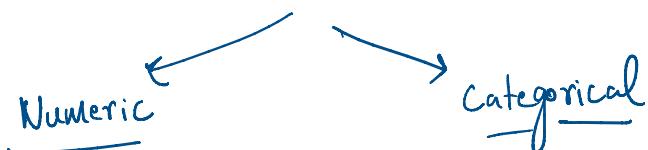


Scaling

Data preprocessing

Make changes / Manipulate the data before giving it to the model



Scaling Numerical

KM

5.

6

7

M

5000
6000
7000

	age	Salary
20	50	45000
60	59	45100

$$\frac{5000 - 75}{5000 - 5} = s$$

Distance based algo
 → Linear
 → Logistic



$$\sqrt{(59 - 50)^2 + (45100 - 45000)^2}$$

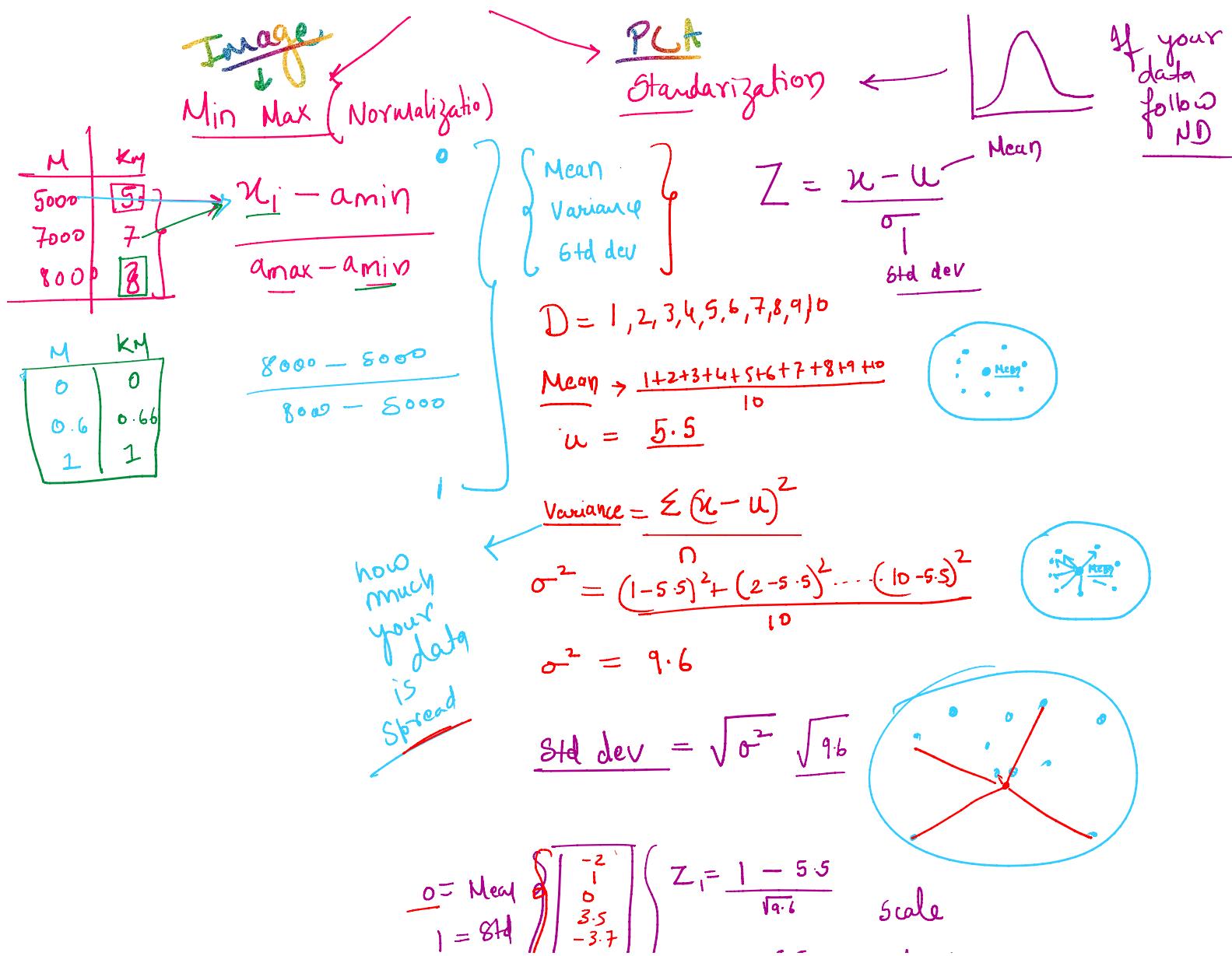
$$\{ = \sqrt{19^2 + 100^2}$$

Model will give more Imp

<u>age</u>	<u>Salary in Thousand</u>
50	45
59	45.1

$$\frac{\sqrt{(59-50)^2 + (45.1-45)^2}}{\sqrt{9^2 + 0.1^2}}$$

Scaling : will Bring all the values in same range



$$\begin{aligned}
 & \text{Mean} = 0 \\
 & \text{Std} = 1 \\
 & \text{Scale} \\
 & \text{Mean} = 0 \\
 & \text{Std} = 1
 \end{aligned}$$

$\frac{0 - 0}{\sqrt{9.6}} = 0$
 $\frac{3.5 - 0}{\sqrt{9.6}} = \frac{3.5}{\sqrt{9.6}}$
 $\frac{-3.7 - 0}{\sqrt{9.6}} = \frac{-3.7}{\sqrt{9.6}}$
 $Z_1 = \frac{1 - 0}{\sqrt{9.6}} = \frac{1}{\sqrt{9.6}}$
 $Z_2 = \frac{2 - 0}{\sqrt{9.6}} = \frac{2}{\sqrt{9.6}}$
 $Z_{10} = \frac{10 - 0}{\sqrt{9.6}} = \frac{10}{\sqrt{9.6}}$

Encoding

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

Computer cannot understand anything other than number

Encode the Categorical Data

Dummy / one hot encoding

Gender	Male-Gender		Female-Gender	
	Male	Female	Male	Female
M	1	0	0	1
F	0	1	1	0
H	1	0	0	1
M	1	0	0	0

Label encoding

Grade	ordinal	
	1	2
A	1	2
B	2	1
A	3	1
C	4	0

$2 > 1 > 0$

$\begin{cases} \text{Male} = 0 \\ \text{Female} = 1 \end{cases}$

$$\begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix} \quad \begin{matrix} 1 & 0 \\ 0 & 1 \end{matrix}$$

Nominal

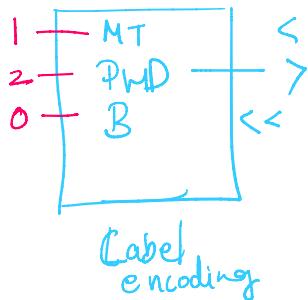
Black = Green

Weather W-S W-R W-W W-D

S	1	0	0	0
R	0	1	0	0

W	0	0	1	0
D	0	0	0	1
S.	1	0	0	0
R				
D				
S				

	Color	E	P	S	CS
M	B1	-	-	-	-
D	g2	-	-	-	-
A	y3	-	-	-	-
E	B4	-	-	-	-



- Outlier
- Confusion Matrix
- • Missing Value
- Feature Selection

Coding

↳ Dataset

- Reg → Classific
- Miss
 - Std
 - out
 - ENC
- Linear DBC Logis Random

