

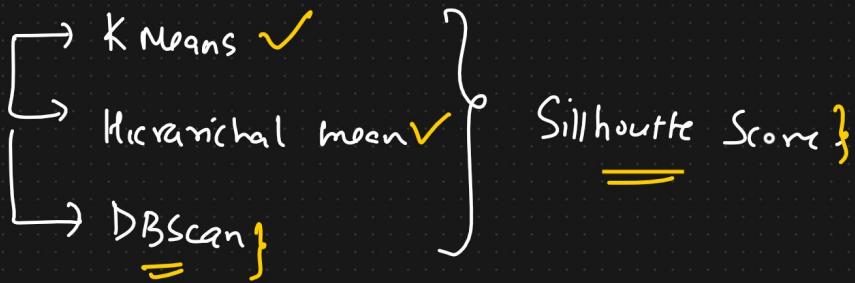
Agenda

① Gradient Boost ✓

→ CatBoost

② Xgboost ✓

③ Unsupervised ML

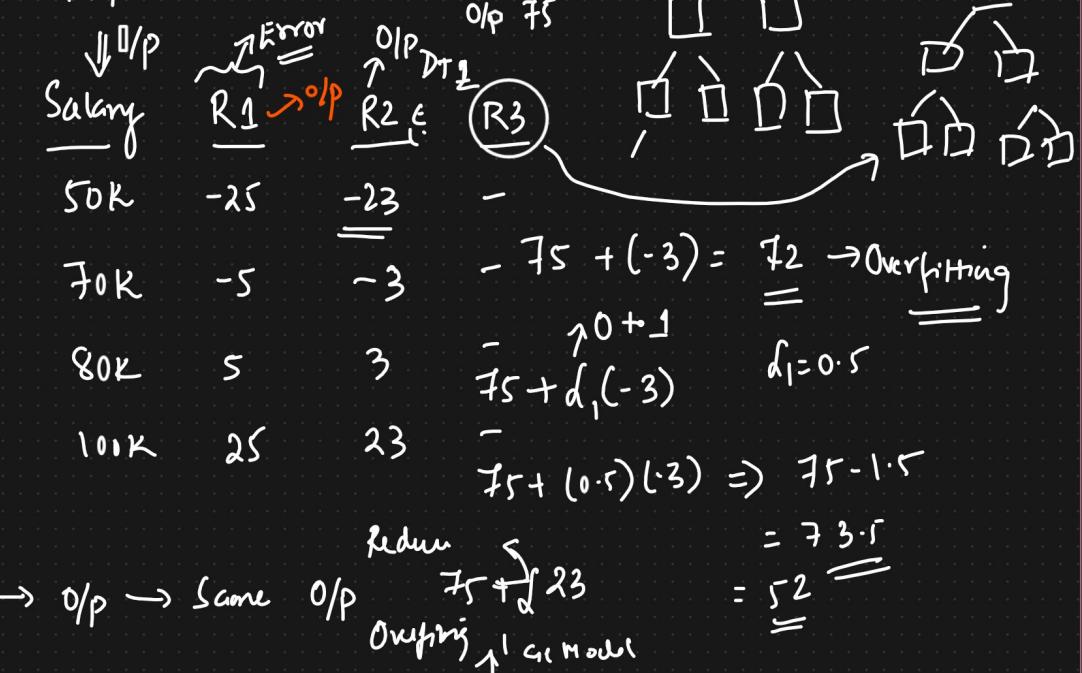


④ Gradient Boosting

→ Classification
Regression

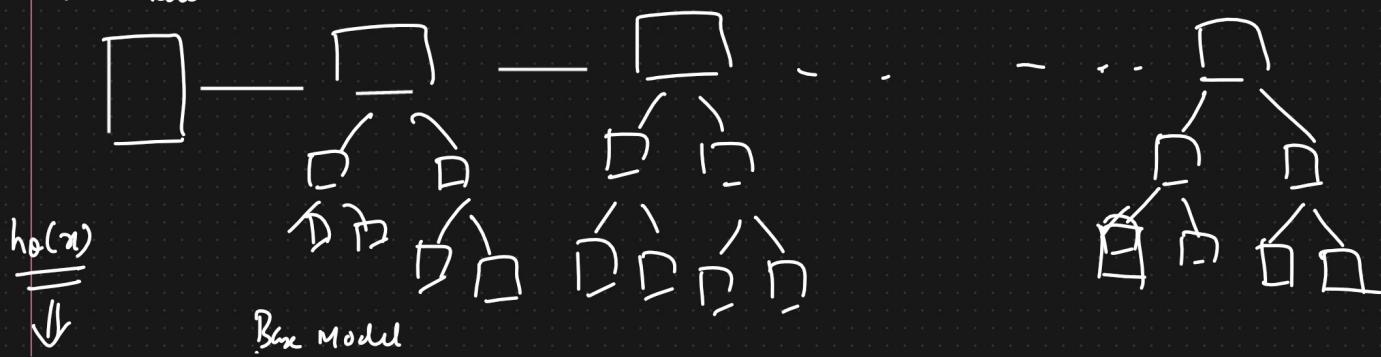
→ Exp → Degree

	2 BE	3 Masters	5 Masters	6 PhD	Test → 3.5 Masters
	50K	70K	80K	100K	
↓ I/P	-25	-5	5	25	
Error	R1 → 0/P	R2 → 0/P	R3 → 0/P	R4 → 0/P	
	-23	-3	3	23	



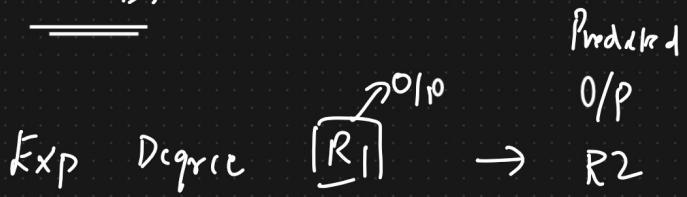
$$\text{Bar Model} = \frac{50 + 70 + 80 + 100}{4} = 75K$$

Bar Model



$$f(n) = h_0(x) + \alpha_1 h_1(x_1) + \alpha_2 h_2(x_2) + \dots + \alpha_n h_n(x_n)$$

Hypothesis Input DT



$$f(n) = \text{Base model} + \alpha_1(DT_1) + x_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)$$

↓

To reduce overfitting

Base model

f_1	f_2	O/P	R_1	O/P	R_2	O/P	R_3	O/P	R_4	O/P
50	-15	-16	50	-15	50	-16	50	-16	50	-16
60	-5	-4	60	-5	60	-4	60	-4	60	-4
70	5	7	70	5	70	7	70	7	70	7
80	15	13	80	15	80	13	80	13	80	13

$\frac{50+60+70+80}{4} = 65$

$\alpha_1 = 0.5$

$$65 + \alpha_1(-16) = 65 + 10 \cdot 5 (-16)$$

