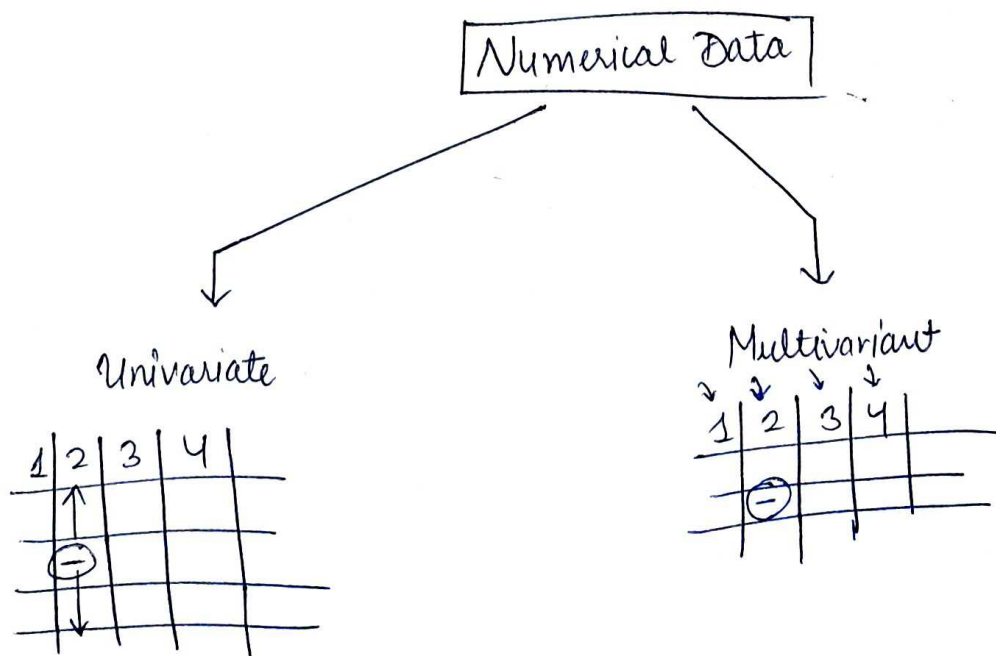
(2)

- 1) It can exclude a large fraction of the original dataset (if missing data is abundant)
- 2) Exclude information could be informative for the analysis (if data is not missing at random)
- 3) When using our model in production, the ~~production~~ model will not know how to handle missing data.

## When to Use CCA?

1. MCAR
2. 5% < (less than 5%)

## Handling Missing Numerical Data



## Mean / Median

## Imputation

Age

27

32

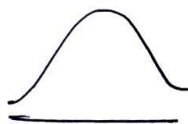
Na

27

Mean

Median

Mean →



Benefit

1. Simple

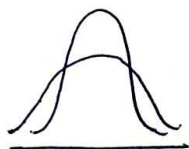
2. 5% >

Median →



## Disadvantage:

1.) Distribution shape



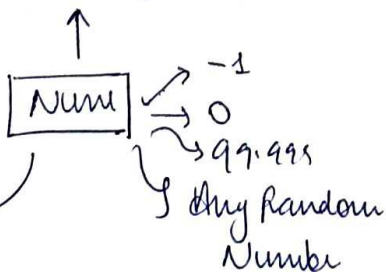
2.) Outliers

3.) Covariance / Correlation changes

## 3. Arbitrary Value Imputation

↳ category → NA →

Missing



Data is not missing at random

Benefit → easy to apply

## Disadvantage:

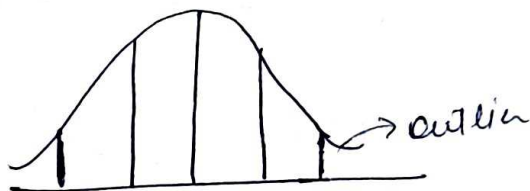
- Distribution shape
- Variance
- Covariance changes



# End of Distribution Imputation

(3)

End of Distribution



Normality

↳ (mean + 3σ)  
(mean - 3σ)

Skewed



(IQR priority)

Q3 - Q1

75th

25th

Q1 - 1.5 IQR  
Q3 - 1.5 IQR

Benefit

→ Easy to apply

- disadvantage:
1. Distribution shape
  2. Outliers
  3. Covariance / correlation changes

