

Data Preprocessing

: Make changes before giving to the Model

KM	M
5	5000
6	6000
7	7000

age	Salary
59	48000
21	100000

Range
↓
Min
Max

$$\frac{20 - 60}{\text{Min} - \text{Max}} \quad \frac{2K - 2CR}{\text{Min} - \text{Max}}$$

Scaling : Will Bring Values in same Range {PCA}

{Image} $\xrightarrow{\text{Normal Distribution}}$ Standardization

KM	M
5	5000
6	6000
7	7000

$$\frac{x - \text{amin}}{\text{amax} - \text{amin}}$$

$$\begin{array}{|c|} \hline \text{Mean} = 0 \\ \text{std dev} = 1 \\ \hline \end{array}$$

$$Z = \frac{x - \mu}{\sigma}$$

$$\frac{7000 - 5000}{7000 - 5000} = \frac{2000}{2000} = 1$$

KM	M
0	0
0.5	0.5

$$D = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

0	0
0.5	0.5
1	1

$$\left\{ \begin{array}{l} \text{Mean} \rightarrow 5.5 \\ \sigma^2 = \frac{\sum (n - \bar{n})^2}{n} \\ = \frac{(1-5.5)^2 + (2-5.5)}{10} \end{array} \right.$$

$$\left\{ \begin{array}{l} \sigma^2 = 9.6 \\ \text{Std Dev} = \sqrt{9.6} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{Mean} = 0 \\ \text{Std Dev} = 1 \end{array} \right.$$

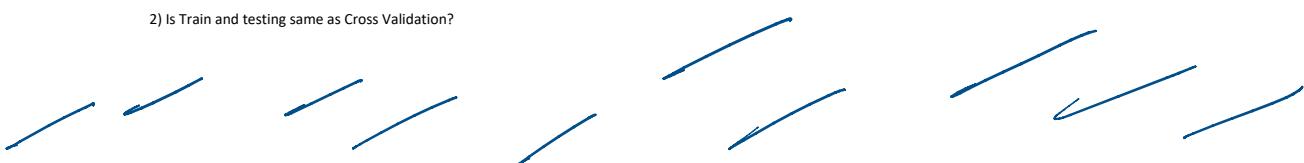
1	-2	5
2	-1	
3	0	
4		
5		
6		
7		
8		
9		
10	-3.7	

$$\begin{aligned} & \frac{1-5.5}{\sqrt{9.6}} \\ &= \frac{2-5.5}{\sqrt{9.6}} \\ & \vdots \\ &= \frac{10-5.5}{\sqrt{9.6}} \end{aligned}$$

what is difference in regression and logistics
Balanced or Imbalanced, train n test
In sigmoid curve

```
ary for linear regression.  
>>> model = sm.OLS(y, x.assign(const=1)) #statsmodels  
>>> results = model.fit()  
>>> print(results.summary())  
Can you please explain all the summarised parameters and their uses. Also, is there any same summary  
function in sklearn module??
```

2) Is Train and testing same as Cross Validation?



