## Today's agenda

1) Recap - Quizzes ✓
2) Employee Attrition Dataset — graphic ✓
3) Purity of Nodes & Entropy ✓
4) Plot for entropy ✓
5) Weighted Entropy ✓
6) Gini Impurity ✓
7) Comparing Gini Impurity with Entropy ✓
8) Code Walkthrough (Time Permits)

_____ X _____ X _____

## Table :

### Splitting Categorical Variables

|  | Gender | | Total |
|---|---|---|---|
| Target Variable | Male | Female | |
| 0 (Stays) | 50 | 10 | 60 |
| 1 (Churn) | 20 | 20 | 40 |
| Total | 70 | 30 | 100 |

I)

| Target Variable | Age | | Total |
|---|---|---|---|
| | < 35 | > 35 | |
| 0 (Stays) | 50 | 10 | 60 |
| 1 (Churn) | 10 | 30 | 40 |
| Total | 60 | 40 | 100 |

II)

Total 100 data points

overtime

② | 2.5 hrs

+ (+) + (+)   ‒ (‒) (+) +
+ +  (+)    ‒ (+) + +
                     (+) +
‒ ‒ ‒   ‒   (+) +
✓ (‒) ‒  ‒   ✓  (+) +
‒ ‒   ‒   ‒    +

| 29 | ① | 35 | ②    Age

+ → churn (leave)
‒ → stay

y
→ ‖ to x axis
→ ‖ to y axis
x

age < 29

Yes            No        age >= 29

Overtime < 2.5         age < 35

Yes      No       Yes      No      age >= 35

I | Stay    II | Churn    III | Stay    IV | Churn

| If else |

overtime

+ + +        ‒ ‒      + +
+ + +        ‒   ‒    + +
2.5 hrs ┈┈┈┈     ‒           + +
‒ ‒ ‒     ‒   ‒      + +
‒ ‒   ‒    ‒        + +
‒ ‒   ‒     ‒        +

        29      35        Age

training data $\left[\begin{array}{l} y_+ = 40 \\ y_- = 60 \end{array}\right.$ ◯ → $\boxed{100 \text{ data-points}}$

Q) $\underline{\text{test data}}$

└→ predict the label

$\boxed{f_1, f_2, f_3, \ldots f_d}$ ✗ $\boxed{\begin{array}{c} y \\ \hline + \\ \hline - \end{array}}$

100 ↻

40 +ve

60 -ve

majority class

$y_+ = 40$
$y_- = 60$ ◯  $\boxed{f_j < \tau_1}$

Yes                          No

$\boxed{\begin{array}{l} y_+ = 10 \\ y_- = 60 \end{array}}$ ◯ ① Left

◯ $y_+ = 30$ ② Right  $\boxed{\text{Node}}$ → definite decision

$y_- = 0$ → $\boxed{\text{pure node}}$ → just 1 class

1
$\frac{10}{70} = \frac{1}{7}$ → +ve

$\frac{60}{70} = \frac{6}{7}$ → -ve

$\frac{30}{30} = 1$ → +ve

$\frac{0}{30} = 0$ → -ve

Purity →

Impurity →

$\boxed{\text{impure node}}$ → more than 1 class
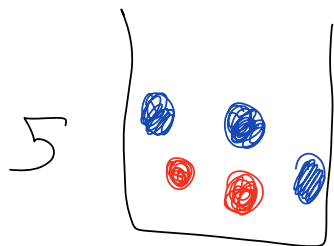
2 → pure

1 → impure

Q) How do we measure impurity & purity?

purity + impurity = 1

$$\boxed{\text{purity} = 1 - \text{impurity}}$$ ✓

$\boxed{\text{Entropy}} \longrightarrow$ measure of randomness
$\downarrow$
impurity
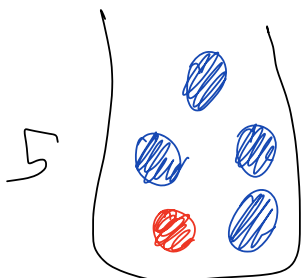


Ⓘ $\qquad p(blue) = \frac{3}{5}, \quad p(red) = \frac{2}{5}$
$\overset{\nearrow}{\underset{P}{}}{}^{P_1} \qquad P_2 = 1 - P_1 \qquad\qquad P_2$

Entropy $= - \left[ p \log_2 (p) + (1-p) \log_2 (1-p) \right]$
$\qquad = - \left[ 0.6 \log_2 (0.6) + 0.4 \log_2 (0.6) \right]$
$\qquad = 0.97$

Ⓘ $\qquad p(blue) = 4/5, \quad \underline{p(red) = 1/5} = p$

Entropy $= - \left[ 0.8 \log_2 (0.8) + 0.2 \log_2 (0.2) \right]$
$\qquad\qquad = 0.72$

