# Attrition Rate of Airtel

Attrition → Yes ⟹ left

→ No ⟹ stayed

The company now has two options - ① Retain   ② Recruite

What analysis can we do on this data?

① Probability of leaving of an employee - classification - Log.Reg.

② Factors' contribution towards attrition - Model's
                                                    Interpretation

$$\vec{w} \cdot \vec{x} + w_0 = 0 \Rightarrow w_1 x_1 + w_2 x_2 + \ldots + w_n x_n + w_0 = 0$$

↳ salary

if $w_i > w_j$ mean feature `i' has higher impact on attrition than feature `j'

# Home work: Standard procedure of creating an ML model

① Acquire the data - Data Ingestion Pipeline

② Pre processing - ⓐ clean
                        ↳ Duplicates
                            ↳ Missing values
                  ⓑ Encoding
                  ⓒ Feature scaling
                      ↳ standardize | Normalize
                  ⓓ Rebalancing data
                  ⓔ Treatment of outliers
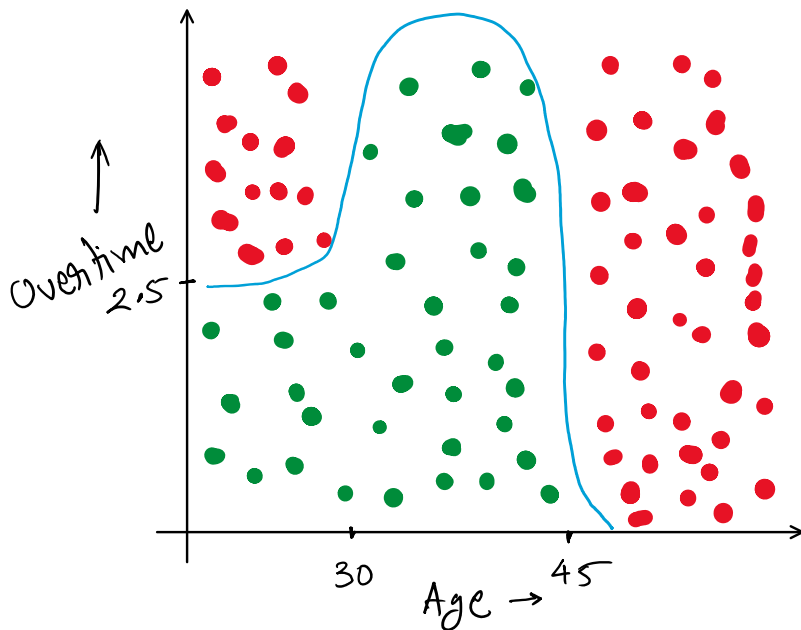
ⓔ Treatment of outliers

ⓕ Feature Engineering
　　└ Reduce the dimensionality
　　└ add new features (more relevant)

ⓖ EDA

③ ML model

△ How about solving this problem with a
different approach.