Encoding { Done

Encoding } Done

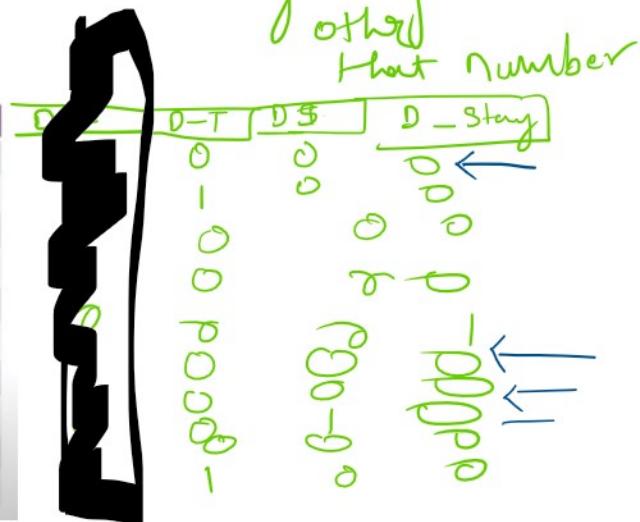
Categorical } String

Computer cannot
understand

anything
other

than number

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis



Categorical

Dummy
encoding

One hot
encoding

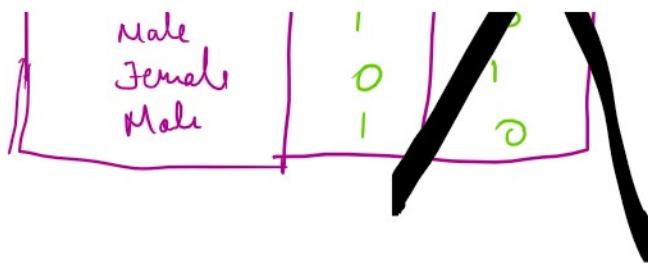
Label
encoding

Gender

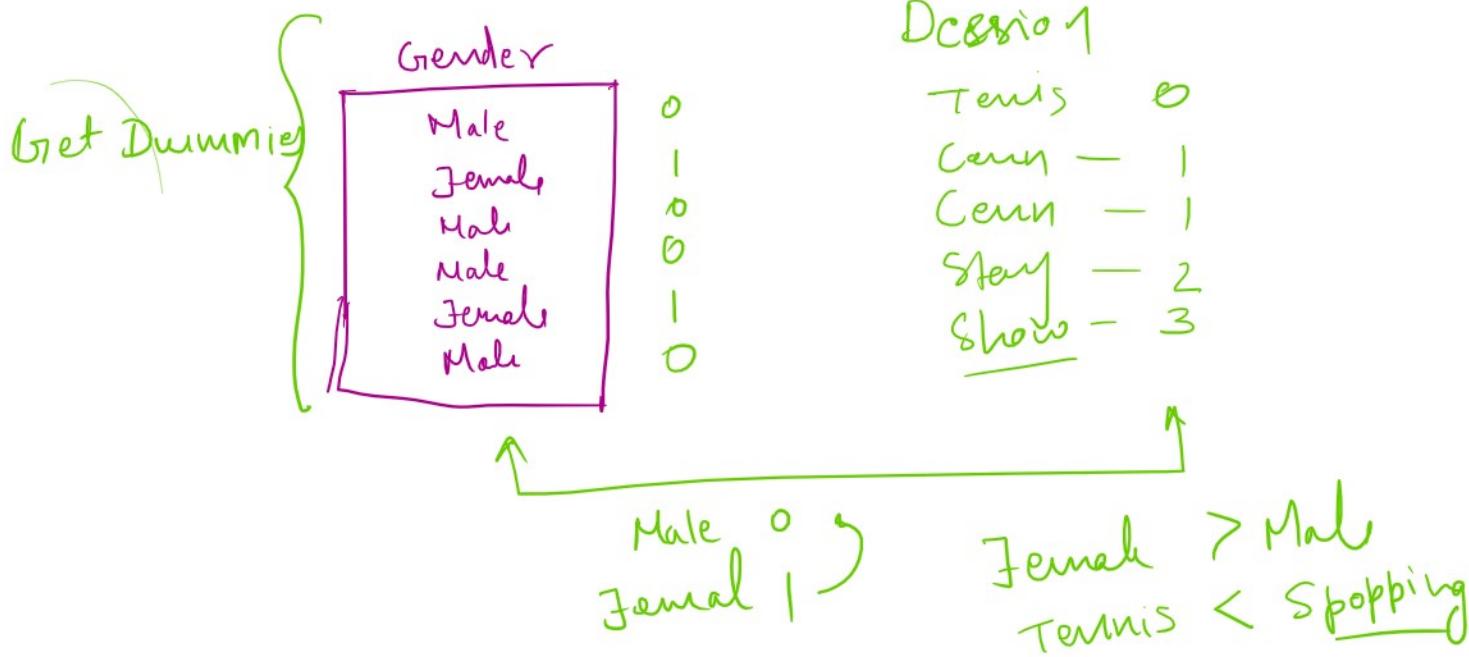
Male-Gender Female-Gender

Gender	Male-Gender	Female-Gender
Male	1	0
Female	0	1
Male	1	0
Male	1	0
Female	0	1

