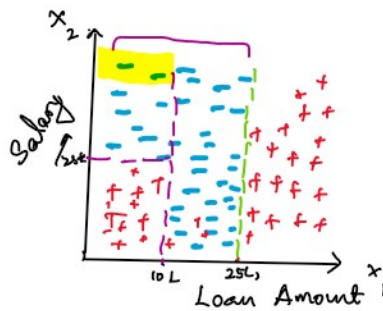## Decision Trees

→ Classification (Binary)

$\boxed{ML}$ → $X \longrightarrow Y$ relation
Cost function
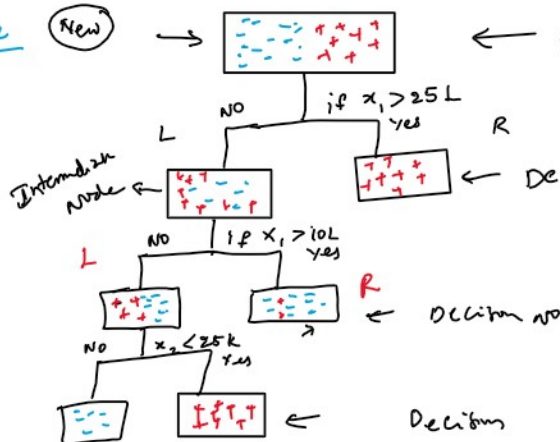Optimization

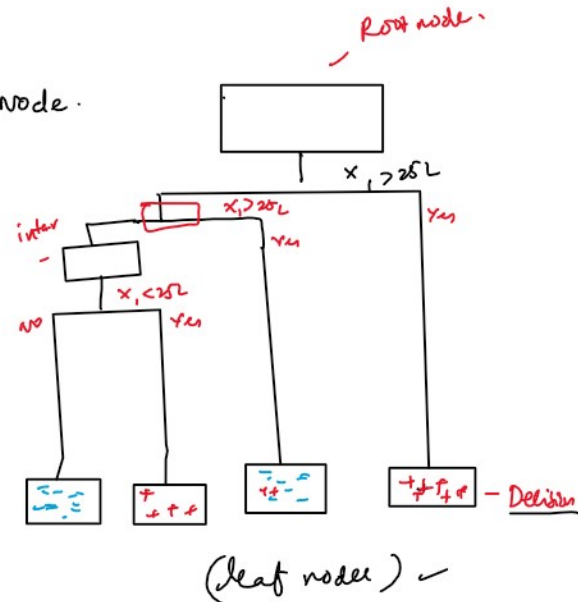Principle :- "Divide and Conquer"



+ - Case (defaulted)
- Paid up ✓

if $X_1 > 25L$ → ⊕ ve ✓

elif if $X_1 > 10L$ → - ve Cases

elm if $X_2 < 25k$ → ⊕ ve

---

$Y$ - Price (100 - 5000)

$\underset{[0-1]}{Nor}$ [X] ← [Y] [100-5000]

$Inn[0-1]$  $\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \cdots w_m x_m$

$w_0$ - high value, $w_1, w_2, \cdots w_m$

. $100, 300 \cdots$

$\underset{[0-1]}{Norm}$ $\overline{[Y]}$  $\hat{y} = w_0 + w_1 x_1 + w_2 x_2 \cdots w_m x_m$

$Y - (M, \sigma)$   $w_0, w_1, w_2 \cdots \{0.1, 0.2\}$

$Z = \dfrac{x_i - M}{\sigma}$    $Z \cdot \sigma + M = x_i$
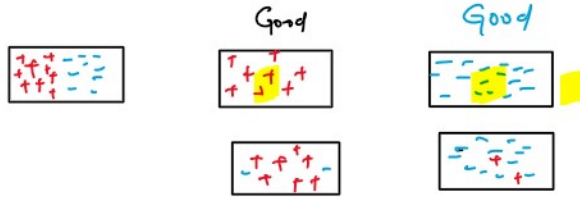
$(\hat{y})$    $Y[0-1] \longrightarrow Y[100-5000]$

$Z = \dfrac{Y_i - M}{\sigma} =$

$\boxed{Y_i = (\sigma \cdot Z + M)}$ ✓

---

Tree (New) → ← Root node.



