# Decision Tree

$\rightarrow$ D.T Intution

$\rightarrow$ Wonking D.T

$\rightarrow$ Entropy / Info . gain

$\rightarrow$ Code

→ Data Scientist Aintel

↳ Employee athition

→ 1/0

1. Find the chances of athition of a Employee

2. What are the key factor responsible for Athition

↳ Feature Importance

Age Overtime | y

Age Overtime | y

27        2

Overtime

- (+) Leave
- (-) Stay

→ Axis parallel
   hyperplane

→ Non-linear
   model



Overtime

Hours

3.5    +  +  +  +  -  -  +  +  +  +
3      +  +  +  +  -  -  +  +  +  +
       +  +  +  +  -  -  +  +  +  +
2.5    +  +  +  +  -  -  +  +  +
2      -  -  -  -  -  -  +  +  +  +
       -  -  -  -  -  -  +  +  +  +
1.5    -  -  -  -  -  -  +  +  +  +
1      -  -  -  -  -  -  +  +  +  +
       -  -  -  -  -  -  +  +  +  +
0.5    -  -  -  -  -  -  +  +  +  +
       -  -  -  -  -  -  +  +  +  +
                              → Age
18              29      35

Years

Training

If (Age < 29):
  If (Overtime < 2.5):
    Stay
  else:
    Leave
else
  If (Age < 35):
    Stay
  else:
    Leave

Overtime

· (+) Leave
· (-) Stay

Hours

3.5
3
2.5
2
1.5
1
0.5

18    29    35    Age

Years

If (Age < 29):

If (Overtime < 2.5):

Stay

else:

Leave

else

If (Age < 35):

Stay

else:

Leave

Age < 29

Yes                    No

Overtime < 2.5              Age < 35

Yes          No        Yes        No

Stay      Leave      Stay      Leave

$\rightarrow$ Which options to choose

$y+ = 100$

$\underbrace{f_1 > c}$   $y- = 100$

$n = 200$

$n_1 = 100$   $n_2 = 100$

$\begin{matrix} y+ = 10 \\ y- = 90 \end{matrix}$   $\begin{matrix} y+ = 90 \\ y- = 10 \end{matrix}$

$-ve$

①

$n = 200$

$\underbrace{f_4 < c_2}$   $y+ = 100$

$y- = 100$

$n_1 = 100$   $n_2 = 100$

$\begin{matrix} y+ = 60 \\ y- = 40 \end{matrix}$   $\begin{matrix} y+ = 40 \\ y- = 60 \end{matrix}$

②

$\rightarrow$ More homogenous = More confidence in prediction

$\longrightarrow$ Entropy

        $\hookrightarrow$ Measures Randomness

        $\hookrightarrow$ Impurity (Heterogenity)

$$H(y) = -\sum_{i=1}^{K} P_i \log P_i$$

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

No. of classes →

K

Prohab. of each class →

Ranges $(0-1)$ →

$$H(y) = -\left[ P_{blue} \log_2 P_{blue} + P_{red} \log_2 P_{red} \right]$$

$P(blue) = \frac{3}{6} = \frac{1}{2}$

$P(red) = \frac{3}{6} = \frac{1}{2}$

$$= -\left[ \frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right]$$
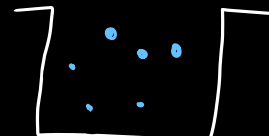
$= 1$

$P_{blue} = \frac{5}{6}$

$P_{red} = \frac{1}{6}$

$$H(y) = -\left[ P_b \log_2 P_b + P_n \log_2 P_n \right]$$

$$= -\left[ \frac{5}{6} \log_2 \left(\frac{5}{6}\right) + \frac{1}{6} \log_2 \left(\frac{1}{6}\right) \right]$$
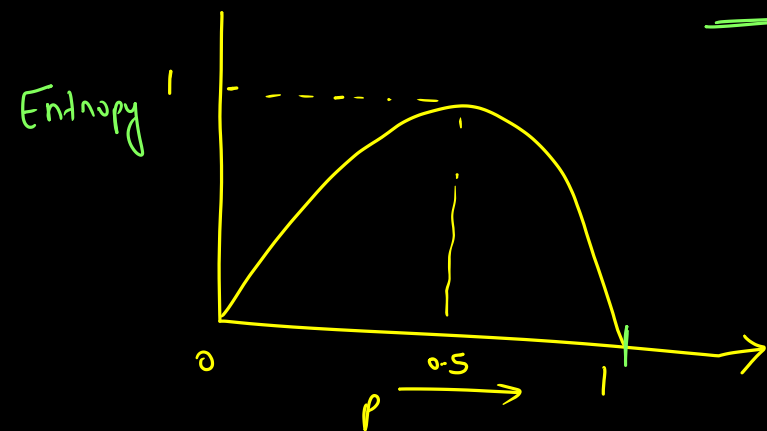
$P_b = 1$