→



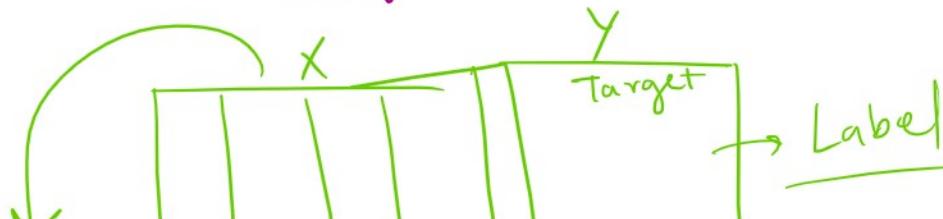
Label encoder

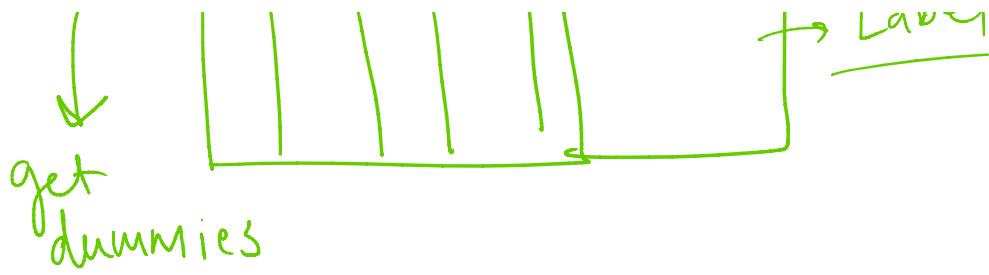


→ Order

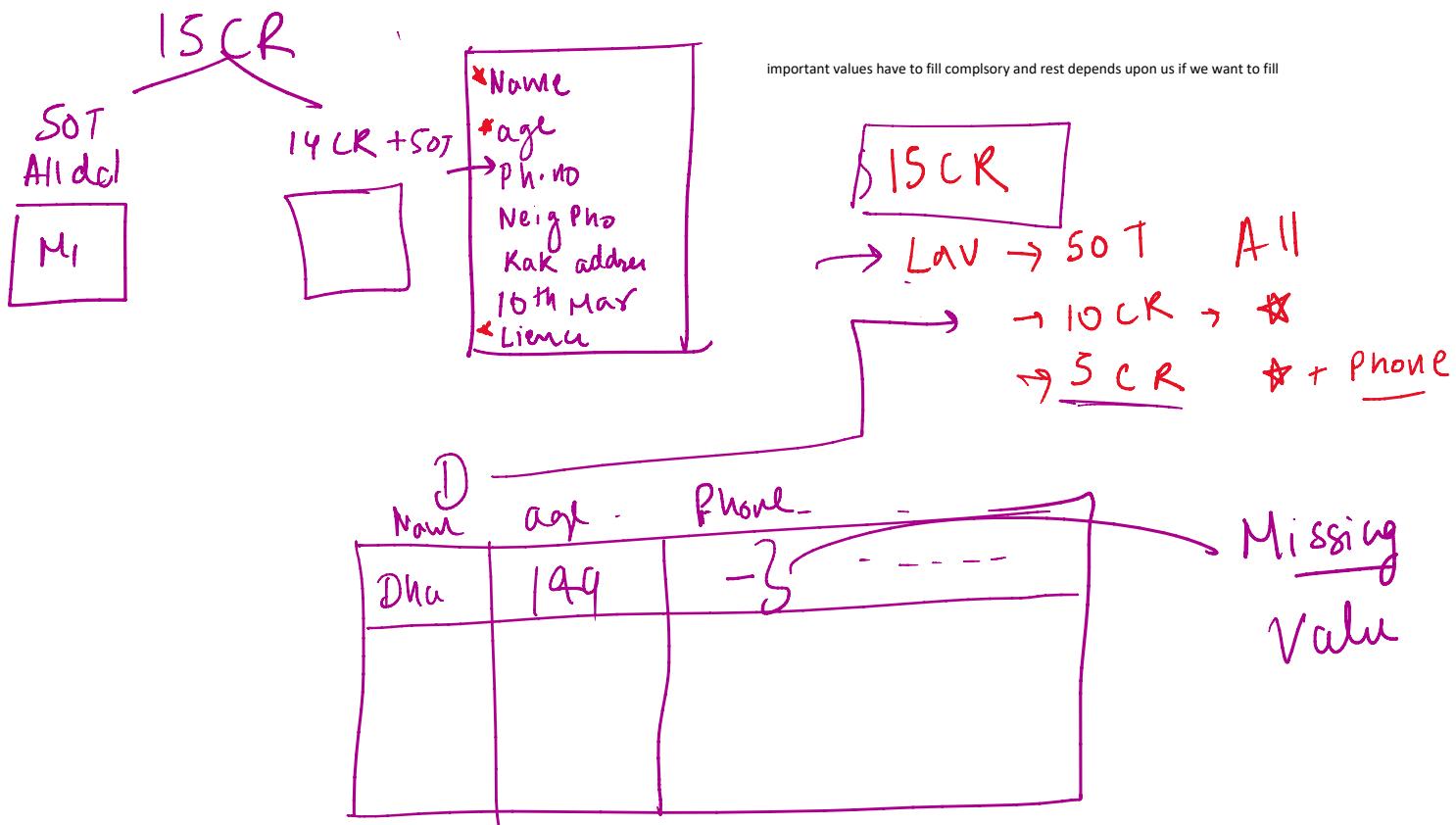
PHD > MTECH > BTech

Target → Label



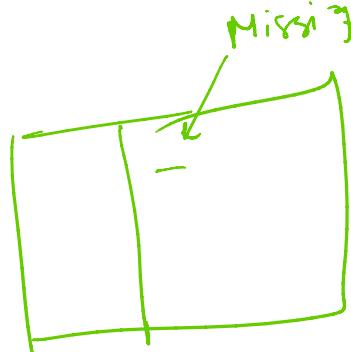


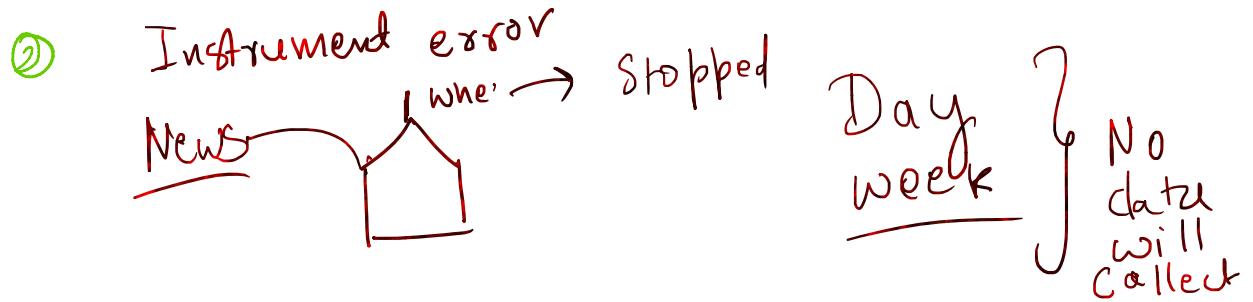
Missing Value



① Human mistake

Rubbish } wrong
→ Noise }





Models which cannot work with Missing Value

* Treatment

Row

SKNo	Name	Occ	Age	Gendry
1	Gajju	DS	27	M
2	Raju	SE	-	M
3	Maje	-	23	F
4	Kohar	-	-	-
5	Ohar	Java	17	F
6	Raju	CS	21	M