L  NO  if $X_1 > 25L$  Yes  R

Intermediate Node ←

← Decision

L  NO  if $X_1 > 10L$  Yes  R

← Decision no

NO  $X_2 < 25k$  Yes

← Decision

---

Root node.



$X_1 > 25L$

inter  $X_1 > 25L$  Yes

NO  $X_2 < 25L$  Yes

NO  Yes

+ - Decision

(leaf nodes)

1. Where should we split the data? (for a given variable ($k$))

2. Which variable should split the data first?

## Learning ( Training )

Cost



Good    Good

$$- y_i \log(p_i) - (1 - y_i) \log(1 - p_i)$$

Metric

50% 50%

Entropy :— $- \sum_{i=1}^{K} p_i \log_2(p_i)$

$P$ is the probability of ith class
$K$ is the number of classes

[$k = 2$]

(Impurity)    —vecun

+ve

$\Rightarrow - \frac{5}{10} \log_2\left(\frac{5}{10}\right) - \frac{5}{10} \log_2\left(\frac{5}{10}\right)$

$\Rightarrow - \frac{1}{2} \log_2 2^{-1} - \frac{1}{2} \log_2 2^{-1}$

$\Rightarrow \frac{1}{2} + \frac{1}{2} = 1$

$- \frac{10}{10} \log_2\left(\frac{10}{10}\right) - \frac{0}{10} \ln( \quad )$

$\Rightarrow \underline{0}$

Best    worst

Cost    $\begin{pmatrix} 0 - 1 \end{pmatrix}$

Parent  (Child node)

$E_0 - E_1$

Information : gain

Split

[20]   $E_0 = 1$

$E_0 : 1 \quad E_1 = 0.65$

(if $x_1 < k$)

$n_L$   $n_R$

L [10]    [10] R

$E_L : 0.7$    $E_R : 0.6$

$E_1 = 0.65$

$I.G = 1 - 0.65$

$\boxed{I.G = 0.35}$

$\rightarrow \left( \frac{n_L}{n} \times E_L \right) + \left[ \frac{n_R}{n} \times 0.6 \right]$

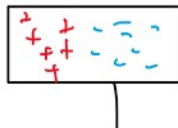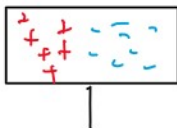$\left[ \frac{10}{20} \times 0.7 \quad + \quad \frac{10}{20} \times 0.6 \right] -$

$x_1 < k$      $x_1 < Q$      $x_1 < R$    Best split ✓
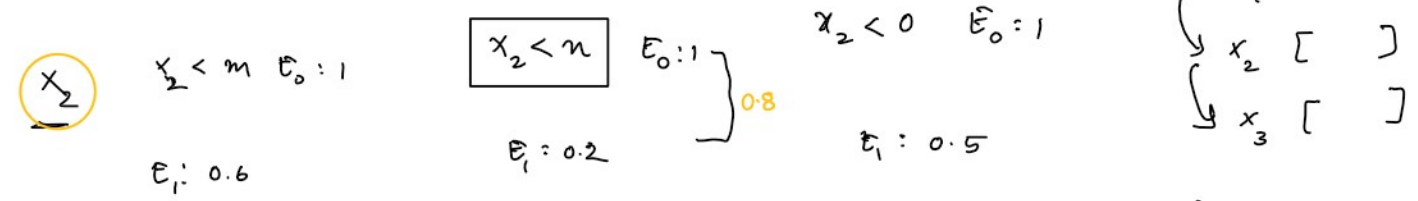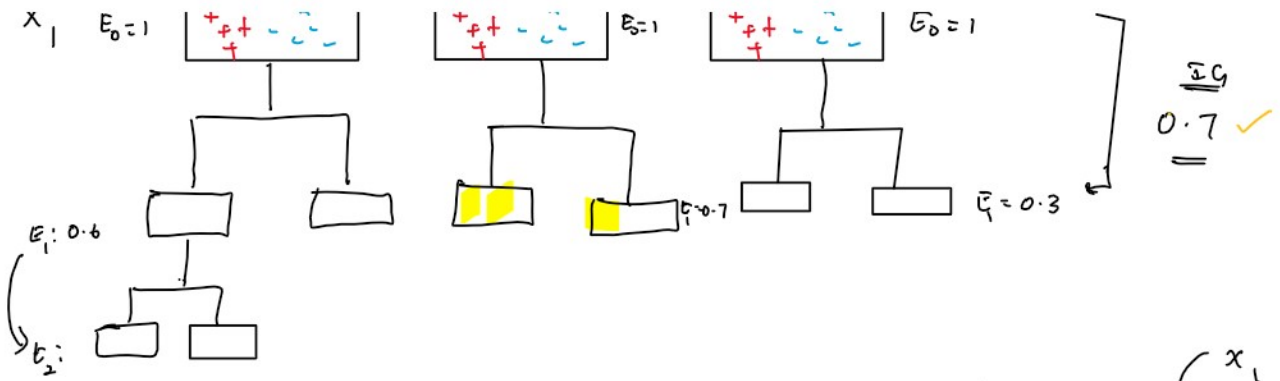
$X_1 \quad E_0 = 1$      $E_0 = 1$      $E_0 = 1$