→ This data is not linearly seperable. ∴ LR won't work

→ Can we use polynomial Logistic Reg.?
　　└ Complex
　　∴ More chances to overfit

→ KNN?
　　└ slow for big datasets

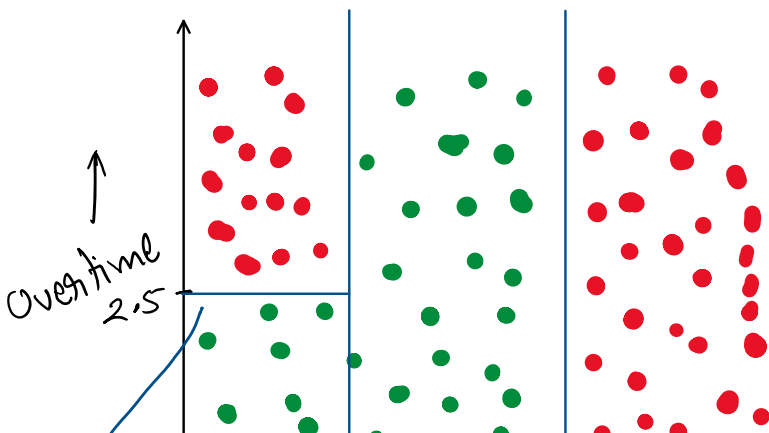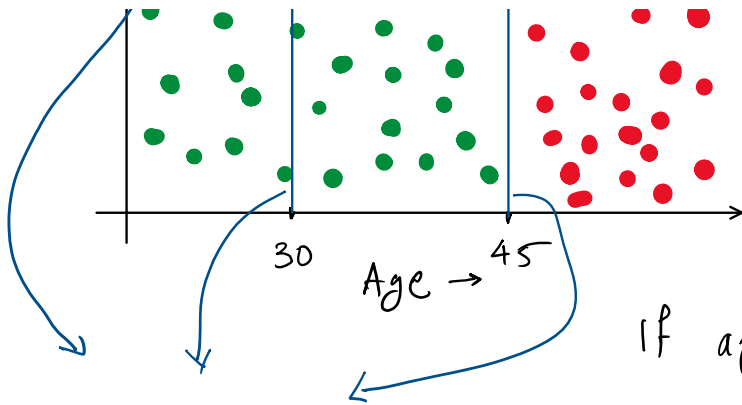→ But Airtel has ~87,000 employees for which KNN might become slow.



Let's ask a few questions to our data:

age ⩾ 45 ⇒ Leave

30 < age < 45 ⇒ Stay

age < 30 ⇒ OT > 2.5 ⇒ Leave

$age < 30 \Rightarrow OT > 2.5 \Rightarrow$ Leave

How can we implement this logic?

Axis Parallel Decision Boundaries

If age $\geqslant 45$: leave

else:
   if $30 < age < 45$: stay
   else:
      if $OT > 2.5$: Leave
      else: stay

If $OT > 2.5$:
   if Age $< 30$: Leave
   else:
      if Age $\geqslant 45$: Leave
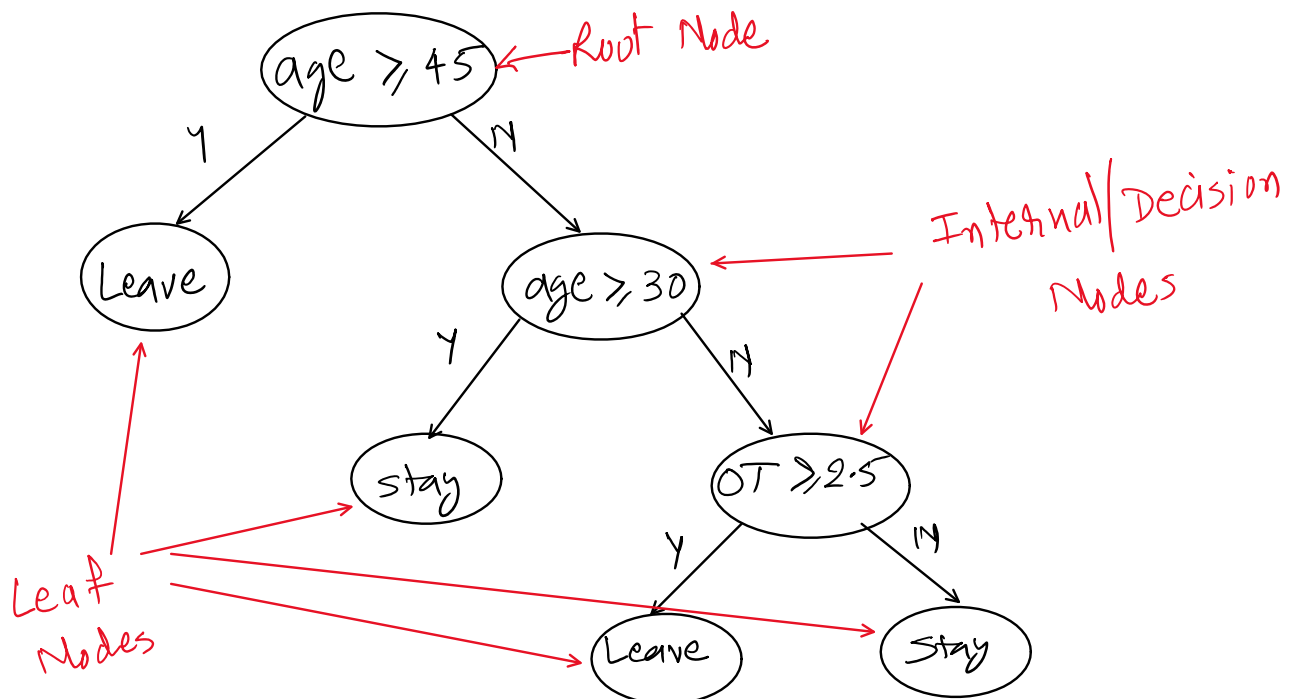      else: stay