①

Delete the missing value

Row

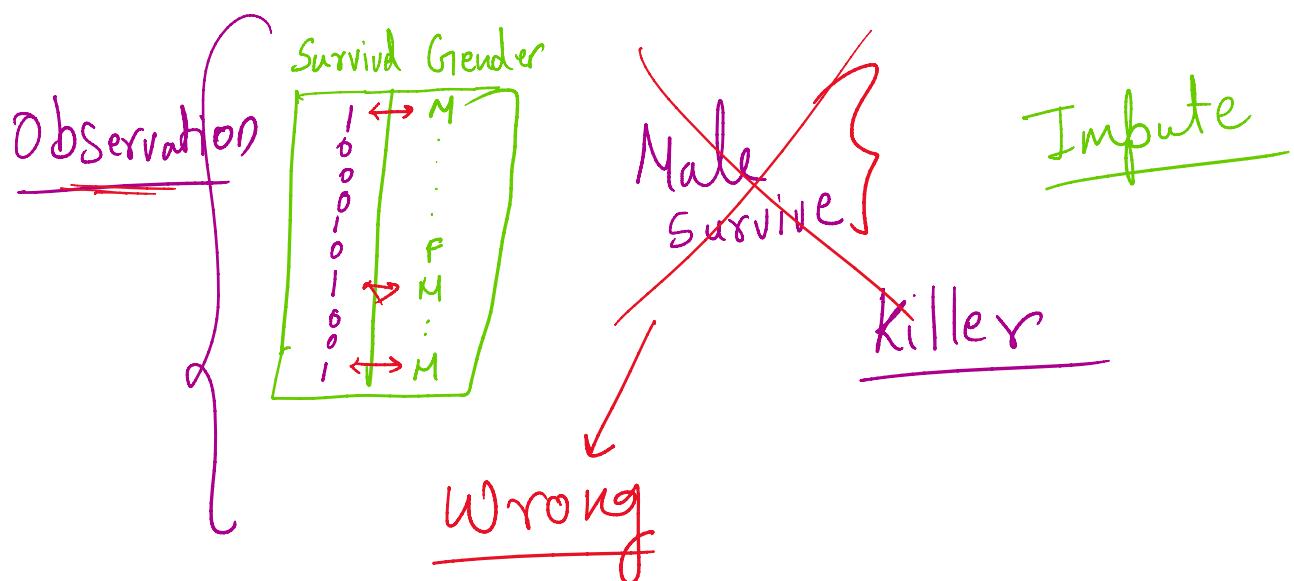
Column

Disadvantage

• Data loss } To be very very careful before

1000 ~ careful before
dropping the row

: Row if has many values missing
then drop the row



Drop Column 40%
10 6%

→ Impute