$P_n = 0$

$$H(y) = -\left[ 1 \log_2 (1) + 0 \log_2 (0) \right]$$

$= 0$

0     0

Entropy

0     0.5     1

P

→ Break until 22:19
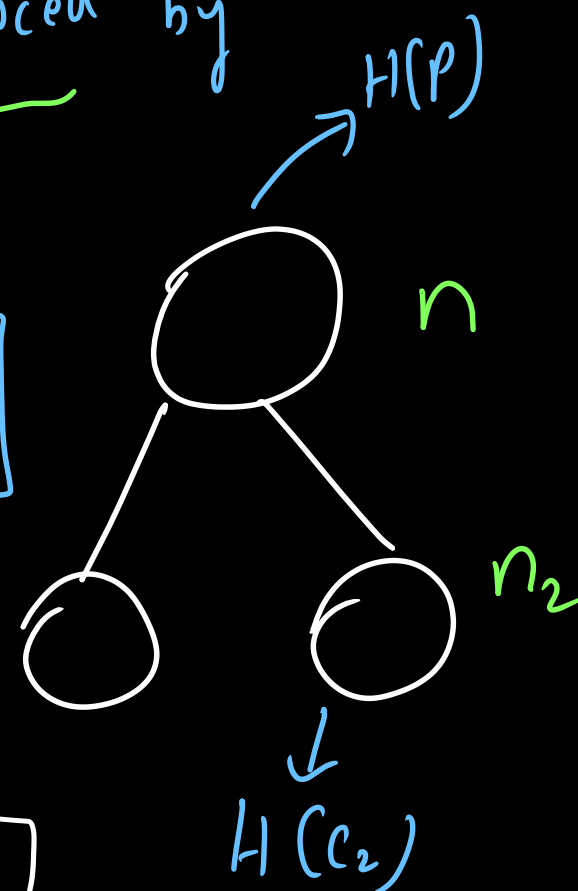
.

# Information Gain

(Age, Gender)