$$= 65 - 8$$

$$= 57$$

② Xgboost Classifier



<u>Salary</u>	<u>Credit</u>	<u>Approval</u>	<u>Res</u>	<u>(credit)</u> ✓
$\leq 50K$	B	0	-0.5	
$\leq 50K$	G	1	0.5	
$\leq 50K$	G	1	0.5	<u>0.250</u>
$> 50K$	B	0	-0.5	
$> 50K$	G	1	0.5	
$> 50K$	N	1	0.5	
$\leq 50K$	N	0	-0.5	$\frac{\text{Similary}}{1.75} = \frac{0.25}{1.75} = 0.142$

① Construct Base Model \rightarrow

$$\underset{\downarrow}{\text{O/P}} = \underline{\underline{0.5}}$$

② Construct Tree with Root Node.

$$\left[-0.5, 0.5, 0.5, -0.5, 0.5, 0.5, -0.5 \right]$$

Salary



③ Calculate Similarity Weight

$$= \frac{\left(\sum \text{Residual} \right)^2}{\sum \text{Pr}(1 - \text{Pr}) + \lambda}$$

$\lambda=0$

$$= \frac{\left[-0.5 + 0.5 + 0.5 - 0.5 \right]^2}{\left[0.5(0.5) + 0.5(0.5) + 0.5(0.5) + 0.5(0.5) \right]} = \frac{\left[-0.5 + 0.5 + 0.5 - 0.5 \right]^2}{\left[0.5(0.5) + 0.5(0.5) + 0.5(0.5) + 0.5(0.5) \right]}$$

$$= 0 + \frac{\left[-0.5 + 0.5 + 0.5 - 0.5 \right]^2}{\left[0.5(0.5) + 0.5(0.5) + 0.5(0.5) + 0.5(0.5) \right]} = \frac{0.25}{0.75} = 0.33$$