→

100

26

21

27

26

23

27

24

130

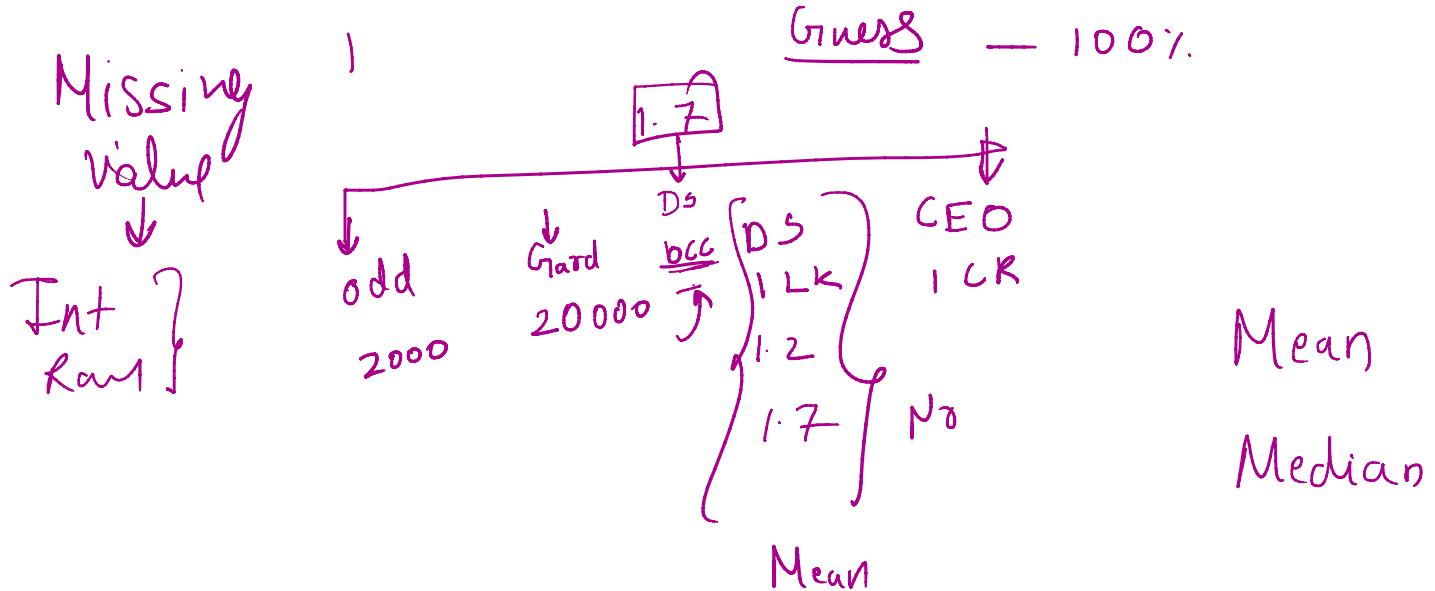
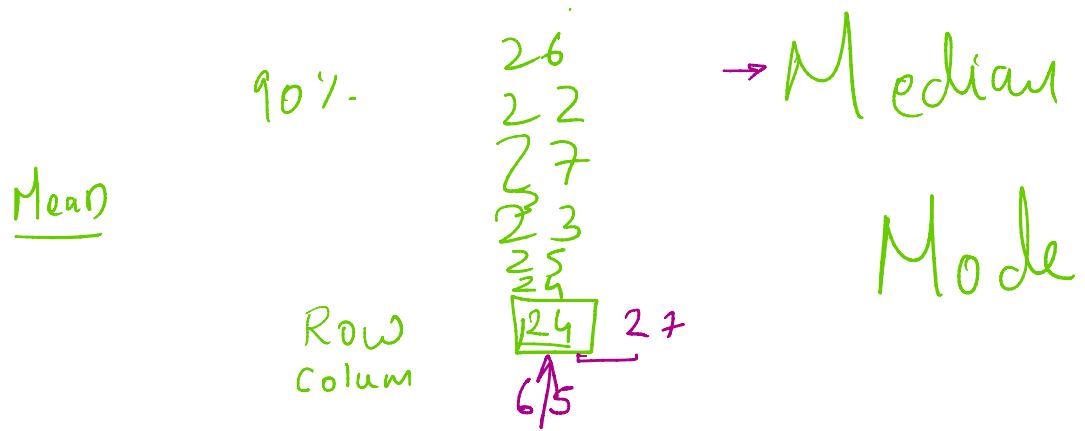
13

90%

26

M
ean

→ Median



Mean —

Median — outlier

GIGO

Mode

Category

Blue
white
green

orange
black

white
white

Mode

white

9° 45