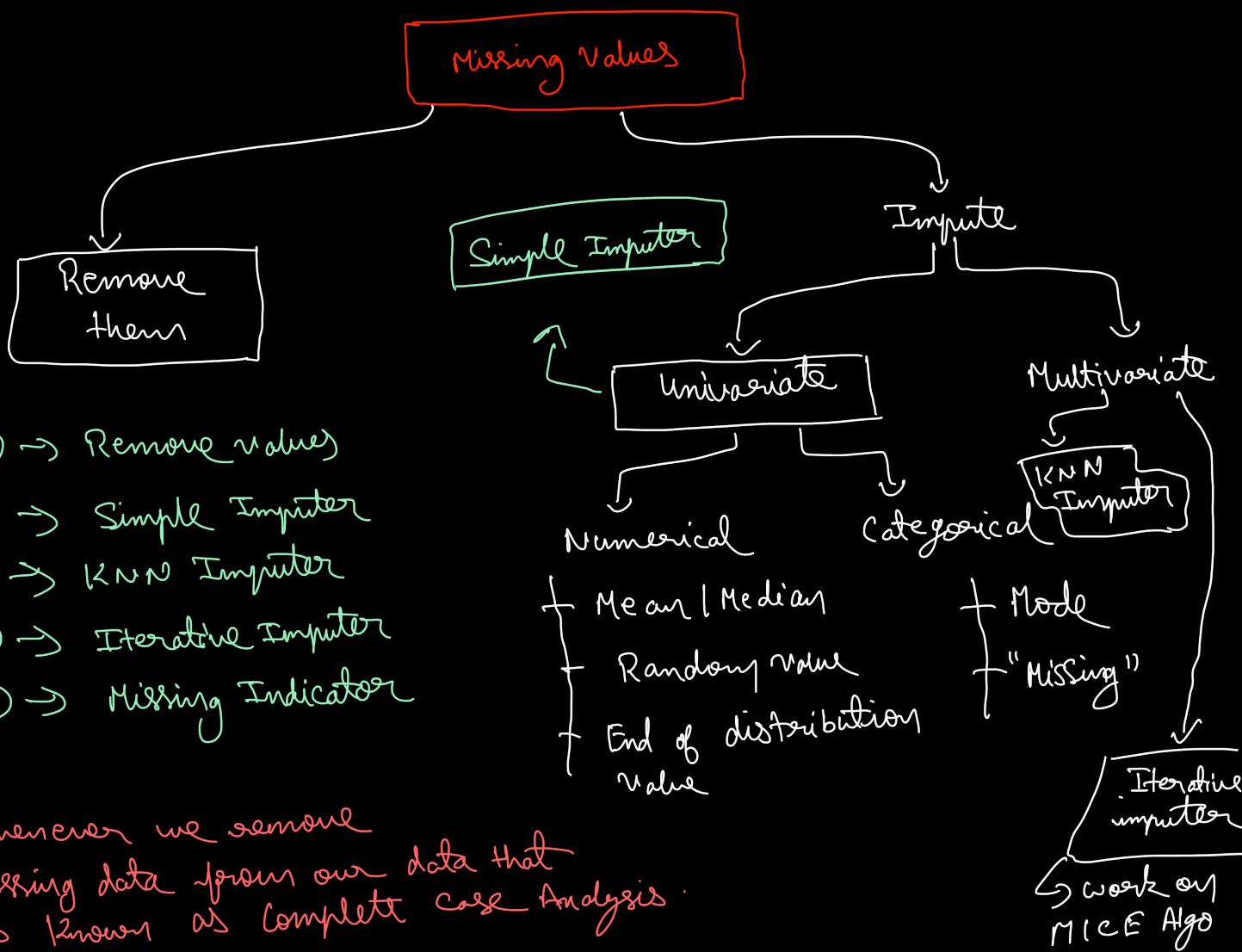


Handling Missing Data



* whenever we remove missing data from our data that is known as Complete Case Analysis.

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

Rows →

Columns →

→ 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

1) Missing Completely at Random

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).

MCAR:
 Missing completely
 At Random

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

When to use CC A?

1. Data should be MCAR.
2. Missing data should less than or equal to 5%.
 - * also this is not rule; it depends
 - * sometimes if more missing values are there, instead of removing records we can remove that particular column.