Logistic Regression log-loss

$$- \left[ y \log_2 (p) + (1-y) \log_2 (1-p) \right]$$

$y = (0,1)$ $\qquad \underset{\downarrow}{} $

$p = (0-1)$ $\qquad p$

KL Divergence : $\left. \begin{array}{c} \text{log-loss} \\ \text{entropy} \end{array} \right]$

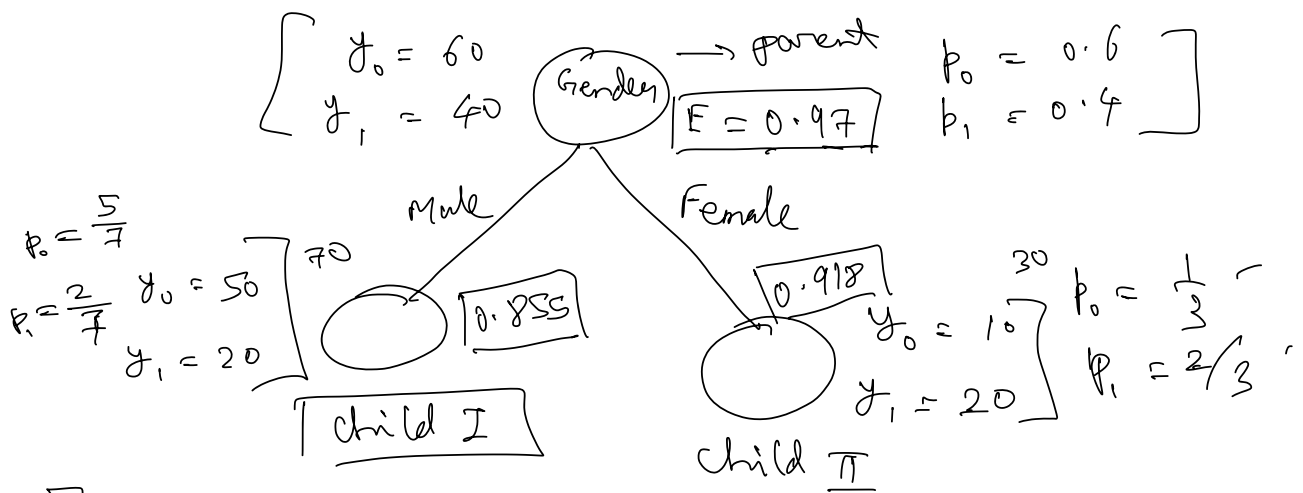$$-\left[ p \log_2 (p + \varepsilon) + (1-p) \log_2 (1-p+\varepsilon) \right]$$

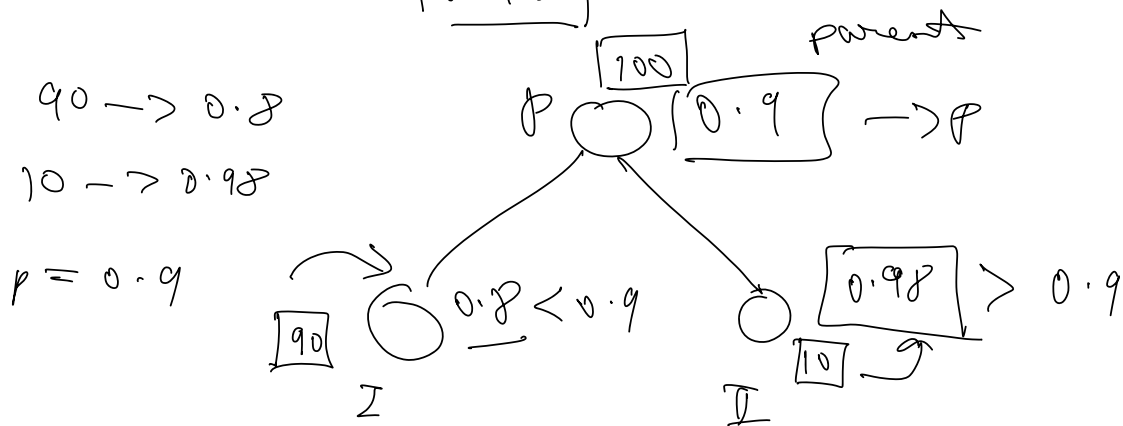$$\varepsilon = 10^{-6}$$

## Graph of Entropy



Entropy $\uparrow 1$

$0$ ——— $2$

$p \longrightarrow$

I) $p = 1$, $\boxed{p = 0}$

$1 - p = 0$

Entropy $= -\left[ 1 \times \log_2 (1) + 0 \times \log_2 (0) \right]$

$\quad\quad\quad\quad = 0$

II) $p = 0.5$

$1 - p = 0.5$

$E = -\left[ 0.5 \times \log_2 (0.5) + 0.5 \times \log_2 (0.5) \right]$

$\quad = -\left[ 1 \times \log_2 (0.5) \right]$

$\quad = -\left[ 1 \times \log_2 (2^{-1}) \right]$

$\quad = -\left[ 1 \times (-1) \right] = 1 \checkmark$

$$\begin{bmatrix} y_0 = 60 & & \rightarrow \text{parent} & p_0 = 0.6 \\ y_1 = 40 & \text{Gender} & \boxed{E = 0.97} & p_1 = 0.4 \end{bmatrix}$$

