

(1)

(12)

| DataSet | T_{10} | T_{100} |
|---------|----------|-----------|
| A | 0.86 | 0.97 |
| B | 0.84 | 0.77 |

1) from the table we can assume T_{100} has low train Error & high test Error implying Overfit.
from this I conclude that T_{10} predicts better on unseen data

2) It's like I would choose T_{10} over T_{100} as the total accuracy of T_{100} is only slightly better than T_{10} . So choosing T_{10} is computationally good decision

⑧

②

a) Split A

$$\begin{array}{r} T \quad F \\ + \quad 25 \quad 25 \\ - \quad 0 \quad 50 \end{array}$$

$$IG = E_{10+91} - E_{split}$$

$$= 0.5 - E_{split}$$

$$E_{split} = E_{A=T} + E_{A=F}$$

$$= \left(\frac{1}{2} \right) \left(\frac{25}{25} \right) + \left(\frac{1}{2} \right) \left(\frac{50}{75} \right)$$

$$= \frac{1}{2} \left(\frac{25}{25} + \frac{50}{75} \right)$$

$$E_{A=T} = 1 - \max\left(\frac{25}{25}, \frac{0}{25}\right)$$

$$= 1 - 1 = 0$$

$$E_{A=F} = 1 - \max\left(\frac{25}{75}, \frac{50}{75}\right)$$

$$= 1 - \frac{50}{75} = \frac{25}{75}$$

$$IG = 0.5 - \frac{1}{2} \left(\frac{25}{25} + \frac{50}{75} \right)$$

$$= 0.5 - \left(\frac{1}{2} \right) \left(\frac{25}{25} + \frac{50}{75} \right)$$

$$= 0.5 - 0. \frac{75}{100} \times \frac{25}{75}$$

$$= \underline{\underline{0.25}}$$

Split at B

| | T | F |
|---|----|----|
| + | 30 | 20 |
| - | 20 | 30 |

$$E_{B=T} = 1 - \max\left(\frac{30}{50}, \frac{20}{50}\right) \Rightarrow 1 - \frac{3}{5} = \frac{2}{5} = 0.4$$

$$E_{B=F} = 1 - \max\left(\frac{20}{50}, \frac{30}{50}\right) \Rightarrow 1 - \frac{3}{5} = 0.4$$

$$IG = 0.5 - \frac{50}{100}(0.4) - \frac{50}{100}(0.4)$$

$$= 0.5 - (0.2) - 0.2 = 0.1$$

$$IG = 0.1$$

$$IG = 0.1$$

Split at C

| | T | F |
|---|----|----|
| + | 25 | 25 |
| - | 25 | 25 |

$$E_{C=T} = 1 - \max\left(\frac{1}{2}, \frac{1}{2}\right) = 0.5$$

$$E_{C=F} = 1 - \max\left(\frac{1}{2}, \frac{1}{2}\right) = 0.5$$

$$IG = 0.5 - \frac{50}{100}\left(\frac{1}{2}\right) - \frac{50}{100}\left(\frac{1}{2}\right)$$

$$= 0.5 - 0.5 = 0$$

$$IG_c = 0$$

Feature "A" is chosen as it has highest Gain.

⑥

When $A=T$

| B \ C | T | F |
|-------|----|---|
| T | 5 | 0 |
| F | 20 | 0 |

Seeing above table as it contains only "+" class label no
Split is required

When $A=F$

| B \ C | T | F |
|-------|---|----|
| T | 0 | 20 |
| F | 0 | 5 |

$$= 1 - \max\left(\frac{25}{75}, \frac{50}{75}\right)$$

$$= 1 - \frac{50}{75}$$

$$\text{Original } A=F = \frac{25}{75} = \frac{1}{3}$$

Split at B

| | T | F |
|-----|----|----|
| B=T | 25 | 0 |
| B=F | 20 | 30 |

$$E_{B=T} = \frac{25}{75} \left(1 - \max\left(\frac{25}{45}, \frac{20}{45}\right) \right) = \frac{1-5}{9} = \frac{4}{9}$$

$$E_{B=F} = 1 - \max\left(\frac{0}{30}, \frac{30}{30}\right) = 0$$

④

(3)

$$E_{\text{original}} = \frac{45}{75} \left(\frac{4}{1} \right) - \frac{30}{75} (0)$$