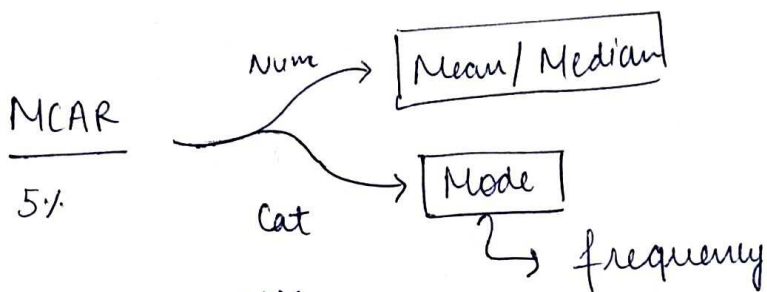
## Handling categorical Missing Data

Most frequent

Missing

→ Random

## Most Frequent Value Imputation



Benefit → easy to apply

Disadvantage

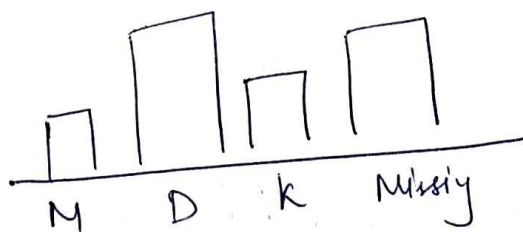
↳ Distribution Shape (change)

## Missing Category Imputation

City  
Mumbai  
Delhi  
Kolkata

(10%) → Outlier

New column → Missing



# Random Imputation

4

Age

26  
32  
56  
41  
NA

NA

Random Numbers



Gender

M  
F  
F  
M  
NA

Benefit

easy to apply  
Data distribution and variance same

Random Imputation → Linear Model  
Logistic Model

Disadvantage

Memory heavy for deployment, as we need to store the original training set to extract values from replace the NA is coming observation.

2) Well suited for linear models as it does not distort the distribution, regardless of the % of NA.

## Missing Indicator

<u>Age</u>	<u>Fare</u>	<u>Age - NA</u>
27	32	F
41	35	F
NA	41	T
62	32	F

Automatically select value for Imputation

Use Grid Search CV

→ Try different combination and select prefect one.

Multivariant

KNN Imputer

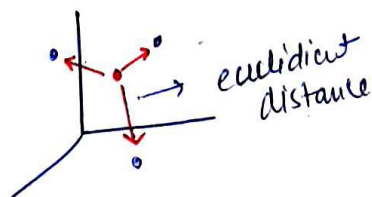
Working:

There are 4 columns  
first column have Null value. So  
find value of null value on the basis of which  
row similar to Null value of row.  
Similarity found with the help of euclidean distant.

S.No	Feature 1	Feature 2	Feature 3	Feature 4
1	33	—	67	21
2	—	45	68	12
3	23	51	71	18
4	40	—	81	—
5	35	60	79	—

Let we have 3-D data  
(33, —, 67)

k = No. of Neighbour





\* In Euclidean distance, we use  $(x, y)$

(5)

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \text{ for 2D and for}$$

3D  $(x, y, z) \rightarrow$

\* But sometimes in data have more than 1 Null value. So, we use Manhattan distance

$$\text{dist}(x, y) = \text{sqrt}(\text{weight} * \text{sq. distance from present coordinates})$$

Formula  $\rightarrow$  Where,  
weight = (Total # of coordinates) / (# of present coordinates)

Let's calculate distance Ignore bcz of Null value present

SNo	Feature 1	Feature 2	Feature 3	Feature 4
1	33	—	67	21
2	—	45	68	12
3	23	51	71	18
4	40	—	81	—
5	35	60	79	—

$$\text{distance} = \sqrt{\frac{3}{2}} (67 - 68)^2 + (21 - 12)^2$$

$\rightarrow$  weight

\* We choose Row 1 and Row 2.