else: if Age $\geqslant 45$: Leave
   else: if Age $< 30$: Stay
      else: stay
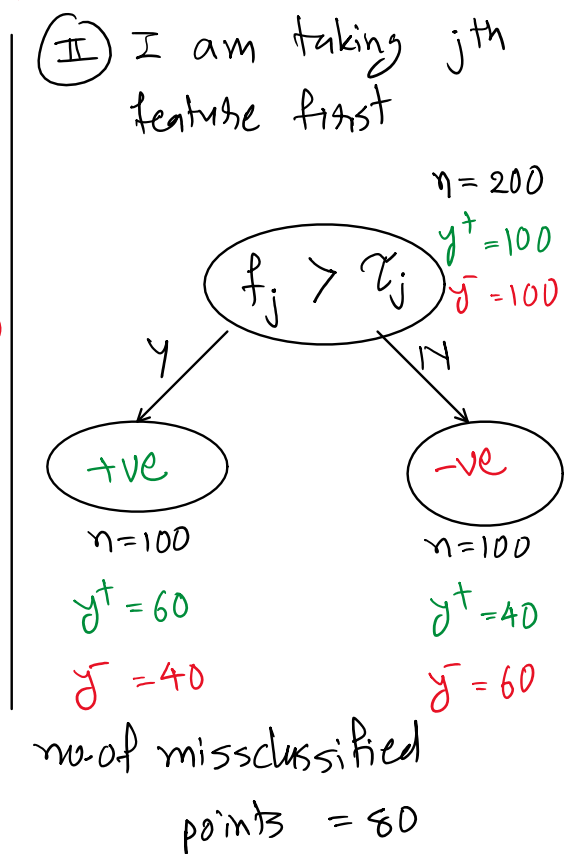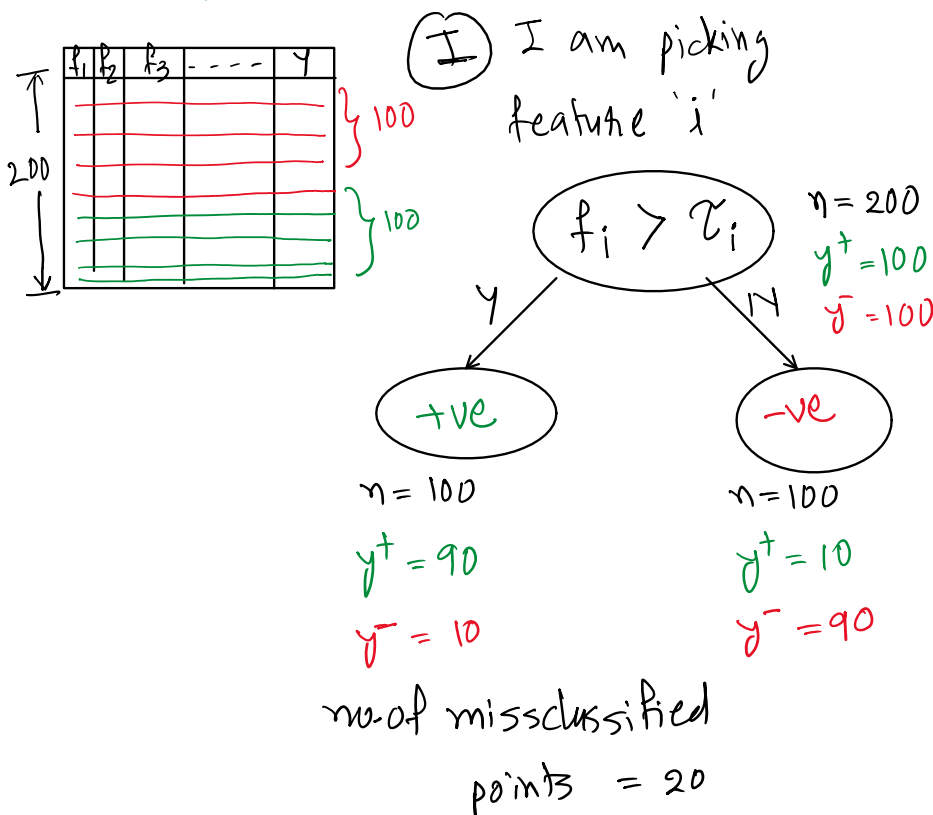
Can we write this logic in a different way as below?