Handling Missing data [Numerical Data | Simple Imputer]

→ Mean / Median Imputation

Age

27

32

(Na)

27

Mean / Median

Benefit

- i) Simple to implement

Disadvantage

1. change in distribution shape
2. outliers occurs
3. covariance / correlation with other variables changes

When to use : When Data is MCAR
Missing Data < 5%.

→ Arbitrary Value Imputation

- | | |
|------------------------------|---------------------|
| <u>Benefit</u> | <u>Disadvantage</u> |
| i) Easy to apply | Same as above |
| ↓ (i) when dat is not MCAR . | |

→ End of Distribution Imputation (extension of arbitrary value imputation)

* In arbitrary value imputation, it is difficult to find best random value to impute. That's why end of distribution imputation comes into picture. Where we pick value from end of distribution.

If data is normally distributed then : We replace missing values with $(\text{mean} + 3\sigma)$ or $(\text{mean} - 3\sigma)$

If data is Skewed: $Q_1 - 1.5 \text{ IQR}$
(IQR Proximity) $Q_3 + 1.5 \text{ IQR}$

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

Benefits : easy to use Disadvantages : Same as above

↳ When to use : when data is not MCAR .

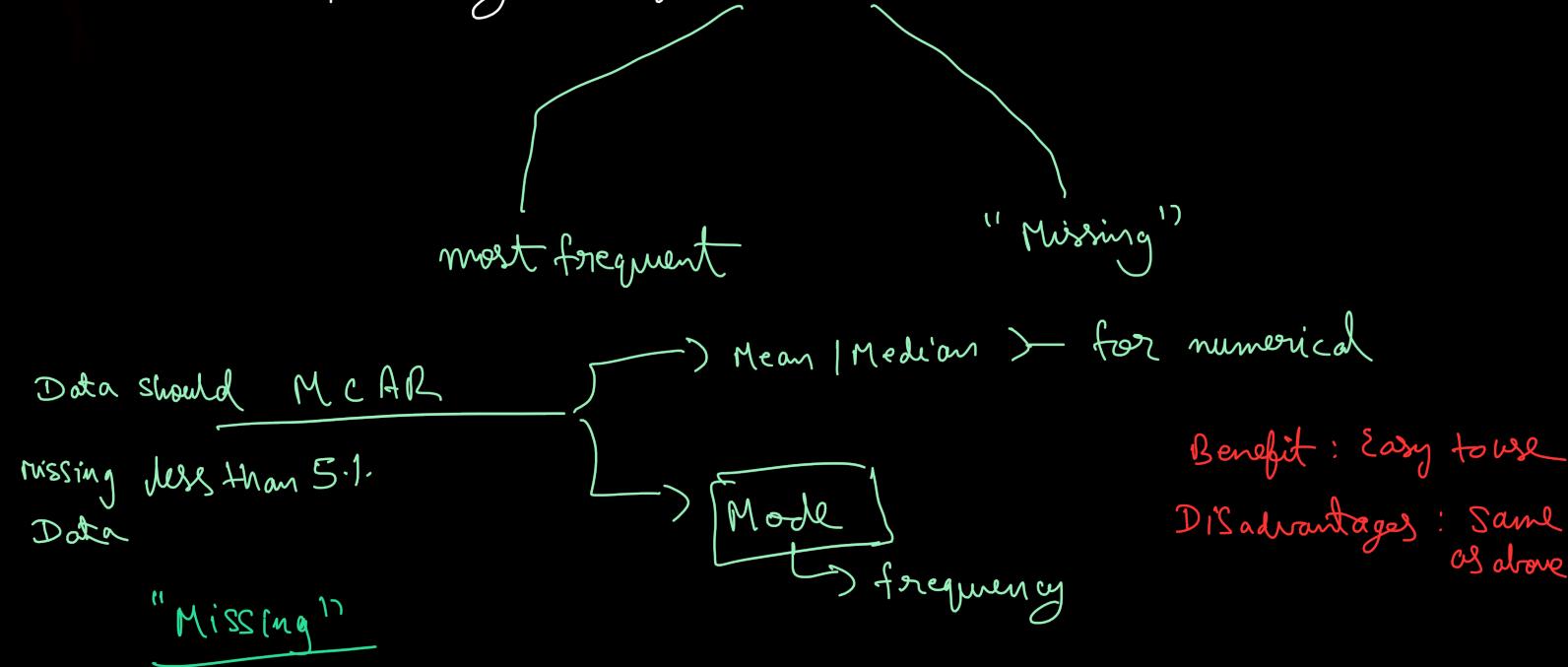
→ Random Sample Imputation



→ Automatically Select best imputation technique



Handling categorical Missing Data



- # Sometime we make a new category such as "missing" and fill NA with that category.
- # So our ML algo come to know that this is missing and learns how to handle.
- # we do this when missing data is not MCAR.
- # And missing data > 10% or more.