We can also check Row 3, and Row 4, Row 5

Total No. of Coordinate = 3

Total No. of Present Co-ordinate = 2

Distance bet<sup>n</sup> Row2 and Row1

$$d = \sqrt{\frac{3}{2} (68-67)^2 + (12-21)^2}$$

$$d = 11.29$$

Distance bet<sup>n</sup> Row2 and Row3

$$d = \sqrt{\frac{3}{3} ((51-45)^2 + (71-68)^2 + (18-12)^2)}$$

$$d = 9$$

Dist<sup>n</sup> bet<sup>n</sup> Row2 and Row4

$$d = \sqrt{\frac{3}{1} ((81-68)^2)}$$

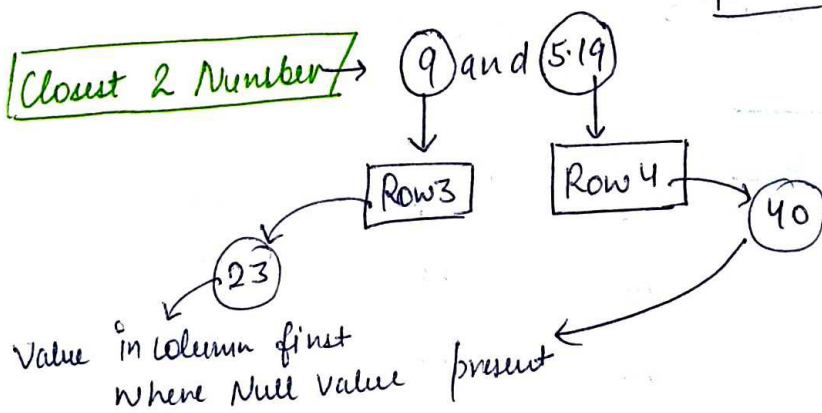
$$d = 5.19$$

Distance bet<sup>n</sup> Row2 and Row5

$$d = \sqrt{\frac{3}{2} ((60-45)^2 + (79-61)^2)}$$

$$d = 14.79$$

Closest 2 Number →



$$\frac{23+40}{2} = \frac{63}{2} = 31.5 \text{ is value of Null Value in column first.}$$

\* If I consider 3 value or closest 3 value

then

$$\frac{23+40+33}{2} = \text{Null Value}$$

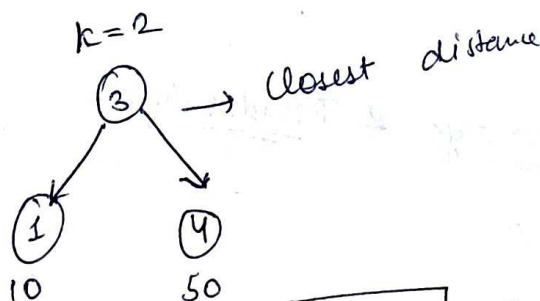
# Advantages & Disadvantages

(6)

1. More accurate
2. More Number of calculation (calculate distance with every point)
3. Upload Train Data on server.

Sklearn  $\rightarrow$  weight = Uniform, distance

1 $\rightarrow$	70			
2 $\rightarrow$				
3 $\rightarrow$	-			
4 $\rightarrow$	50			



Uniform

$$\frac{70 + 50}{2} = 60$$

Null Value

distance

$$\frac{\frac{1}{10} \times 70 + \frac{1}{50} \times 50}{2} = 4$$

Null Value

Iterative Imputer / MICE

MICE stands for Multivariate Imputation by chained Equation.

There are four types of missing data

1. Missing completely at Random (MCAR)
2. Missing At Random (MAR)
3. Missing Not At Random (MNAR)

Assumption → When Missing data is MAR.

↓  
Means when optionally data missing (ex: option → you have to fill date of birth. Some of people fill some of not)

MNAR

→ Intentionally data missing

□ □ □ □  
↔  
Not relate

□ □ □ □

↔  
Missing data relate to other

Advantage & Disadvantage

↳ accurate

↳ slow process  
↳ Train Data on server

How it works?