④ Gain

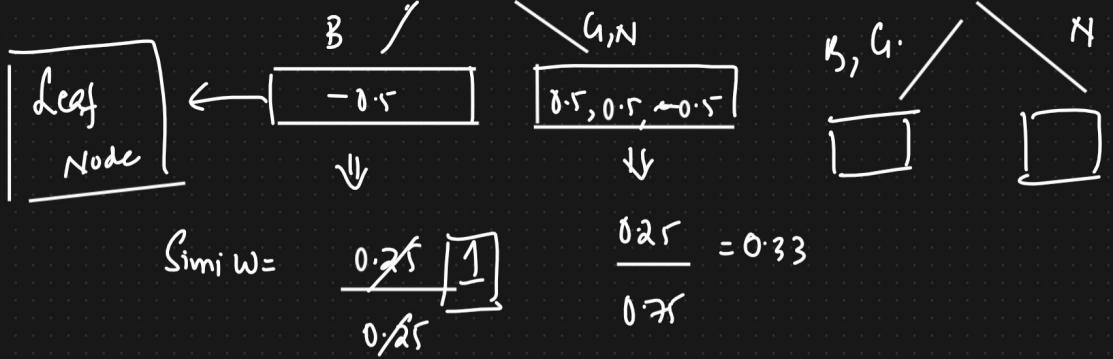
$$\text{Gain} = 0 + 0.33 - 0.142$$

$$\text{Gain} = \underline{\underline{0.188}}$$

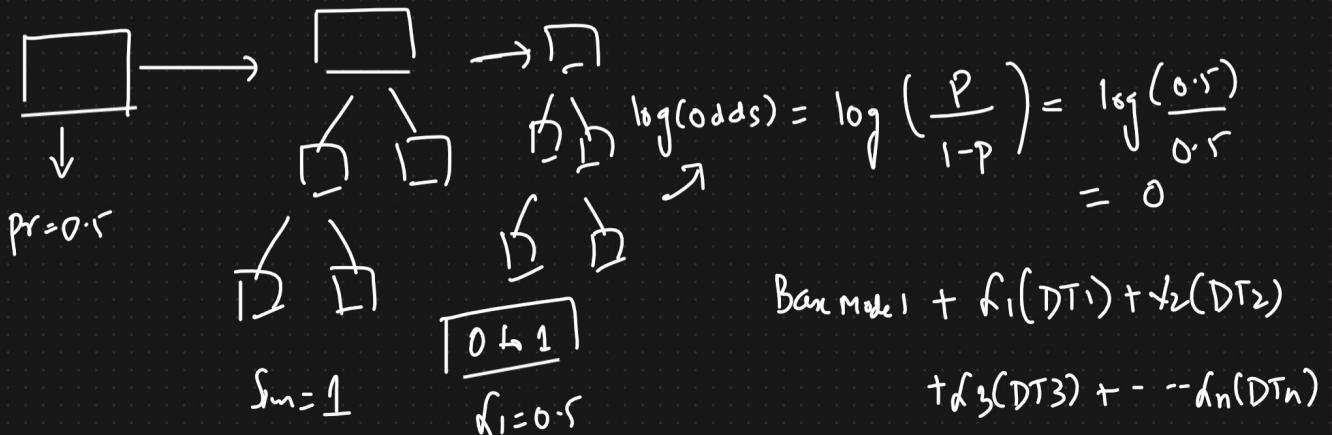
$[-0.5, 0.5, 0.5, -0.5]$

$[-0.5, 0.5, 0.5]$

Credit



$$\text{InfoGain} = 1 + 0.33 - 0 = 1.33$$



$$\underline{x_{\text{gboost}}} = \frac{\downarrow}{\Gamma} \left(0 + f_1(1) \right)$$

$$= \Gamma(0 + 0.5)(1)$$

Sigmoid Activation fn

$$\text{Sigmoid} = \frac{1}{1 + e^{-x}}$$

0-26

1

Classification problem

$$\hat{X}_{q\text{boost}}^{\text{base}} : \text{Base model} + \top \left(f_1(DT_1) + f_2(DT_2) + f_3(DT_3) + \dots f_n(DT_n) \right)$$