Male        Female

$p_0 = \dfrac{5}{7}$

$p_1 = \dfrac{2}{7}$   $\begin{array}{c} 70 \\ y_0 = 50 \\ y_1 = 20 \end{array}$   $\boxed{0.855}$     $\boxed{0.918}$    $\begin{array}{c} 30 \\ y_0 = 10 \\ y_1 = 20 \end{array}$   $p_0 = \dfrac{1}{3}$   $p_1 = 2/3$

$\boxed{\text{child I}}$        child II

$$\begin{aligned}
E_{\text{parent}} &= -\left[ 0.6 \log_2 (0.6) + 0.4 \log_2 (0.6) \right] \\
&= \underline{0.97} \\
E_{\text{child I}} &= -\left[ 0.72 \log_2 (0.72) + 0.28 \log_2 (0.28) \right] \\
&= \boxed{0.855} \\
E_{\text{child II}} &= -\left[ \tfrac{1}{3} \log_2 \left(\tfrac{1}{3}\right) + \tfrac{2}{3} \log_2 \left(\tfrac{2}{3}\right) \right] \\
&= \boxed{0.918}
\end{aligned}$$

$\begin{array}{c} \boxed{100} \\ P \bigcirc \boxed{0.9} \end{array}$   parent $\rightarrow P$

$90 \rightarrow 0.8$

$10 \rightarrow 0.98$

$p = 0.9$     $\boxed{90} \bigcirc \, 0.8 < 0.9$     $\boxed{0.98} > 0.9$

      I            $\boxed{10}$   II

$$\begin{aligned}
\text{weighted entropy} &= \frac{90}{100} \times 0.8 + \frac{10}{100} \times 0.98 \\
&= \boxed{0.818} < \boxed{0.9} \rightarrow P
\end{aligned}$$

$$WE_{children} = \frac{70}{100} \times 0.855 + \frac{30}{100} \times 0.918$$

$$= \boxed{0.874} \quad < \quad \boxed{0.97} \longrightarrow P$$

Information Gain = Reduction in Entropy
(IG)

$$= 0.97 - 0.874$$

Split objective $= \boxed{\sim 0.1} \longrightarrow IG$

$\longrightarrow$ (1)

is maximize IG

## Gini Impurity

$$GI = 1 - [\, p(0)^2 + p(1)^2 \,]$$

$$E = -[\, p \log(p) + (1-p) \log(1-p) \,]$$

$$\boxed{p = p(0)}$$

$\downarrow$

Compute intensive

$\boxed{GI}$  I: $p(0) = 1$, $p(1) = 0$

$$GI = 1 - [\, 1^2 + 0^2 \,] = 0$$

II: $p(0) = 0.5$  $p(1) = 0.5$

$$GI = 1 - [\, 0.5^2 + 0.5^2 \,] = 0.5$$

GI

$0.5$ - - - - - -

$0$        $0.5$        $1.0$        $P(0) \longrightarrow$        $Z$

$E$

$1$ - - - - - -

$0$        $0.5$        $Z$        $1$

$y_0 = 60$                    $E = 0.97$
$y_1 = 40$    Age $< 35$    $\longrightarrow$    $\geq 35$

$y_0 = 50$  $\boxed{60}$                    $\boxed{40}$
$y_1 = 10$          I                    $y_0 = 10$    $P_0 = 1/4$
                                   II    $y_1 = 30$    $P_1 = 3/4$
$P_0 = \frac{5}{6}$
$P_1 = 1/6$