✡ Which questions should we start with/ we ask to our data?

→ This is a v. imp point because the questions we decide will be asked to each & every datapoint and hence can become computationally very expensive if chosen in wrong sequence.
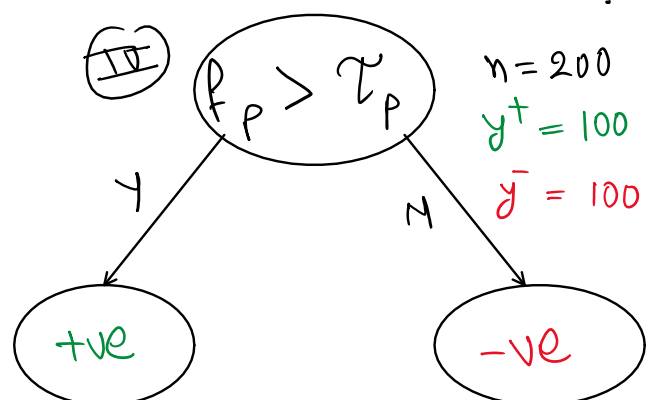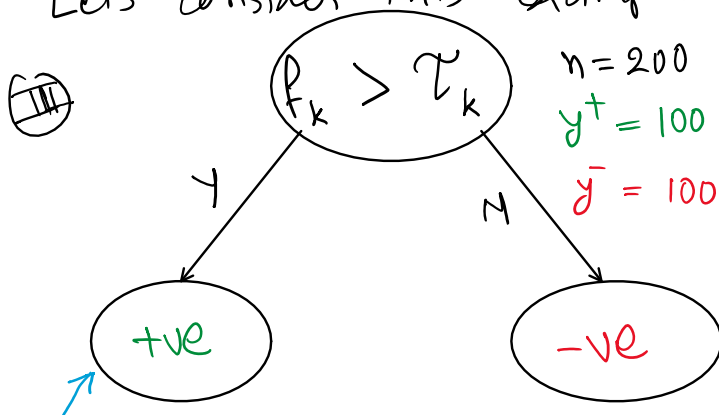
Let's understand this concept by an example.

→ Suppose we have 200 data points with two classes

$y^+$ : 100 datapoints          $y^-$ : 100 data points.