1

$$\log \text{loss} \rightarrow \log \left(\frac{P}{1-P} \right) \downarrow$$

Sigmoid Activation function.

f_1, f_2 Category $\xrightarrow{O/P}$ Continuous $\xrightarrow{O/P \rightarrow \text{Regression}}$

④ Logistic Regressor

$$\frac{(-1 - 1 + 1 + 1 + 1)^2}{5+1} = \frac{1}{6} \approx 0.166 \quad \boxed{\text{Base Model}} \rightarrow 51K$$

	Exp	Gap	Salary	Res1
→ 2	No	40K	-11	
→ 2.5	Yes	42K	-9	
3	No	52K	1	
4	No	60K	9	
4.5	Yes	62K	11	

$$\begin{aligned} \frac{1}{6} &= 0.166 \quad \boxed{\text{Exp}} \\ \lambda &= 1 \quad [-1] \\ \psi &= \boxed{[-9, 1, 9, 11]} \end{aligned}$$

↓

$$\begin{aligned} \text{Similarity weight} &= \frac{121}{1+1} = \frac{121}{2} = 60.5 \\ &= \frac{121}{5} = 24.2 \end{aligned}$$

$$\begin{aligned} ② \text{Similarity weight} &= \frac{(\sum \text{Residual})^2}{\sum p_i(1-p_i) + \lambda} \\ &\Downarrow \\ &= \frac{(\sum \text{Residual})^2}{\text{No. of residuals} + \lambda} \end{aligned}$$

$$\begin{aligned} \text{Gain} &= 60.5 + 24.2 - 0.166 \\ &= 84.534 \end{aligned}$$

No. of residuals + λ

$$\boxed{\quad} \longrightarrow \boxed{\text{Exp}} = 0.166$$

$$\begin{aligned} \text{Op} &= -10 \quad \lambda = 10 \\ \text{Op} &= \frac{-11 - 9}{2} = -10 \quad \boxed{[-11, -9]} \quad \boxed{[1, 9, 11]} \quad \lambda = 3 \end{aligned}$$

$$\text{Sim}(w+) = \frac{(-11 - 9)^2}{2+1}$$

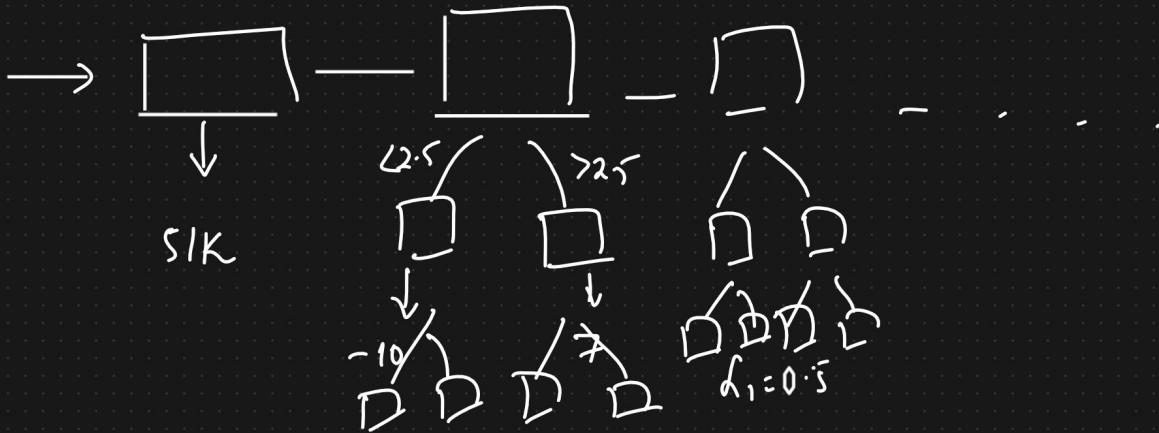
$$\text{Sim}(w-) = \frac{(1+9+11)^2}{4} = \frac{(21)^2}{4} = 110.25$$