$$\frac{1}{3} - \frac{4}{15} = \frac{15-12}{3(15)} = \frac{1}{15}$$

$$IG_B = \frac{1}{15}$$

Split at C

| | T | F |
|---|----|----|
| + | 0 | 25 |
| - | 25 | 25 |

$$E_{C=T} = 1 - \max\left(\frac{0}{25}, \frac{25}{25}\right) = 0 //$$

$$E_{C=F} = 1 - \max\left(\frac{25}{50}, \frac{25}{50}\right) = 0.5$$

$$IG_C = E_{\text{original}} - \frac{25}{75} (0) - \frac{50}{75} (0.5)$$

$$= \frac{1}{3} - \frac{2}{3} (0.5)$$

$$IG_C = 0$$

Feature 'B' will be chosen as it has highest IG.

(c)

$$\text{Error rate} = 0.2$$

$$50 \times 0.2 \times 100 = 20$$

20 instances are misclassified.

(d)

$C = T$

| | A | B | + | - |
|---|---|---|----|----|
| T | T | T | 5 | 0 |
| F | T | F | 0 | 20 |
| T | F | T | 20 | 0 |
| F | F | F | 0 | 5 |

$$\epsilon_{\text{split}} = \epsilon_{\text{original}} = 1 - \max\left(\frac{25}{50}, \frac{25}{50}\right) = \frac{25}{50} = \frac{1}{2}$$

Split A

| | T | F |
|---|----|----|
| + | 25 | 0 |
| - | 0 | 25 |

$$\epsilon_{A=T} = 1 - \max\left(\frac{25}{25}, \frac{0}{25}\right)$$

$$= 1 - 1 = 0$$

$$\epsilon_{A=F} = 1 - \max\left(\frac{0}{25}, \frac{25}{25}\right)$$

$$= 1 - 1 = 0$$

(7)

$$- \text{Congrat} = \frac{25}{50}(0) - \frac{25}{50}(0)$$

$$\frac{1}{2} - 0 - 0 = \frac{1}{2}$$

$$\text{Congrat} \quad IGA = \frac{1}{2} = 0.5$$

Split B

| | Congrat | F |
|---|--------------------|----|
| T | 5 | 20 |
| - | 20 | 5 |

$$\text{Congrat} = \frac{1}{2}$$

$$C_{B=T} = 1 - \max\left(\frac{5}{25}, \frac{20}{25}\right)$$

$$= 1 - \frac{20}{25} = \frac{5}{25}$$

$$C_{B=F} = 1 - \max\left(\frac{20}{25}, \frac{5}{25}\right)$$

$$= 1 - \frac{20}{25} = \frac{5}{25}$$

$$IG_B = \frac{1}{2} - \frac{25}{50} \left(\frac{5}{25} \right) - \frac{25}{50} \left(\frac{5}{25} \right)$$

$$= \frac{1}{2} - \frac{2}{10} = \frac{5-2}{10} = \frac{3}{10}$$

$$IG_B = 0.3$$

$IG_A > IG_B$ So A is chosen for next Split

for $C=F$

| A | B | + | - |
|---|---|----|----|
| T | T | 0 | 0 |
| F | T | 25 | 0 |
| T | F | 0 | 0 |
| F | F | 0 | 25 |

When A

$$E_{\text{original}} = 1 - \max\left(\frac{25}{50}, \frac{25}{50}\right) \\ = 1 - \frac{25}{50} = \frac{1}{2}$$

When $A=F$

| | T | F |
|---|---|----|
| + | 0 | 25 |
| - | 0 | 25 |

$$E_{A=T} = 1 - \max\left(\frac{0}{50}, \frac{0}{50}\right) \\ = 1$$

$$E_{A=F} = 1 - \max\left(\frac{25}{50}, \frac{25}{50}\right) \\ = 1 - \frac{1}{2} = \frac{1}{2}$$

$$I.G. = \frac{1}{2} - \frac{0}{50}(1) = \frac{50}{50}\left(\frac{1}{2}\right)$$

$$I.G._A = 0$$

Split A & B
When

B=T

| | T | F |
|---|----|----|
| T | 25 | 0 |
| F | 0 | 25 |

$$\text{Entropy} = 0.5$$

$$E_{B=T} = 1 - \max\left(\frac{25}{50}, \frac{0}{50}\right) = 1 - \frac{1}{2} = 0.5$$

$$E_{B=F} = 1 - \max\left(\frac{0}{50}, \frac{25}{50}\right) = 1 - \frac{1}{2} = 0.5$$

$$IG_B = 0.5 - 0 - 0 = 0.5$$

$IG_B > IG_A$ so Attribute B is chosen.

