9:05 pao
milte hai

DT

Entropy
⇓
Impurity

Age > 30
30 c & 50 c'

Gender =

$$H(y) = -p \log_2 p - (1-p) \log_2 (1-p)$$

entropy (y):



y. value_counts()

$\downarrow$
#C        #NC

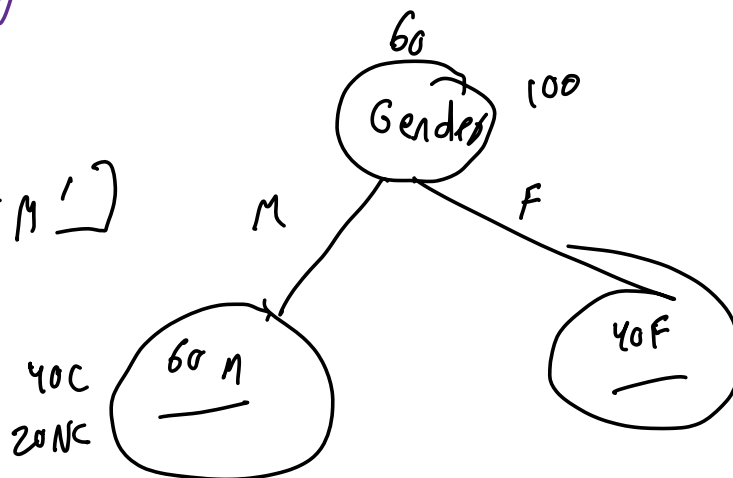$$-\frac{3}{5}\log\left(\frac{3}{5}\right) - \frac{2}{5}\log\left(\frac{2}{5}\right)$$
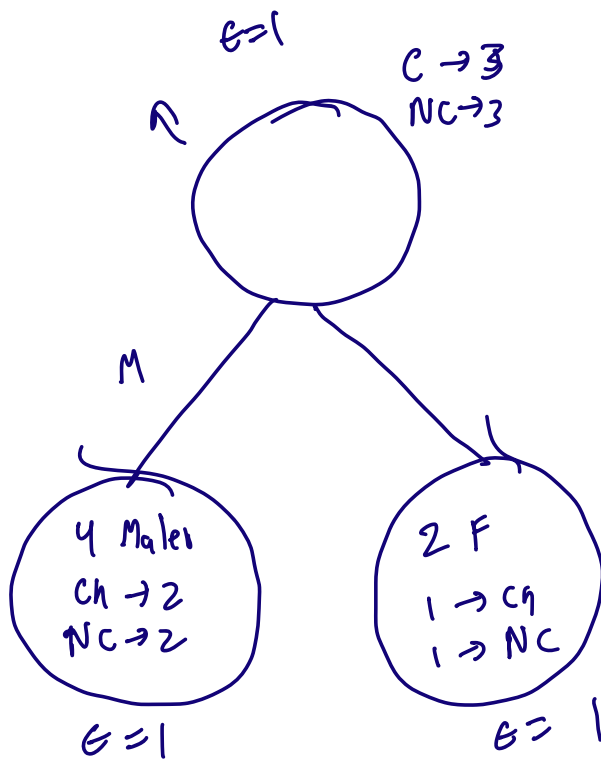
y. shape[0]

| Age | exp | Choose |
|-----|-----|--------|
| 20  | 10  | Y      |
| 30  | 20  | N      |
| 40  | 30  | Y      |
| 60  | 80  | N      |
| 40  | 20  | Y      |

entropy (y):

y [feature == 'M']



60
Gender    100

M          F

40C
20NC        60 M          40 F

E=1

C → 3
NC → 3

| Gender | Age | Attrition |
|--------|-----|-----------|
| M | 20 | Yes |
| F | 30 | No |
| M | 40 | Yes |
| F | 50 | Yes |
| M | 60 | No |
| M | 80 | No |

M

4 Males

Ch → 2
NC → 2

G = 1

2 F

1 → Ch
1 → NC

G = 1

Gini Index

$$GI(y) = 1 - \sum_{i=1}^{k} (p(y_i))^2$$

$-p \log p$
⇓
$p^2$

2 class:

$$1 - [p(y_i = 1)^2 + p(y_i = 0)^2]$$

Case 1:

$$p(+) = 0.5$$
$$p(-) = 0.5 \implies \text{Entropy} = 1$$

$$= 1 - [(0.5)^2 + (0.5)^2]$$

$$= 1 - [0.25 + 0.25]$$

$$= 0.5$$

**Case 2:**

$$P(+) = 1$$
$$P(-) = 0 \implies E = 0$$

$$= 1 - [(1)^2 + (0)^2]$$

$$= 0$$

$$[0 - 1] \rightarrow \text{Entropy}$$
$$[0 - 0.5] \rightarrow \text{Gini}$$

## Split Numerical Column

Age < 35 20

3 0

4 0

| $f_1$ | $f_2$ | $y_i$ |
|---|---|---|
| $c_1$ | 2.2 | 1 |
| $c_2$ | 2.6 | 1 |
| $c_1$ | 3.5 | 0 |
| $c_4$ | 3.8 | 0 |
| $c_5$ | 4.6 | 1 |
| $c_1$ | 5.3 | 0 |

① $f_2 \leq 2.2$

$IG_1$

② $f_2 \leq 2.6$

$IG_2$

③ $f_2 \leq 3.5$

$IG_3$ ✓

$x \leq Q_1$     $x \leq Q_2$     $x \leq Q_3$     $x \leq Q_4$     ✓

$IG_1$

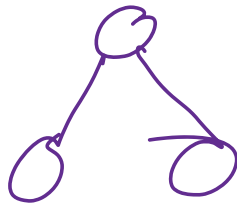XG Boost

| Age | Aff |
| --- | --- |
| 20 | Y |
| 18 | N |
| 20 | X |
| 19 | X |
| 21 | N |



Age ≤ 20

Y            N

Ruge depth

⇓

overfit

Decision stump →

Shallow tree

Deep tree

depth ⇒ 2 → cv_error
     ⇒ 3 →  ⟶
     ⇒ 4 →  ⟶   }

Train    CV    Test

O-35

Age ≤31

≤M

Age ≤18

2 minutes
reg ability