Yes \swarrow No \searrow

$$-11 \quad \boxed{1} = \frac{D}{9} \cdot \frac{400}{3} = 133.3$$

$$\text{If gain} = 133.3 + 110 \cdot 25 - 0.166$$

$$\approx 243 =$$



$$51 + l_1(-10) + l_2() + l_3() + l_4() + \dots + l_n(D) =$$

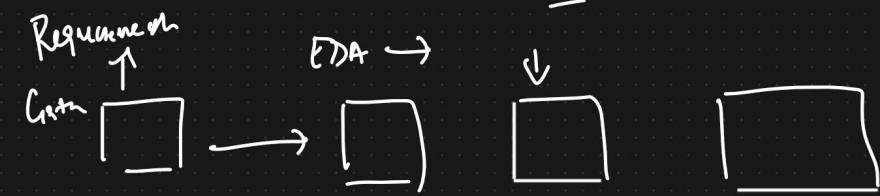
$$51 + 0.8(-10) = 51 - 8 = \underline{\underline{43}}$$

Interview → 2 interviews

logistic

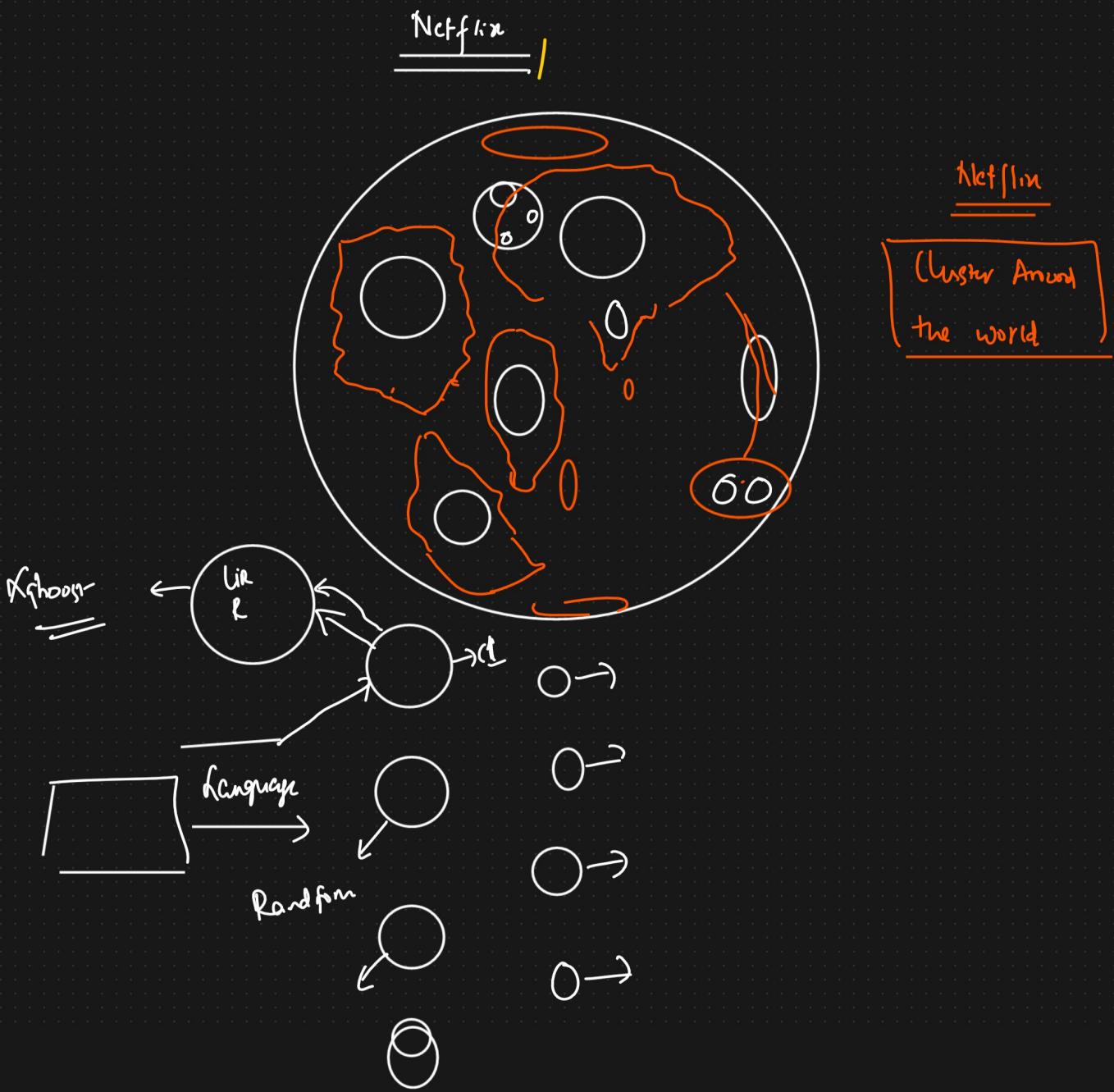
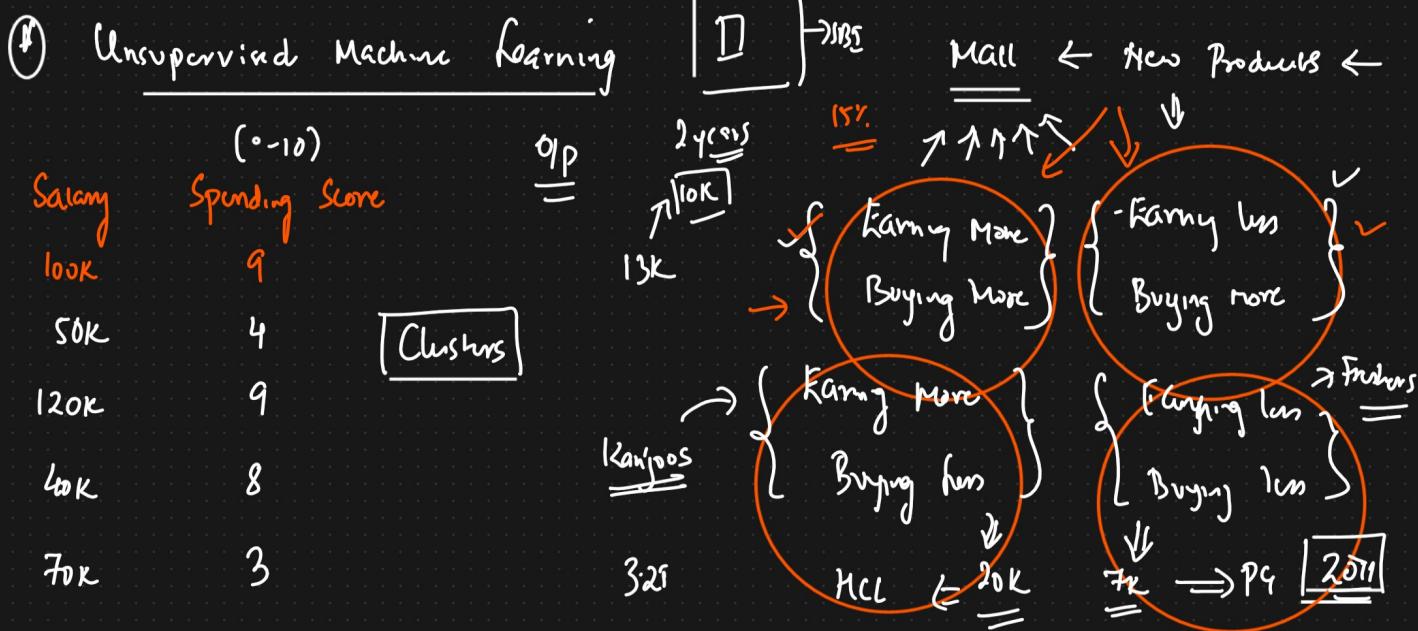
{ ① Projects ② Theoretical ③ Lifecycle of DS Project }