I) I am picking feature 'i'

$f_i > \tau_i$

$n = 200$
$y^+ = 100$
$y^- = 100$

Y → +ve          N → -ve

$n = 100$          $n = 100$
$y^+ = 90$          $y^+ = 10$
$y^- = 10$          $y^- = 90$

no. of missclassified points = 20

II) I am taking jth feature first

$f_j > \tau_j$

$n = 200$
$y^+ = 100$
$y^- = 100$

Y → +ve          N → -ve

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

no. of missclassified points = 80

Lets consider this example:

III)

$f_k > \tau_k$

$n = 200$
$y^+ = 100$
$y^- = 100$

Y → +ve          N → -ve

IV)

$f_p > \tau_p$

$n = 200$
$y^+ = 100$
$y^- = 100$

Y → +ve          N → -ve

| +ve | -ve | +ve | -ve |
|-----|-----|-----|-----|
| $n = 80$ | $n = 120$ | $n = 90$ | $n = 110$ |
| $y^+ = 80$ | $y^+ = 20$ | $y^+ = 85$ | $y^+ = 15$ |
| $y^- = 0$ | $y^- = 100$ | $y^- = 5$ | $y^- = 95$ |

no. of misclassified points = 20

Pure Homogenous Region

Pure Homogenous Node

Slightly less homogenous

✳ We need to quantify "homogenity" of a region/node

Ans- Entropy - measure of impurity/Heteroginity

Entropy high = high Heteroginity = Low Homogenity = Less Purity

Entropy of a node 'y' is denoted by $H(Y)$ & given by:

$$H(Y) = -\sum_{i=1}^{m} P(y_i) \cdot \log_2 P(y_i)$$
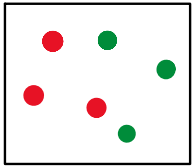
In our example, our labels are : {+ve, -ve}

$$\therefore P(-ve) = 1 - P(+ve)$$

$$H(Y) = -\left[P(y^+) \cdot \log_2 P(y^+) + P(y^-) \cdot \log_2 P(y^-)\right]$$

$$H(Y) = -\left[P(y^+) \cdot \log_2 P(y^+) + (1 - P(y^+)) \cdot \log_2 (1 - P(y^+))\right]$$
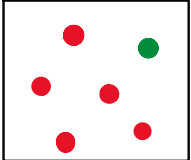
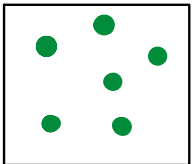$\rightarrow$ Recall Logloss ?

# Example-

$P(-ve) = \frac{1}{2}$     $H(y) = 1$
$P(+ve) = \frac{1}{2}$

$P(-ve) = \frac{5}{6}$     $H(y) = 0.65$
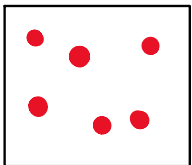$P(+ve) = \frac{1}{6}$