$\underline{\underline{I.G}}$

$x_1$  $E_0 = 1$     $E_0 = 1$     $E_0 = 1$

$\overline{IG}$

$0.7$ ✓

$E_1 : 0.6$

$E = 0.7$     $\overline{G} = 0.3$

$t_2:$

$x_2$   $x_2 < m$  $E_0 : 1$

$\boxed{x_2 < n}$  $E_0 : 1$  $\Big\}$ $0.8$

$x_2 < 0$   $E_0 : 1$

$x_1$ [      ]
$x_2$ [      ]
$x_3$ [      ]

$E_1 : 0.6$        $E_1 : 0.2$       $t_1 : 0.5$

$x_2 < 20,$  $x_1 < 24$

**Split points**

Sweetness

Y   Switch Calc

$x_1$
1
$x_1 < 2 \rightarrow$ 2
$x_1 < 3 \rightarrow$ 3
4
$\boxed{x_1 < 6}$   5
$0.6$   6
8
9
$x_1 < 10 \rightarrow$ 10

Y
1    $\leftarrow x_1 < 2$
0    $\leftarrow x_1 < 3$
1
1
1
1    $\leftarrow x_1 < 6$
0    $\leftarrow x_1 < 7$
1    $\leftarrow x_1 < 8$
0
0
0

$x_2$
20
5th $\rightarrow$ 36
10th $\rightarrow$ 45
15k $\rightarrow$ 56
$\rightarrow$ 1
$\rightarrow$ 1
$\rightarrow$ 1

$\boxed{x_2 < 30}$  $E_h = 0.8$
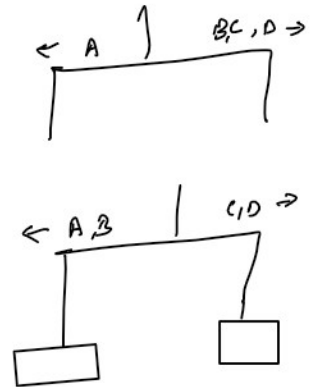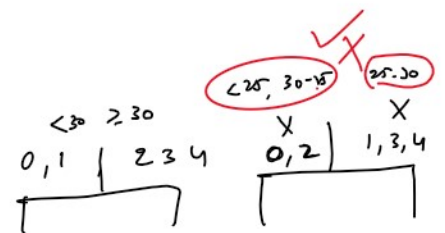
(break) 10:28

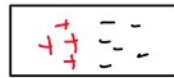$X, \{A, B, C, D\}$

$\rightarrow A \mid B, C, D$

$\rightarrow A, B \mid C D$

$\rightarrow AC \mid B D$

$\rightarrow AD \mid B D$

$\leftarrow A \mid B, C, D \rightarrow$

$\leftarrow A, B \mid C, D \rightarrow$

Ordinal variable $\rightarrow$
$\begin{bmatrix} 0 & 1 & 2 & 3 & 4 \end{bmatrix}$
$<25$  $25-30$  $30-35$  $35-40$  $40+$