Projects: ① Goal ⇒ Y. 3% Automating something { Business }

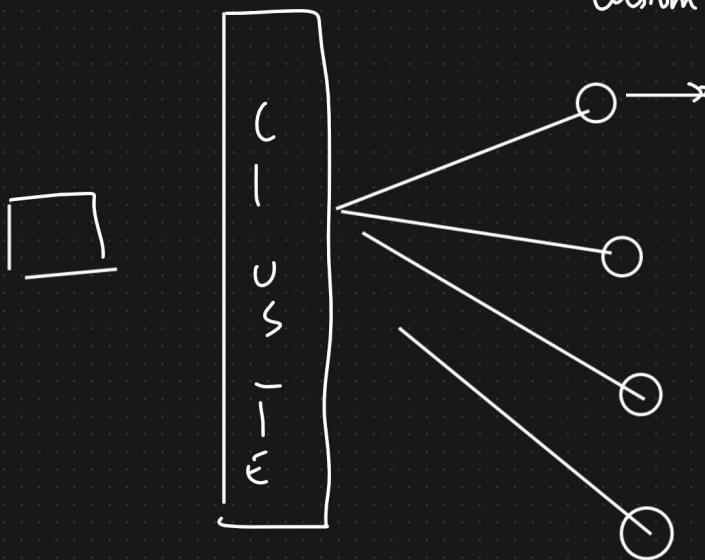


BA { DATA Collection } FE FS ⇒

{ 3rd party API's
↓
Required Dataset of (CATBOOST) } ←



Custom Ensemble Techniques



Salary Spending Slope

100K	9
90K	8
70K	4
120K	8
10K	1
50K	4

Discount to be

O/P

15%

12%

5%

15%

13%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

1%

4%

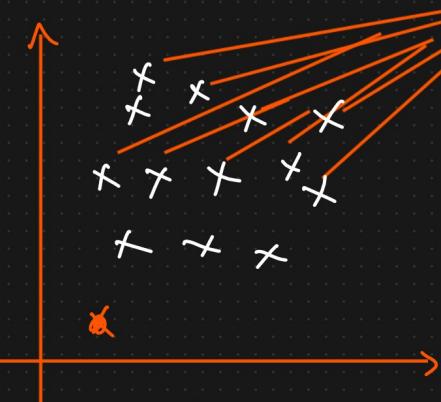
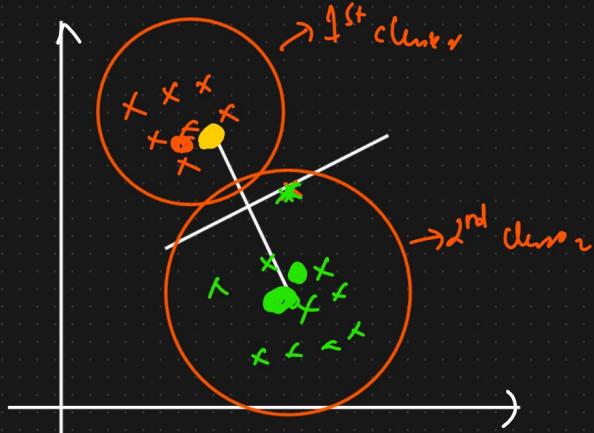
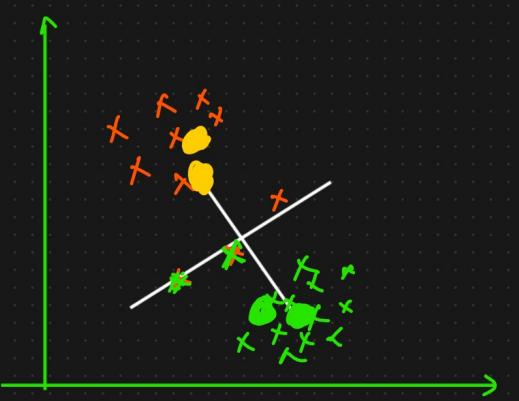
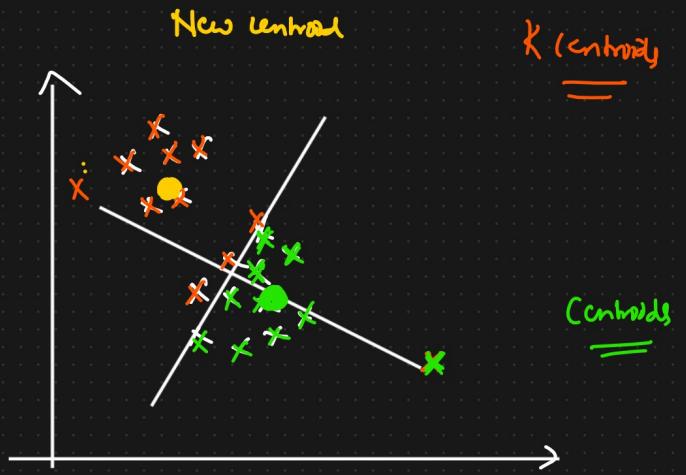
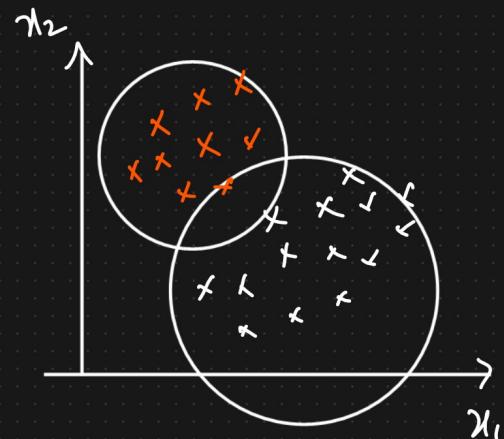
1%