① The Greedy Algorithm does not always give the best result.

Question 5

for Split A

① When ~~B=T~~ When A=F

| | | |
|---|-----|-----|
| | B=T | B=F |
| + | 3 | 1 |
| - | 1 | 5 |

| | | |
|-----|-----|-----|
| | A=T | A=F |
| B=T | 4 | 0 |
| B=F | 3 | 3 |

$$\epsilon_{\text{original}} = -[0.4 \log_2 0.4 + 0.6 \log_2 0.6]$$

$$= 0.9209$$

IG at Split A

$$\epsilon_{\text{original}} - (\epsilon_{A=T}) - (\epsilon_{A=F})$$

$$\epsilon_{A=T} = -\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7} = 0.985$$

$$\epsilon_{A=F} = -0 - \frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$IG_A = \epsilon_{\text{orig}} - \epsilon_{A=T} - \epsilon_{A=F}$$

$$= 0.9209 - \frac{7}{10}(0.985)$$

$$= 0.9209 - 0.6895$$

$$IG = 0.2814$$

11

$$E_{B=T} = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.811$$

$$E_{B=F} = -\frac{1}{6} \log \frac{1}{6} - \frac{5}{6} \log \frac{5}{6} = 0.65$$

$$IG_B^{\text{original}} = \frac{4}{10} (0.811) + \frac{6}{10} (0.65)$$

$$= \cancel{0.9709} 0.9709$$

$$IG_B = 0.2535$$

$IG_A > IG_B$ So feature A is chosen

⑤ $G_{\text{original}} = 1 - (\text{positive prob.})^2 - (\text{negative prob.})^2$

$$= 1 - (0.4)^2 - (0.6)^2$$

$$= 0.48$$

Gain when Split At "A"

$$G_{A=T} = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 0.4897$$

$$G_{A=F} = 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0$$

$$IG_A = 0.48 - \frac{7}{10} (0.4897) - \frac{3}{10} (0) = 0.137$$

Gain when split at 3

(12)

$$G_{B=1} = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.375$$

$$G_{B=2} = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 = 0.277$$

$$\begin{aligned}\Delta &= 0.48 - \left(\frac{4}{10}\right)(0.375) - \frac{6}{10}(0.277) \\ &= 0.48 - 0.3162 \\ &= 0.1638\end{aligned}$$

$I_{G_B} > I_{G_A}$ so feature B is chosen

(1) Yes, As we can see from the figure as the respective measures and monotonically increases, have different values even though it has same range, which is shown in the figure.

3

a)

Probability of + class = $4/9$ Probability of - class = $5/9$

$$\begin{aligned}
 H(x) &= - \sum (P_i \times \log_2 P_i) \\
 &= - \left(\left(\frac{4}{9} \right) \times \log_2 \left(\frac{4}{9} \right) \right) - \left(\left(\frac{5}{9} \right) \times \log_2 \left(\frac{5}{9} \right) \right) \\
 &= 0.991
 \end{aligned}$$

b)

for feature a_1

$a_1 = \text{True}$ 3 Positive 1 Negative
 $a_1 = \text{False}$ 1 Positive 4 Negative

$$\begin{aligned}
 &\left(\frac{4}{9} \times \left(-\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) - \left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) \right) - \left(\frac{5}{9} \times \left(-\frac{1}{5} \log_2 \left(\frac{1}{5} \right) - \frac{4}{5} \log_2 \left(\frac{4}{5} \right) \right) \right) \\
 &= 0.7616
 \end{aligned}$$

$$IG = \text{Entropy}_{\text{Total}} - \text{Entropy}_{a_1}$$

$$= 0.991 - 0.7616$$

$$= 0.2294$$

for feature a_2

$a_2 = \text{True}$ 2 Positive 3 Negative
 $a_2 = \text{True}$ 2 Positive 3 Negative

$$\left(\left(\frac{5}{9} \right) * \left(-\left(\frac{2}{5} \right) * \log_2 \left(\frac{2}{5} \right) - \left(\frac{3}{5} \right) * \log_2 \left(\frac{3}{5} \right) \right) + \frac{4}{9} \left(\left(\frac{2}{4} \right) * \log_2 \left(\frac{2}{4} \right) + \left(\frac{2}{4} \right) * \log_2 \left(\frac{2}{4} \right) \right) \right)$$