$<25$
0  |  1 2 3 4

$<30$  $\geq 30$
0,1  |  2 3 4

$<25, 30-35$  ✓  $25-30$ ✗
$0, 2$  |  $1, 3, 4$

3, 5    E - 1

(6, 4)    ($\tilde{E}$ - 0.8)
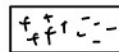
$$\underline{\text{Entropy}} - \sum_{i=1}^{K} P_i \log_2 (P_i)$$

$$= \boxed{- \frac{4}{10} \log_2 \left(\frac{4}{10}\right) - \frac{6}{10} \log_2 \left(\frac{6}{10}\right)}$$

$$= (0.8)$$

$\underline{\text{Gini Index}}$ :   Measure of   Impurity

$$\underline{\text{Gini}} \rightarrow 1 - \sum_{i=1}^{K} (P_i)^2 \qquad \begin{array}{l} P - \text{Prob of a class} \\ i - \text{number of classes} \end{array}$$

$$= 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2\right] \qquad\qquad = 1 - \left[1^2 + 0^2\right]$$

$$= 1 - \left[\frac{1}{4} + \frac{1}{4}\right] \qquad\qquad\qquad = 1 - 1$$

$$= 1 - \frac{1}{2} = \boxed{0.5} \qquad\qquad\qquad = \underline{0} \qquad \text{Best}$$

(1980)

$$\underline{\text{Gini}} - \left[0 - 0.5\right]$$



$\underline{\text{Gini}}$ - Computationally faster

① $E_0$

Root

$x_1 < K$

RoR

$E_1$

$x_2 < P$ Salary

$E_2$

$x_3 < \eta$ dums    (Pure node)

$\rightarrow$ Stopping :-

① when we have a pure node.

② when there is no further I.G

④

$E_2$



$x_4 > m \mid cm$     $x_3 < \eta \, dum$     (Pure node)

$B$

$E_3 : 0.1$

No | Gender = Male
Yes

$E_4^*$  C-1     --
6      2  ✗

① → Highly prone to overfitting ✓

→ Greedy

## Early Stopping Criteria
## (Pruning) →

✓✓ ① Depth of the tree ✓  ✓

Max-depth = 3

→✓ ② Min no. of observation to split

⑩  ✗

→✓ ③ Min no. of observations in
leaf node = ⑤

Entropy

(0-1)

++  --
++  --

$P_+ = P$
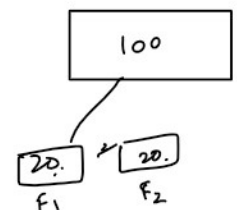$P_+ = 1-P$
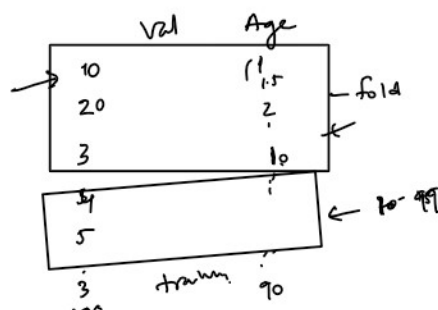
$P_+$     $P_-$

$-y_i \log(P_i) - (1-y_i)\log(1-P_i)$

$-P\log_2(P) - (1-P)\log_2(1-P)$  ←
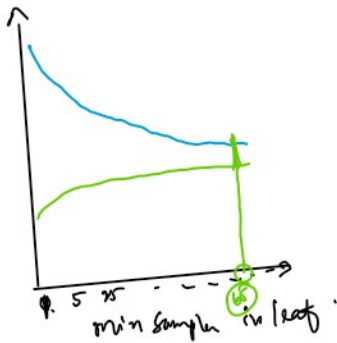
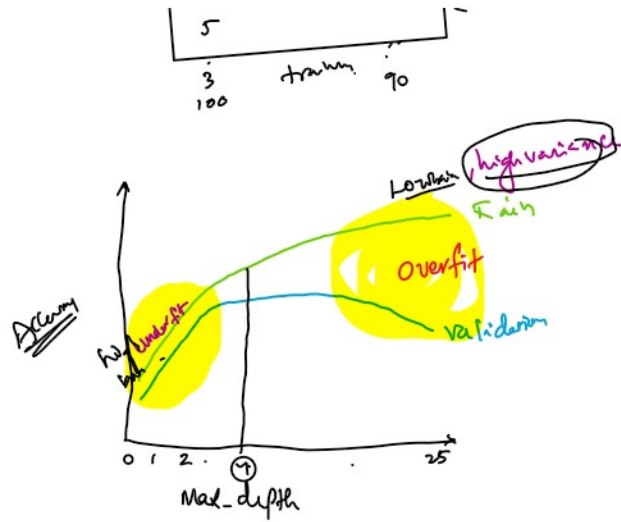$$\text{Entropy} = -\sum_{i=1}^{K} P_i \log_2(P_i)$$