4%

1%

4%

① K Means Clustering



$K=1 \text{ to } 20$
 $\uparrow \uparrow \uparrow$ WCSS of within cluster sum of square }

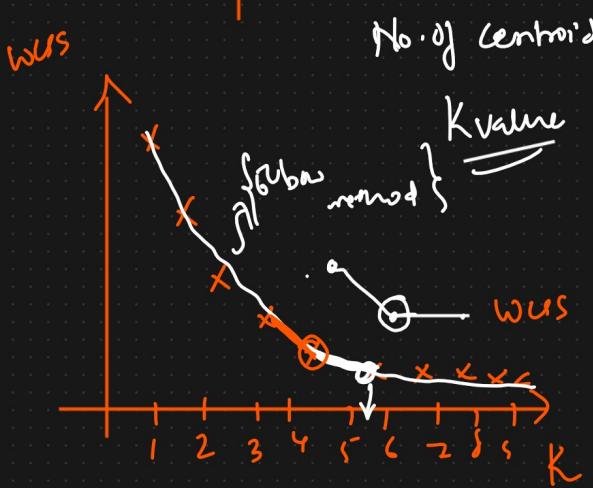
K-2

$\uparrow \uparrow \uparrow$ WCSS of will decrease }

$\downarrow \downarrow \downarrow$ K=3

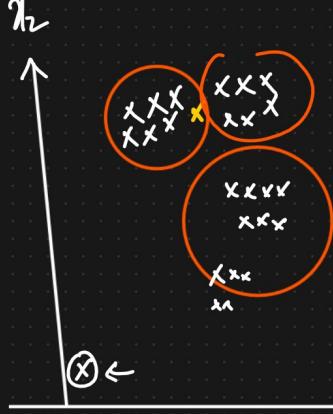
Kmeans ++

$K=3 \downarrow$



$\downarrow \downarrow \downarrow$ K=4

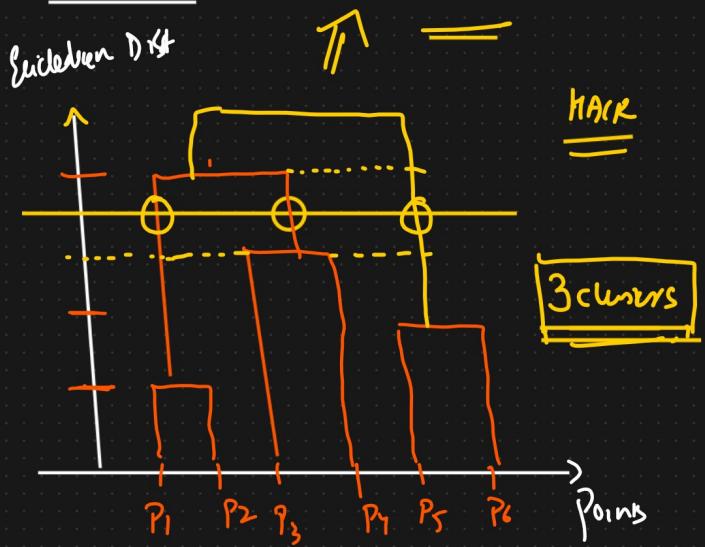
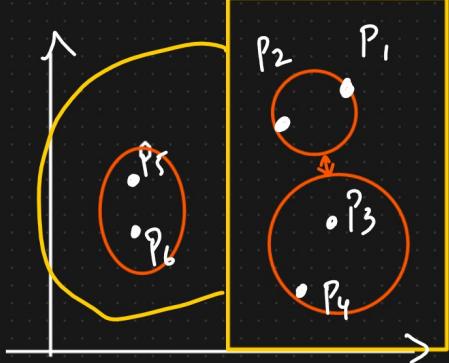
$\downarrow \downarrow \downarrow$



Kmeans ++

Dendrogram

② Hierarchical Mean Clustering



{ ① We need to find the longest vertical line such that none of the horizontal line passes through it. }