$$= 0.9838$$

$$IG = \text{Entropy}_{\text{Total}} - \text{Entropy}_{a_1}$$

$$= 0.991 - 0.9838$$

$$= 0.0072$$

③

| 1 | 2 | 4 | 5 | 6 | 7 | 8 |
|-----|---|-----|-----|-----|-----|-----|
| 0.5 | 2 | 3.5 | 4.5 | 5.5 | 6.5 | 7.5 |
| 0 | 4 | 1 | 3 | 1 | 3 | 2 |
| 0 | 5 | 0 | 5 | 1 | 4 | 1 |

$$\text{Split 1} = - \left[\left(\frac{4}{9} \right) \log_2 \left(\frac{4}{9} \right) + \left(\frac{5}{9} \right) * \log_2 \left(\frac{5}{9} \right) \right]$$

$$IG = 0$$

$$\text{Split 2} \stackrel{\text{WTA}}{=} - \frac{1}{9} \left[7 \log_2(1) + 0 \log_2(0) \right] + \left(\frac{8}{9} \right) \left[\left(\frac{3}{8} \right) \log_2 \left(\frac{3}{8} \right) + \left(\frac{5}{8} \right) \log_2 \left(\frac{5}{8} \right) \right]$$

$$IG = 0.1428$$

$$\text{Split 3} = -\left(\frac{2}{9}\right) \left[\frac{1}{2} \log\left(\frac{1}{2}\right) + \frac{1}{2} \log\left(\frac{1}{2}\right) \right] - \left(\frac{7}{9}\right) \left[\frac{3}{7} \log\left(\frac{3}{7}\right) + \left(\frac{4}{7}\right) \log\left(\frac{4}{7}\right) \right] \text{①}$$

$$\text{I.G} = 0.00248$$

$$\text{Split 4} = -\left(\frac{3}{9}\right) \left[\left(\frac{2}{3}\right) \log\left(\frac{2}{3}\right) + \frac{1}{3} \log\left(\frac{1}{3}\right) \right] + \left(\frac{6}{9}\right) \log\left(\frac{2}{6}\right) + \left(\frac{1}{6}\right) \log\left(\frac{4}{6}\right)$$

$$\text{I.G} = 0.0727$$

$$\text{Step-5} = -\left(\frac{5}{9}\right) \left[\frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right] - \frac{4}{9} \left[\frac{2}{4} \log\left(\frac{2}{4}\right) + \frac{2}{4} \log\left(\frac{2}{4}\right) \right]$$

$$\text{I.G} = 0.0071$$

$$\text{Step-6} = -\left(\frac{6}{9}\right) \left[\left(\frac{3}{6}\right) \log\left(\frac{3}{6}\right) + \frac{3}{6} \log\left(\frac{3}{6}\right) \right] - \frac{3}{9} \left[\frac{1}{3} \log\left(\frac{1}{3}\right) + \frac{2}{3} \log\left(\frac{2}{3}\right) \right]$$

$$\text{I.G} = 0.01$$

$$\text{Step-7} = -\frac{8}{9} \left[\left(\frac{4}{8}\right) \log\left(\frac{4}{8}\right) + \frac{4}{8} \log\left(\frac{4}{8}\right) \right] - \frac{1}{9} \left[\left(\frac{0}{1}\right) \log\left(\frac{0}{1}\right) + \left(\frac{1}{1}\right) \log\left(\frac{1}{1}\right) \right]$$

$$\text{I.G} = 0.1022$$

from all the splits "Split 2" has highest I.G.

(d) The best split is the one which has ^{highest} ~~best~~ Information gain ⁽¹⁶⁾
from given " a_1 ", " a_2 ", " a_3 " " a_1 " has Information gain.

(e) misclassification error: $1 - \max\left(\frac{7}{9}, \frac{2}{9}\right) = 1 - \frac{7}{9} = \frac{2}{9}$
misclassification error: $1 - \max\left(\frac{5}{9}, \frac{4}{9}\right) = 1 - \frac{5}{9} = \frac{4}{9}$

(f) $G_{(a_1)} = 1 - \left(\frac{7}{9}\right)^2 - \left(\frac{2}{9}\right)^2 = 0.34$
 $G_{(a_2)} = 1 - \left(\frac{4}{9}\right)^2 - \left(\frac{5}{9}\right)^2 = 0.493$

