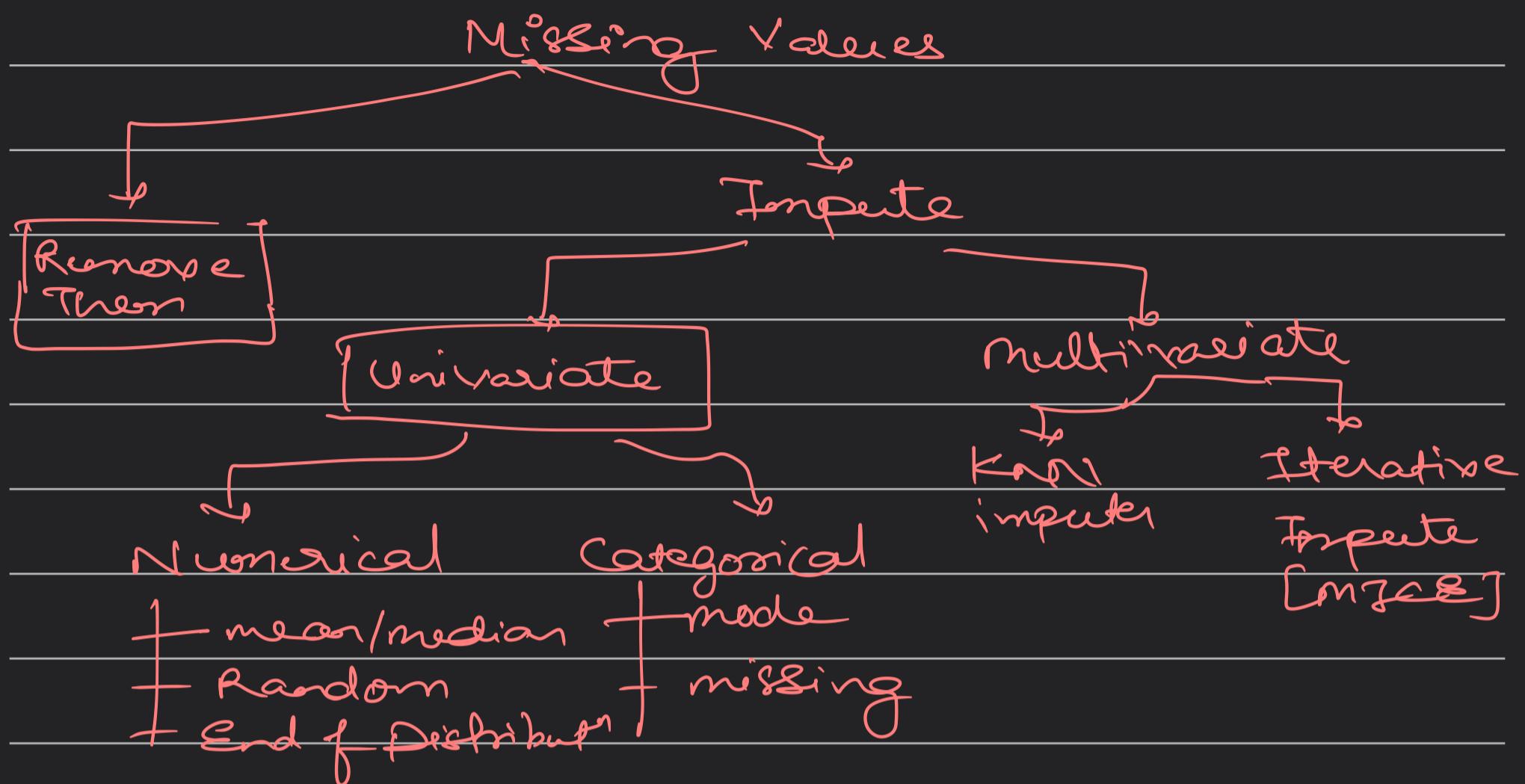


Handling missing data - Imputation
filling missing data.



Removing the Data (Complete Case Analysis)
Approach

Also called List-wise deletion of cases
consist in discarding observations where values
in any of the variables are missing.
→ Columns → Rows

Complete case analysis literally means analysing
only those cases for which information is
present in each column.

No missing value in any of the column
hr needs.

Assumption for CCA

- Data should be missing completely at random.

for ex 1000 values so are empty
but these 50 values should be
randomly missing. No starting, middle
or end, randomly only.