→ Used to decide which feature to split on

{ → Measure how much entropy is reduced by making a split

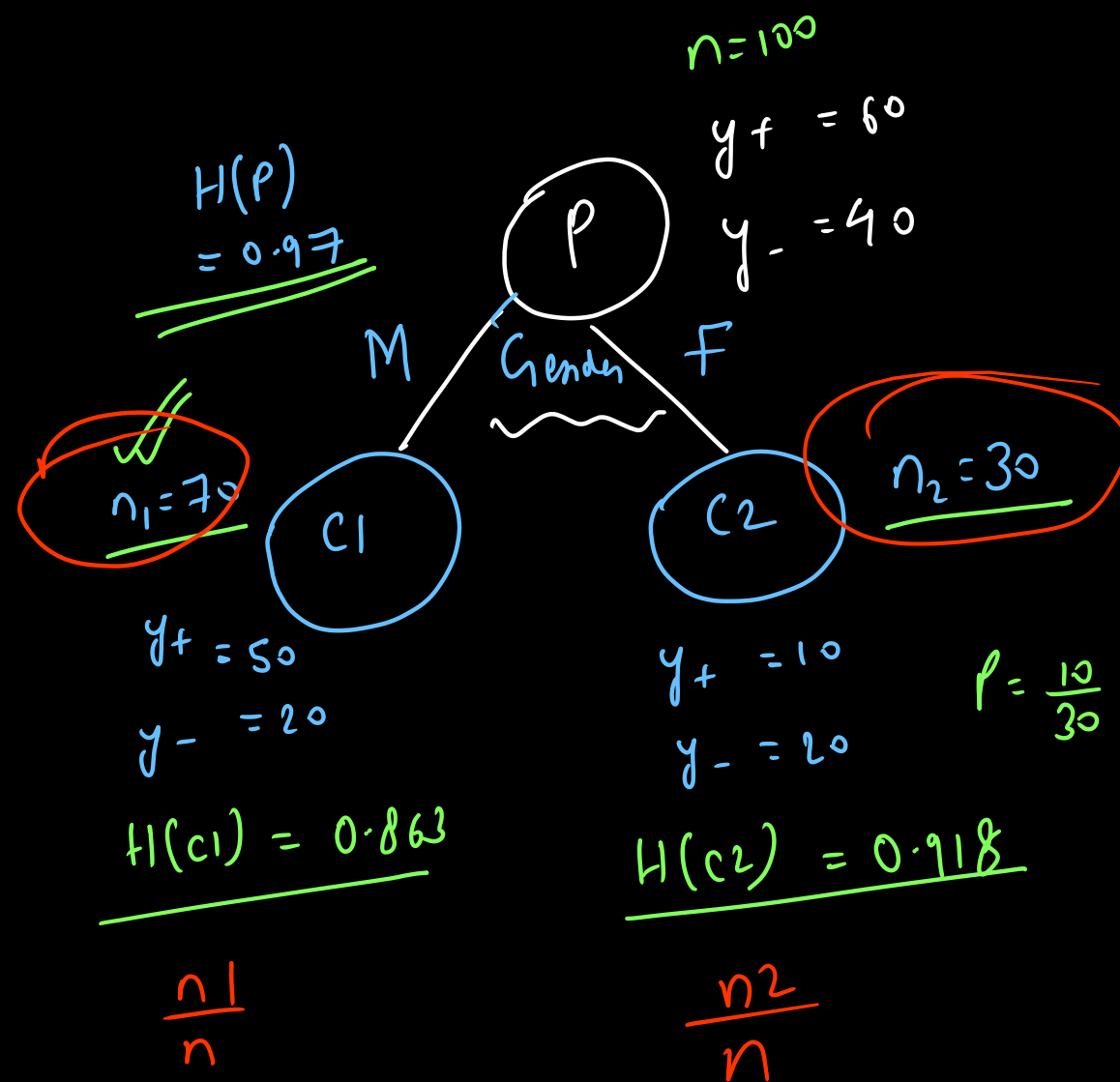→ $IG = H(P) - \left[ \frac{n_1}{n} H(C_1) + \frac{n_2}{n} H(C_2) \right]$

$H(P)$

$n$

$n_1$

$H(C_1)$

$n_2$

$H(C_2)$

→ High $IG$ → high reduction in entropy

(Heterogenous) → More Homogenous

✓ $n = 100$

$y+ = 60$
$y- = 40$

| Age | Gender | $y$ |
|-----|--------|-----|
| 20 | M | $y+$ |
| → 30 | F | $\vdots$ |
| 5 | F | $y-$ |