$P(-ve) = 0$     $H(y) = 0 \rightarrow$ Purest Node
$P(+ve) = 1$
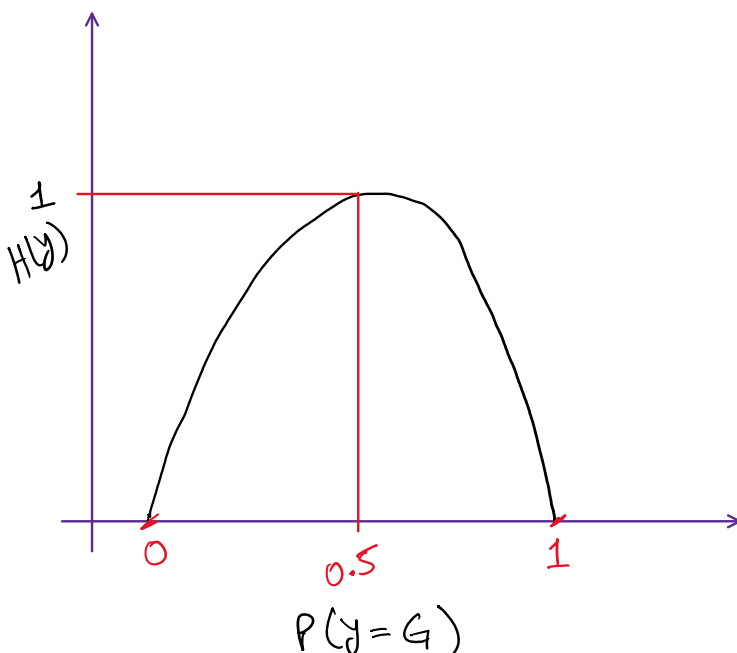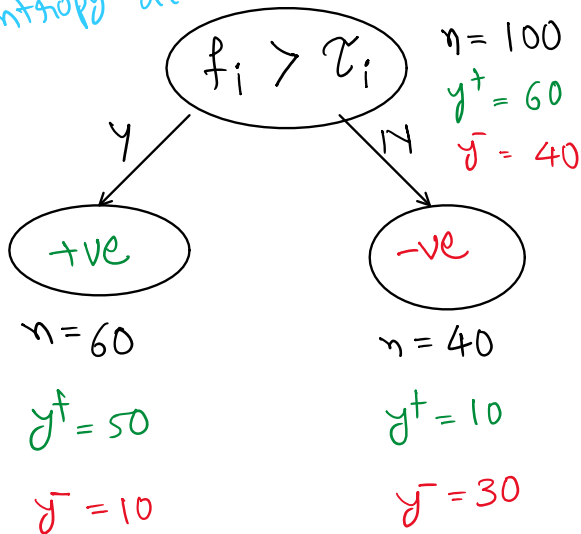
$P(-ve) = 1$     $H(y) = 0 \rightarrow$ Purest Node
$P(+ve) = 0$



$P(y=G)$

$\bigstar$ Coming back to our question - Which feature should I consider first ? / which question should I ask first

Entropy at Parent level = $H_p$



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

Y → +ve
N → -ve

+ve: $n = 60$, $y^+ = 50$, $y^- = 10$

-ve: $n = 40$, $y^+ = 10$, $y^- = 30$

Entropy at children level = $H_c$

Drop in Entropy = $H_p - H_c$



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

Y → +ve
N → -ve

+ve: $n = 70$, $y^+ = 50$, $y^- = 20$

-ve: $n = 30$, $y^+ = 10$, $y^- = 20$

① Entropy at the parent level:

$P(G) = 6/10$     $P(R) = 4/10$
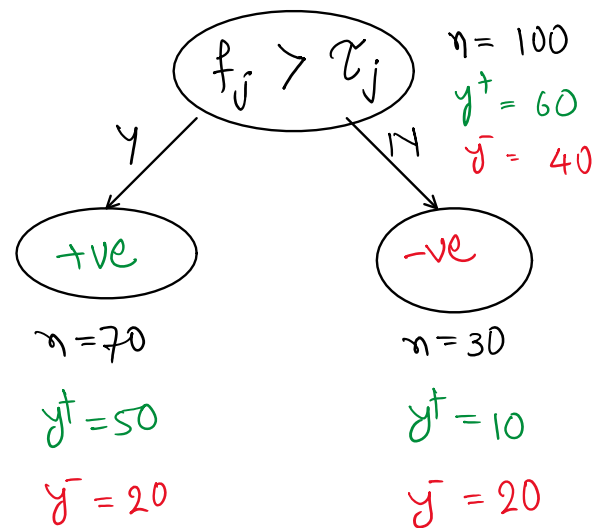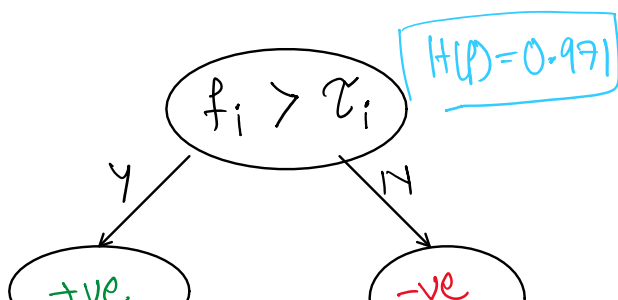
$H(P) = 0.971$