Advantages:

- ① easy to implement, no data manipulation is required
- ② Data before and after removal is same

Disadvantages

- ① can exclude , large part of observations
- ② Excluded values could be useful for analysis
- ③ When using models in production, model will not know how to handle missing data.

• When to use CCA?

We use CCA when missing data is less than 5% and second data should be missing at random.

What if 95% of missing data?

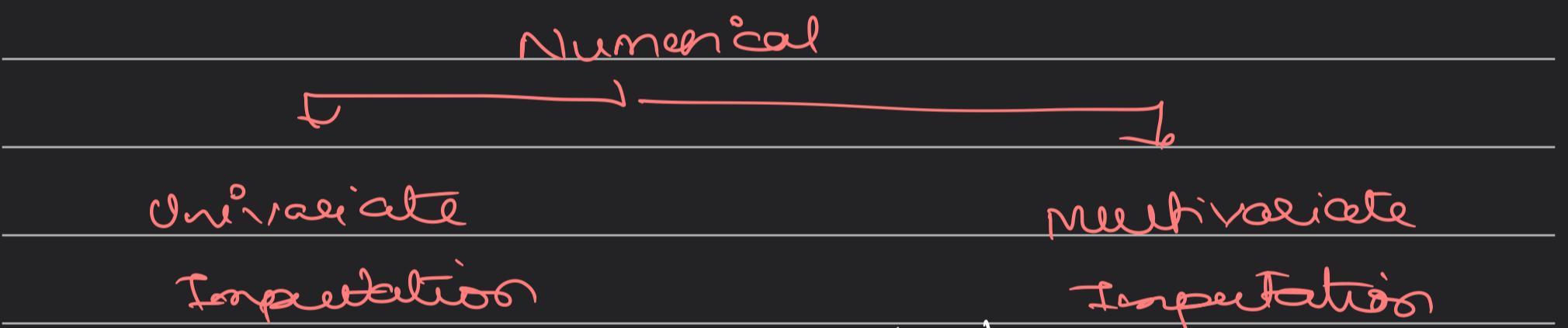
ii) We drop those columns where we have missing 95% data in columns.

• Thing to remember when dropping their values?
nestagoon

Well the distribution[↑] before & after
should some what match.
This tell us that MCRA existed and
as CCA is successful for that
variable.

For categorical data the ratio of different categories in variable should be maintained before & after CCA.

Handling Missing Numerical Data.



Let's say we have $\begin{array}{|c|c|c|c|} \hline 1 & 2 & 3 & 4 \\ \hline \end{array}$ and column 2 has missing values, and we use some statistical method (mean, max, min, etc) to fill the values in column "2" using values present in column 2 we are doing "univariate Imputation"

Similarly if we use multiple columns to fill missing values in a particular column its called multivariate imputation

• Techniques for Univariate Imputation

a) Mean / median Imputation

When to use mean or median

If Data is normally distributed

use mean or median any of them.

If data is somewhat skewed we

use median data for Imputation

It's very simple to use, but not used in production.

Not reliable if missing data > 10%

→ changes the shape of distribution.

→ we also get extra outliers

→ changes or creates changes in co-variance or correlation btw variables.

When to use

1) Data is MCAR

2) less than 5% data is missing

b) Arbitrary value Imputation

Missing values gets replaced with word "Missing", this help my

model to differentiate between

observation's that had data and the ones that didn't.

In numerical data we use a "data end value" that is not present in our variable

The whole idea is create difference observation that had data with others that didn't have it.

PDF graph gets distort.

variable gets effected

Covariance & correlation gets effected.

• End of Distribution Implications

An extension of last technique, we just use one "end value" either from starting or ending to find the value

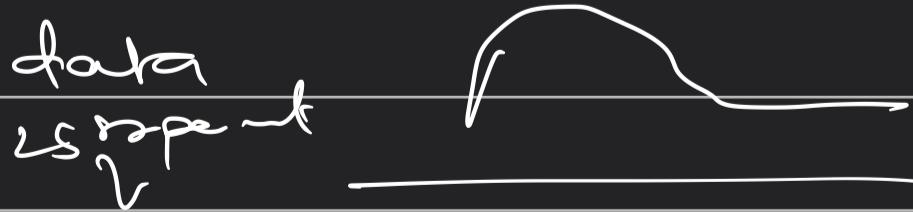
for a normally distributed data



$$\text{Value} = \text{mean} + 3\sigma$$

$$\text{Value} = \text{mean} - 3\sigma$$

for skewed data



$$\text{Value} = Q_1 - 1.5 \times IQR \rightarrow \frac{(Q_3^{th} - Q_1^{th})}{(Q_3 - Q_1)}$$

$$\text{or Value} = Q_3 + 1.5 \times IQR$$

\downarrow
75th percentile

Handling missing categorical data.

most frequent value



missing Category
Imputation

Missing Frequent Values: → we use mode of the data categories, here we have assumption

- 1.) Missing Category at Random
- 2.) Mode should have higher frequency compared to other data in category

changes the distribution of the data.

missing Category Imputation

Here we create another variable called "Missing" and put the missing categories here for observations

Done when data missing > 10%. Also if data is "not" missing at random.