1. Actual Dataset

	R&D spend	Administration	Marketing Spend	Profit
21	8.0	15.0	30.0	11.0
37	4.0	5.0	20.0	9.0
2	15.0	10.0	41.0	19.0
14	12.0	16.0	26.0	13.0
44	2.0	15.0	3.0	7.0

2. Removing the target column

	R&D spend	Administration	Marketing spend
21	8.0	15.0	30.0
37	4.0	5.0	20.0
2	15.0	10.0	41.0
14	12.0	16.0	26.0
44	2.0	15.0	3.0



3. Introduced some fake nan values

⑦

	R&D Spend	Administration	Marketing Spend
21	8.0	15.0	30.0
37	<b>NAN</b>	5.0	20.0
2	15.0	10.0	41.0
14	12.0	<b>NAN</b>	26.0
44	2.0	15.0	<b>NAN</b>

Step 1 → Fill all the NAN values with mean of respective col

	R&D Spend	Administration	Marketing Spend
21	8.00	15.00	30.00
37	<u>9.25</u> mean	5.00	20.00
2	15.00	10.00	41.00
14	12.00	mean → <u>11.25</u>	26.00
44	2.0	15.00	mean → <u>29.25</u>

Step 2 → Remove all column missing values

	R&D	Admins	MS
21	8.0	15.00	30.0
37	NAN	5.00	20.0
2	15.0	10.00	41.00
14	12.0	11.25	26.00
44	2.0	15.00	29.25

	Admins	MS
21	15.00	30.00
2	10.00	41.00
14	11.25	26.00
44	15.00	29.25

Test data

Training Data

Output Data

Input data

Step3:- Predict the missing value of col1 using other cols

	R&D Spend	Administration	Marketing Spend
21	8.00	15.00	30.00
37	23.14	5.00	20.00
2	15.00	10.00	41.00
14	12.00	11.05	26.00
44	2.00	15.00	29.25

Step4:- Remove all col2 missing Value

	R&D	Admin	MS
21	8.0	15.0	30.0
37	23.14	5.0	20.0
2	15.0	10.0	41.0
14	12.0	NAN	26.0
44	2.0	15.0	29.25

Input Value (red box)

Output Value (blue box)

Test Value (green box)

Step5:- Predict the missing values of col2 using other cols

	R&D Spend	Administration	Marketing Spend
21	8.0	15.0	30.0
37	23.14	5.0	20.0
2	15.0	10.0	41.0
14	12.0	11.06	26.0
44	2.0	15.0	29.25

predicted Value

Step 6: Remove all col3 values

(8)

	R&D	Admin	MS
21	8.0	15.0	30.0
37	23.14	5.0	20.0
2	15.0	10.0	41.0
14	12.0	11.06	26.0
44	2.0	15.0	NAN

	R&D	Admin	MS
21	8.0	15.0	30.0
37	23.14	5.0	20.0
2	15.00	10.0	41.0
14	12.0	11.06	26.0
44	2.0	15.0	31.56

Predicted value

Iteration 0 (Mean)

	R&D	Admin	MS
21	8.0	15.0	30.0
37	9.25	5.0	20.0
2	15.0	10.0	41.0
14	12.0	11.25	26.0
44	2.0	15.0	29.25

Iteration 1 (Predict)

	R&D	Admin	MS
21	8.0	15.0	30.0
37	23.14	5.0	20.0
2	15.0	10.0	41.0
14	12.0	11.06	26.0
44	2.0	15.0	31.56

Difference

	R&D	Admin	MS
	0.0	0.0	0.0
	13.89	0.0	0.0
	0.0	0.0	0.0
	0.0	-0.19	0.0
	0.0	0.0	2.31

Iteration 1 (Predict)

	R&D	Admin	MS
21	8.0	15.0	30.0
37	23.14	5.0	20.0
2	15.0	10.0	41.0
14	12.0	11.06	26.0
44	2.0	15.0	31.56

Iteration 2 (New Predict)

	R&D	Admin	MS
21	8.0	15.0	30.0
37	25.78	5.0	20.0
2	15.0	10.0	41.0
14	12.0	11.22	26.0
44	2.0	15.00	31.56

Difference

	R&D	Admin	MS
	0.0	0.0	0.0
	0.64	0.0	0.0
	0.0	0.0	0.0
	0.0	0.16	0.0
	0.0	0.0	0.0

\* This process → continuous → Difference (all 0.0)

In Iteration 2 we do same process → first Remove the value into NAN. And then predict the value and then find the difference bet<sup>n</sup> Iteration 1 and Iteration 2. Until all difference be 0.0.