$H(P)$
$= 0.97$

$n = 100$
$yf = 60$
$y- = 40$

P

M   Gender   F

$n_1 = 70$    C1

$n_2 = 30$    C2

$y+ = 50$
$y- = 20$

$y+ = 10$
$y- = 20$

$P = \frac{10}{30}$

$H(c_1) = 0.863$

$H(c_2) = 0.918$

$\frac{n_1}{n}$
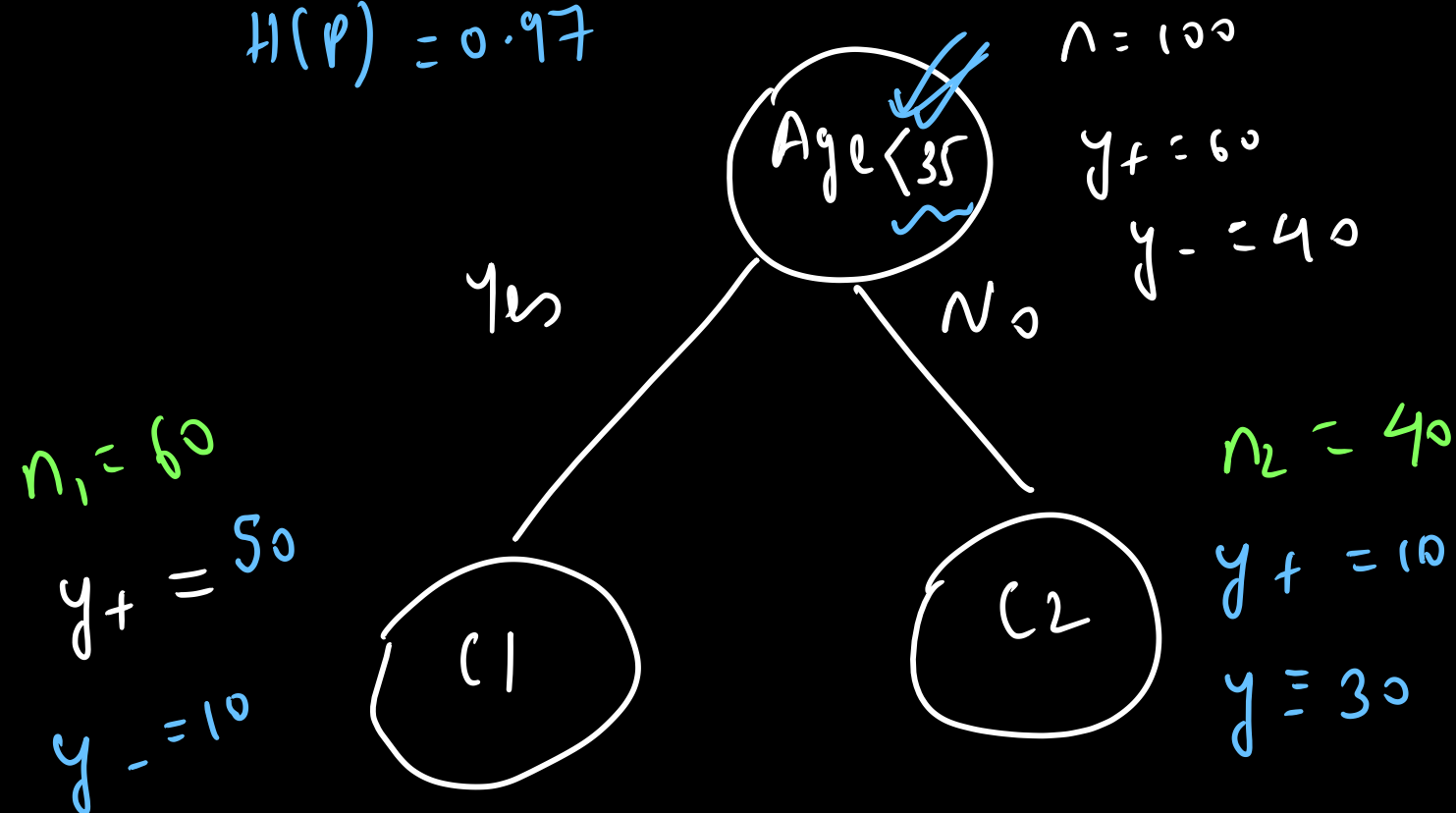
$\frac{n_2}{n}$

$IG_{(Gender)} = H(P) - \left[ \frac{n_1}{n} H(c_1) + \frac{n_2}{n} H(c_2) \right]$

$= 0.97 - \left[ 0.7 \times 0.863 + 0.3 \times 0.918 \right]$

$= 0.0914$

$H(P) = 0.97$

Age < 35

$n = 100$
$y+ = 60$
$y- = 40$

Yes

No

$n_1 = 60$
$y+ = 50$
$y- = 10$

C1

C2

$n_2 = 40$
$y+ = 10$
$y = 30$

$H(C_2) = 0.8112$

$H(C_1) = 0.65$

$$IG(Age<35) = H(P) - \left[ \frac{n_1}{n} H(c_1) + \frac{n_2}{n} H(c_2) \right]$$

$$= 0.97 - \left[ \frac{60}{100} \times 0.65 + \frac{40}{100} \times 0.8112 \right]$$

$$= 0.2565$$

Age        Gender        y

                M           y+
10

20              F

                            y-
15

7

35

40

7  → Age <=7

10       Age <=10

15         |

20         |

35       Age <=35

40

0 - 10

10 - 20

20 - 30

30 - 40