Entropy at children level:

### Left Question (feature)

Left child : $H(y) \Rightarrow P_{-g} = 5/6$ & $P_{-r} = 1/6$

$\therefore H(y) = 0.65$

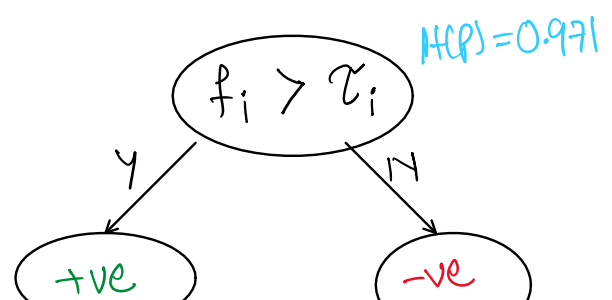Right child : $H(y) \Rightarrow P_{-g} = 1/4$ & $P_{-r} = 3/4$
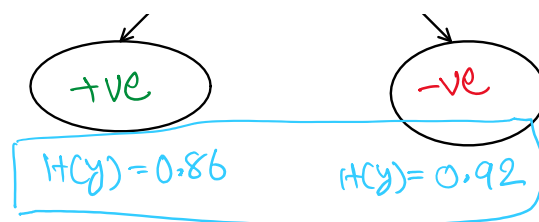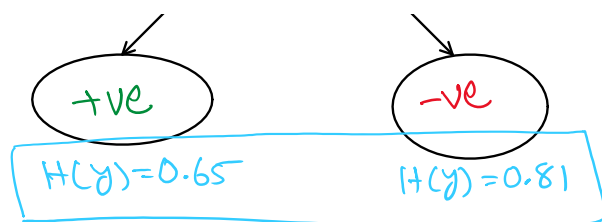
$H(y) = 0.81$

### Right Question (feature)

Left child : $H(y) \Rightarrow P_{-g} = 5/7$ & $P_{-r} = 2/7$

$\therefore H(y) = 0.86$

Right child : $H(y) \Rightarrow P_{-g} = 1/3$ & $P_{-r} = 2/3$

$\therefore H(y) = 0.92$



$H(P) = 0.971$

$f_i > \tau_i$

Y → +ve
N → -ve



$H(P) = 0.971$

$f_i > \tau_i$

Y → +ve
N → -ve

$\boxed{+ve}$     $\boxed{-ve}$          $\boxed{+ve}$     $\boxed{-ve}$

$H(y) = 0.65$     $H(y) = 0.81$        $H(y) = 0.86$     $H(y) = 0.92$

⤷ Need of combining these values

※ Overall Entropy at the children level :

$$= \frac{n_1}{n} \cdot H_{C_1} + \frac{n_2}{n} \cdot H_{C_2}$$

$H_{C_1} = \frac{60}{100} \times 0.65 + \frac{40}{100} \times 0.8$       $H_{C_j} = \frac{70}{100} \times 0.86 + \frac{30}{100} \times 0.92$

$H_{C_i} = 0.71$                                       $H_{C_j} = 0.88$

Entropy drop with feature-i              Entropy drop with $j^{th}$

$\quad = 0.971 - 0.71$                        feature $= 0.971 - 0.88$

$\quad = 0.261$                                       $= 0.091$

$\boxed{\text{This drop in Entropy is also know as Information Gain}}$