→ Here after this Randomness Comes into data because we are not imputing here we are creating a different category.

Missing Indicator | Random Sample | Imputation |

Random Imputation : Impute by any random value of that particular column.

Age
- 26
- 32
- 56
- 41
[Na]
=
[Na]

Random numbers

Benefit }
~~~~ Preserves the variance of the variable (But why?)  
and shape of distribution. And it good for  
linear algorithm.

Why: Data points which are present large in numbers,  
their probability to get selected randomly is high  
that's why shape of distribution and variance  
get preserved.

- ~~~~ well suited for linear models as it does not  
distort the distribution, regardless of the % of NA
- ~~~~ No impact over outliers

### Disadvantages

- Memory heavy for deployment, as we need to store the  
original training set to extract values from and  
replace the NA in coming observations.
- ~~~~ Covariance with other variables get disturbed  
because of randomness.
- If more data is missing then it is not good  
technique. In case of categorical variable, this  
can change the ratio of categories or even  
in case of numerical variable also variance  
can be changed.

## # Missing Indicator

↳ Make completely new column for every column which has missing data . In which we keep two values True and False . True for those record which has value and False for missing record .

| Age | Fare | AgeNa |
|-----|------|-------|
| 22  | 32   | F     |
| 41  | 35   | F     |
| Na  | 41   | T     |
| 62  | 32   | F     |

- \* There is assumption with this technique that Machine Learning algorithm learns to differentiate between missing values and non-missing values . And because of this algorithm performs better some time .
- \* Atleast we should try it once .

## # Automatically Select Value for Imputation

- ↳ By Grid Search CV
- ↳ Go to Code

# Univariate Imputation : when we fill missing values by taking help of that particular variable only .

# Multivariate Imputation : when we fill missing values of any column by taking help of others variables in data .

## K NN Imputer

## Multivariate Imputation

↳ use KNN algorithm

↳ To calculate distance : It uses non-Euclidean distance

→ Imputation for completing missing values using k-Nearest Neighbors.

→ Each sample's missing are imputed using the mean value from n neighbors nearest neighbors found in the training set. Two samples are close if the features that neither is missing are close.

Non Euclidean distances : Calculate the euclidean distances in the presence of missing values.

When calculating the distance between a pair of samples, this formula ignores feature coordinates with a missing value in either sample and scales up the weight of the remaining coordinates:

$$\text{dist}(x, y) = \sqrt{\text{weight} * \sum \text{distance from present coordinate}}$$

where, weight = Total number of coordinates / Number of present coordinates.

for example, the distance between [3, na, na, 6] and [1, na, 4, 5] is

$$\sqrt{\frac{4}{2} (3-1)^2 + (6-5)^2}$$

4 = Total number of coordinates

2 = number of present coordinates

\* If all the coordinates are missing or if there are no common present coordinates then NaN is returned for that pair.

## Advantages & Disadvantages

- ↳ More accurate
- ↳ More number of calculation
- ↳ Prediction slow [Every time when we do prediction, to calculate missing values we have to use the training data kept on Server, that is also a problem.  
(Storage problem)]

## # Iterative Imputer / MICE

MICE stands for Multivariate Imputation by Chained Equations

### Assumptions

#### category of missing data

- MCAR
  - MAR  $\Rightarrow$  Data should Missing At Random.
  - MNAR
- This is assumption

### Advantage & Disadvantages

- ↳ accurate

- ↳ Slow
- ↳ need to store training data on Server to impute when we do prediction

$\Rightarrow$  How it works?

- Fill all the NaN values with mean of respective Column.
- Remove all Col1 missing values (starting from left)
- Predict the missing values of Col1 using other columns By using any machine learning algorithm. Treat Col1 missing as output feature and rest as input features.
- do same for all columns in which missing values are there.

- Then find difference final data and data with mean imputed. Our task is to reduce this difference close to zero.
- After one iteration, we take first iteration original output as base and repeat the all above mentioned steps. Till we get difference close to zero.
- Go to code for more understanding or read some blog over google.