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

$$= 0.65$$

$$E_{II} = - \left[ \frac{1}{4} \log_2 \left( \frac{1}{4} \right) + \frac{3}{4} \log_2 \left( \frac{3}{4} \right) \right]$$

$$= 0.811$$

$$W E_{I \& II} = \frac{60}{100} \times 0.65 + \frac{40}{100} \times 0.811$$

$$= 0.7144$$

$$I.G = 0.97 - 0.7144 = \boxed{0.2566} \longrightarrow ?$$

$0.2566 > \boxed{0.1}$

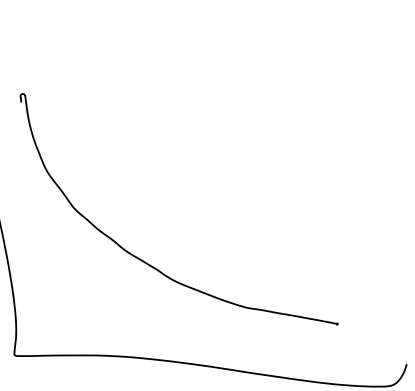## GI for Gender

$$GI_P = 1 - [0.6^2 + 0.4^2]$$
$$= \boxed{0.48}$$

$$GI_I = 1 - \left[ \left(\frac{5}{7}\right)^2 + \left(\frac{2}{7}\right)^2 \right]$$
$$= 0.408$$

$$GI_{II} = 1 - \left[ \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right]$$
$$= 0.444$$

$$W\,GI_{I\&II} = \frac{70}{100} \times 0.408 + \frac{30}{100} \times 0.444$$
$$= \boxed{0.4188}$$

$$IG = 0.48 - 0.4188 = \boxed{0.0612}$$

## GI for Age < 35

$$GI_P = 0.48$$
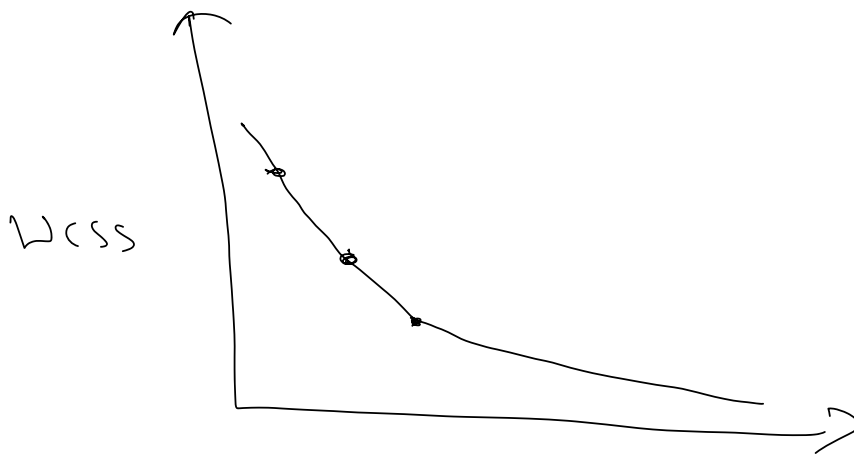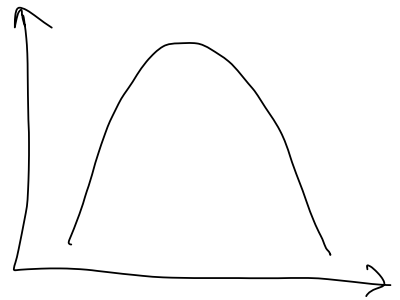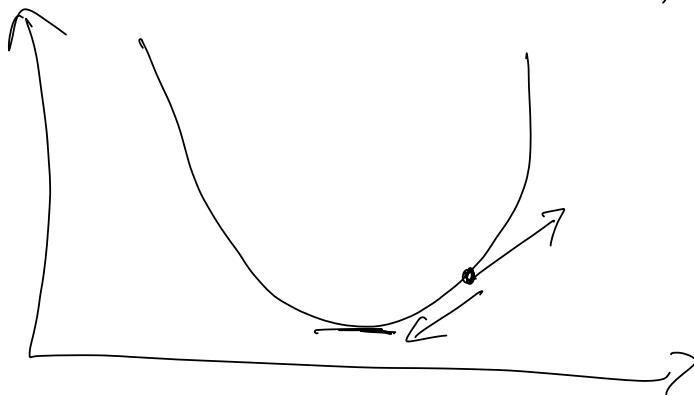
$$GI_I = 1 - \left[ \left(\frac{5}{6}\right)^2 + \left(\frac{1}{6}\right)^2 \right]$$
$$= 0.277$$

$$GI_{II} = 1 - \left[ \left(\frac{1}{4}\right)^2 + \left(\frac{3}{4}\right)^2 \right]$$
$$= 0.375$$

$$W\,GI_{I\&II} = \frac{6}{10} \times 0.277 + \frac{4}{10} \times 0.375 = 0.3162$$

$$IG = 0.48 - 0.3162 = \boxed{0.1638} \checkmark$$

"age < 35 ⟶ better split"





WCSS





$k \longrightarrow$

plot wcss vs $\boxed{\text{initialization}}$

$\downarrow$

$d$ dimensional vector

$\boxed{WCSS}$ vs $d_1, d_2, d_3 \cdots$

WCSS

$d \longrightarrow$