7

$$\text{misclassification Error}_x = 1 - \max\left(\frac{40}{80}, \frac{40}{80}\right) = 1 - \frac{1}{2} = 0.5$$

$$\text{misclassification Error}_y = 1 - \max\left(\frac{60}{100}, \frac{40}{100}\right) = 1 - \frac{60}{100} = 0.4$$

$$\text{misclassification Error}_z = 1 - \max\left(\frac{70}{100}, \frac{30}{100}\right) = 1 - \frac{70}{100} = 0.3$$

We pick z as it has lowest error rate i.e. 0.3

$x=0 \quad z=0$

$$\text{misclassification error} = 1 - \max\left(\frac{15}{60}, \frac{45}{60}\right) = 1 - \frac{45}{60} = 0.25$$

$$\text{misclassification error} = 1 - \max\left(\frac{15}{40}, \frac{25}{40}\right) = \frac{3}{8} = 0.375$$

$x=1 \quad z=0$

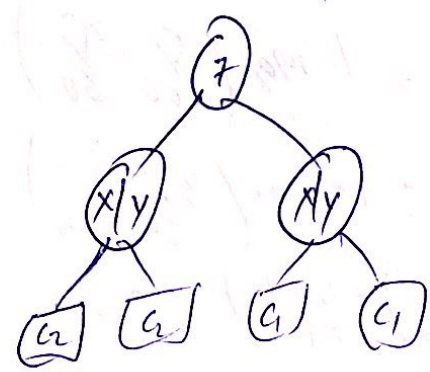
$$\text{misclassification error} = 1 - \left(\frac{15}{60}, \frac{45}{60}\right) = 0.25$$

$$\text{misclassification error} = 1 - \left(\frac{15}{60}, \frac{25}{40}\right) = 0.375$$

As the both classification error are same the next split

can be either x or y

$$\text{weighted avg} = \left(\frac{60}{100} \times \frac{1}{4}\right) + \left(\frac{40}{100} \times \frac{3}{8}\right) = 0.3$$



$$y=1 \quad z=1$$

$$\text{misclassification Error} = 1 - \max\left[\frac{45}{60}, \frac{15}{60}\right] = \left[1 - \left(\frac{3}{4}\right)\right] = \frac{1}{4} = 0.25$$

$$y=0 \quad z=1$$

$$\text{misclassification Error} = 1 - \max\left[\frac{25}{40}, \frac{15}{40}\right] = 1 - \left(\frac{5}{8}\right) = 0.375$$

$$\text{Weighted Avg.} = \left[\left(\frac{60}{100}\right) \times \left(\frac{1}{4}\right)\right] + \left[\left(\frac{40}{100}\right) \times \left(\frac{3}{8}\right)\right] = 0.3$$

$$\leftarrow \text{Total Weighted Error} = 0.3 + 0.3 = 0.6$$

⑤

$$x=0 \quad y=0 \quad \text{Misclassify Error} = 1 - \max\left(\frac{55}{60}, \frac{5}{60}\right) = 0.083$$

$$x=0 \quad y=1 \quad \text{Misclassify Error} = 1 - \max\left(\frac{5}{60}, \frac{55}{60}\right) = 0.083$$

$$\text{Weighted Average: } \frac{1}{2}(0.083) + \frac{1}{2}(0.083) = \underline{0.083}$$

$$x=0 \quad z=0 \quad \text{Misclassify Error} = 1 - \max\left(\frac{45}{60}, \frac{15}{60}\right) = \frac{1}{4}$$

$$x=0 \quad z=1 \quad \text{Misclassify Error} = 1 - \max\left(\frac{15}{60}, \frac{45}{60}\right) = \frac{1}{4}$$

$$\text{Weighted Average} = \left(\frac{1}{2}\right)\left(\frac{1}{4}\right) + \left(\frac{1}{2}\right)\left(\frac{1}{4}\right) = \underline{\frac{1}{4} = 0.25}$$

$$x=0 \quad y=0 \quad \text{Misclassify Error} = 1 - \max\left(\frac{5}{60}, \frac{35}{60}\right) = \frac{1}{8}$$

$$x=1 \quad y=1 \quad \text{Misclassify Error} = 1 - \max\left(\frac{35}{60}, \frac{15}{60}\right) = \frac{1}{8}$$

$$\text{Weighted Average} = \frac{1}{2}\left(\frac{1}{8}\right) + \frac{1}{2}\left(\frac{1}{8}\right) = 0.125$$

The lowest error is for $x=1$ for "y"

(11)

(c) No, greedy heuristic does not create optimum result.