Just a randomness gets created in data.

* Technique used for Both Univariate
Imputation of Numerical and Categorical
Data.

Random value
Imputation

missing
Indicator

Automatic
selection of
Parameters.

Random value Imputation

Filling missing values through random
members, we randomly select
from the members present in our
variable, randomly.

distribution

Helps preserve the variance of
variable. We do this ^A using
pandas.

* Memory heavy for deployment as we
need data for dataset to set and
extract values & replace NA in
coming observation for prediction

Good for linear models as it preserves
the distribution.

Not good for Decision ~~Based~~ tree based
algorithms.

ster

Co-variance with data gets a bit distributed.

• Missing Indicator

For any column that has a missing data, we create a new column for each such column and keep value True or False for every observation

Age	Fore	Age - Missing
37	22	False
21	35	False
NA	Up	True
40	100	False

already implying

feels

One the mostly used technique, but
we have estimate that the ML algo
learns to different data with missing
observations & non-missing observations

• Auto Select value

We use something called Grid search
CV where scikit learn uses auto
matically selects different imputation
technique and select the best one.

Multivariate Imputation

KNN Imputer

Iterative Imputer

KNN Imputer

works on the KNN Algorithm, you are like your neighbours, your behaviour is in accordance with your neighbour.

We feel a missing value in row by taking value from row which is most similar to our missing value row.

This similarity is calculated using Euclidean distance.

Each row can be used as a coordinate where distance of other rows is from missing row is calculated & data from closest row is used to fill the missing values.

Here K represent number of neighbour
Ex $K = 3$ nearest ones

Missing Value \rightarrow Value row 1, + value row 2 + value row 3

we screen algo to get distances.

Adv & Dis Adv *

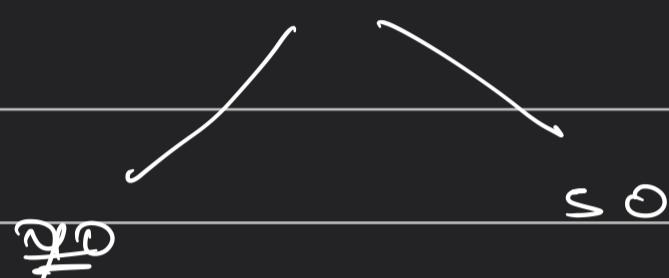
- 1) More Accurate
- 2) more no. of calculation
- 3) Deploying on production is heavy

as number of value gets calculated from whole X-train data.

v) Good for smaller & medium size dataset.

concept of uniform & Distance

we calculate weight either by taking mean $K=2$



$$\text{Uniform} \quad \frac{70 + 50}{2} = 60$$

Distance = multiply by reciprocal of distance with value of data

Every one has same weightage

closest neighbour has more value given, farthest one has less.

• Iterative Impute (MICE)

This technique use something called M.I.C.E algorithm

Multivariate Imputation By Chained Equations (MICE)

We use MICE under some assumption

few categories of missing data.

- 1) Missing Completely at Random
- 2) " at Random
- 3) missing not at Random
 - In some people didn't filled data
 - Data removed consciously.

for few data couldn't get collected

Assumption

we use mice in missing at random

This technique is quite accurate,

but is slow as we MLE Algo to

predict missing values

we need to put training data on server. Problematic for bigger data.

How mice worse?

worse on input variables⁶

Step 1 → replace NAN values with mean values of each column respect of

Step 2 → move left to right and whenever we had missing values we replace again with NAN.

Step 3 → we use ML algo to predict

this missing values for that column using other columns present in the data.

Step 4 → same from left to Right
fill missing values with NaN and repeat step 3

Step 5 → After getting done for each row. we calculate difference for each iteration

Before & After. [mean - prediction]

Step 6 → we repeat this iteration until all difference becomes 0. At each iteration, the previous iteration becomes base,

when difference becomes 0 the ML algo cannot improve more and we got most accurate value to replace this missing value.