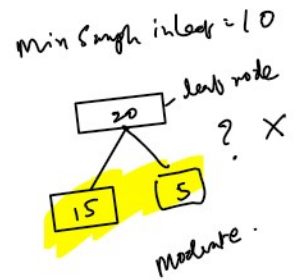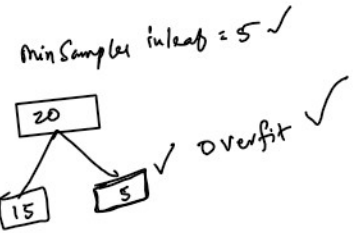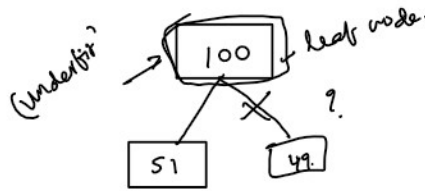$$\text{Gini} = 1 - \sum_{i=1}^{K} (P_i)^2$$

==Start : 9:05==

① folds = KFold ( n: split = 5 , Shuffle : True , random-state = ① )

val     Age

10      11.5
20      2       — fold
3       10

4
5       ← 10-99

3     trainin   90
100

100

20.     20.
F1      F2

5
3
100  traiin  90

Plot



Low bias, high variance
Train
Overfit
validation
High bias
underfit
Accuracy
0 1 2 . (4) . . 25
Max_depth

min Samples in leaf = 5 ✓



20
15    5
✓  Overfit ✓

min Samples in leaf = 10

leaf node
20
15    5   ? ✗
Moderate.



0. 5 25
min samples in leaf

min sample in leaf = 50

(underfit?)

100   + leaf node.
51      49.   ?

→ Minimum Samples to Split

min samples = 50
100
deeper

min = 100
120

min split: 200
120   ?
✗ ↑
Stopping underfitting

high (?) — underfit

low (?) — overfit

max features = 4

Split points < 'best' approach
'Random' ←
age {10 — 99}

{x₁, x₂ ... x₁₀}

{x₁, x₃, x₈, x₉}   age < 21
age < 36
age < 61

age < 20
age ≤ 25

$\{x_1, x_9, x_7, x_6\}$

$age < 61$

max. features = 'auto'

Classification

Regression

$\sqrt{m}$ — m no. of features

$\frac{m}{3}$ — m no. of features

Classification

Regression ?

Entropy
Gini
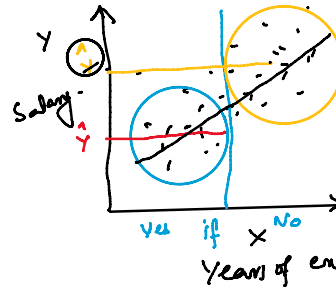$\Bigg\}$ for regression?

$-\sum p_i \log_2(p_i)$   $\times$

if $x < 5$ $\hat{y} = mean(y_L)$

$x > 5$ $\hat{y} = mean(y_R)$

MSE

if years < 5 avg salary $15L$

$\geq 5$ avg Sal $30L$

Y
Salary

Yes   if  x  No

if $x < 5$

Years of emp.

Rootnode
$5L - 95L$   $\mu = 37$

m=100

$Y$   $var(Y)$   $20,000$

year of emp < 5

yes
$60$   $\mu = 15L$

No.
$40$   $\mu = 65L$
$65 - 95L$

$5600$

$5 - 43L$

$50L$

$6000$

$100b$
$75 - 95L$

$84L$

Multiclass Classification  -  Entropy,  Gini work?

$$-\sum_{i=1}^{K} P_i \log_2(P_i)$$

K-classes

$$\boxed{10 \quad 20 \quad 30}$$

$$-\frac{10}{60}\log_2\left(\frac{10}{60}\right) - \frac{20}{60}\log_2\left(\frac{20}{60}\right) - \frac{30}{60}\log_2\left(\frac{30}{60}\right)$$

(Binary)
Classification and Regression Trees (CART)     $\rightarrow$     Ensembles