Assignment: Decision Trees & Random Forest

ID: 2021701010

Question 5:
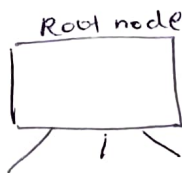
| | Review | Smell | Taste | Portion |
|---|---|---|---|---|
| 1 | Negative | Woody | Sweet | Small |
| 2 | Negative | Fruity | Salty | Large |
| 3 | Negative | Fruity | Salty | Large |
| 4 | Positive | Fruity | Sour | Small |
| 5 | Positive | Woody | Sour | Small |
| 6 | Negative | Woody | Sweet | Large |
| 7 | Positive | Woody | Sour | Large |
| 8 | Positive | Fruity | Salty | Small |
| 9 | Positive | Fruity | Salty | Small |
| 10 | Negative | Woody | Sweet | Large |

↑
To predict

Q

↳ Compute a 'decision- Tree' with the goal to predict the food review based on its (smell), taste & portion size.

a) Compute the entropy of each rule in the first stage.

Ans

Root node



The "entropy" at root node $= -\sum_{i=1}^{k} P(C_i) \log_2 [P(C_i)] = H_{root}$

→ where $C_i = i$th class

$P(C_i) = \dfrac{\# \text{ of samples in Class } i}{\# \text{ total samples}}$

And the entropy after splitting (based on some feature & a particular question about that feature)

$= \underset{\text{no-split}}{\sum_{j=1}} P(D_j) H(D_j)$ , where $P(D_j) = \dfrac{|D_j|}{|D|} = \dfrac{\# \text{ of samples in } D_j}{\text{total samples in } D}$

& $H(b_j)$ = Entropy of the node containing $b_j$.
= Entropy of each child nodes

& ~~Information gain~~:

& no.split = total # splits (= total child nodes formed).

- ## Information gain : = $\Delta z$

$$\Delta i = \left(\text{Entropy before splitting}\right) - \left(\text{Entropy after splitting}\right)$$

$$= \left(\text{Entropy of parent node}\right) - \left(\text{weighted combinations of Entropy of child nodes}\right).$$

$\times$

## Calculation time :

- Entropy at root node (Entropy before splitting)

$$= -P(C_1) \log_2 [P(C_1)] + - P(C_2) \log_2 [P(C_2)]$$

where $C_1$ = +ve class $\}\to$ for Review
$C_2$ = -ve class

$$= - p_1 \log_2 p_1 - p_2 \log_2 p_2$$

$$= -\left(\frac{5}{10}\right) \log_2 \left(\frac{5}{10}\right) - \left(\frac{5}{10}\right) \log_2 \left(\frac{5}{10}\right)$$

$$= 1$$

- Entropy "if data are splitted on the basis of Smell" :

Smell: (10) Data
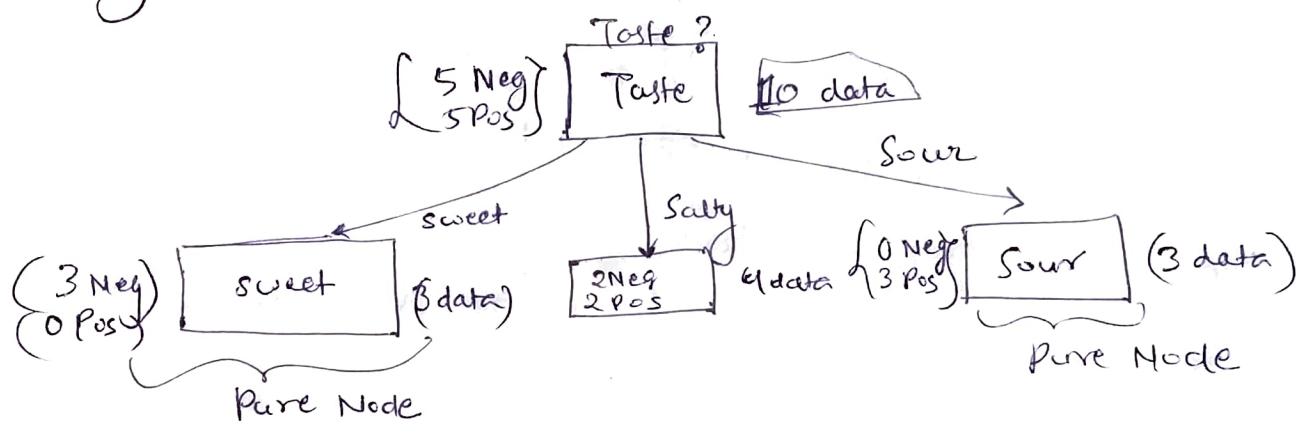Woody?

5 Neg
5 Pos

Woody Yes

Woody No (= fruity)

{ 3 Neg
2 Pos }  Yes   5 (Data)

{ 2 Neg
3 Pos }  No.   5 (Data)

$$\text{Total\_H} = \underbrace{\frac{5}{10}}_{P(D_1)} \underbrace{\left[-\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)\right]}_{H(D_1)} + \underbrace{\frac{5}{10}}_{P(D_2)} \underbrace{\left[-\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)\right]}_{H(D_2)}$$

$$= \quad -\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \frac{2}{5}\log_2\left(\frac{2}{5}\right)$$

$$= \quad -\frac{3}{5}(-0.737) - \frac{2}{5}(-1.322)$$

$$= \quad 0.971$$
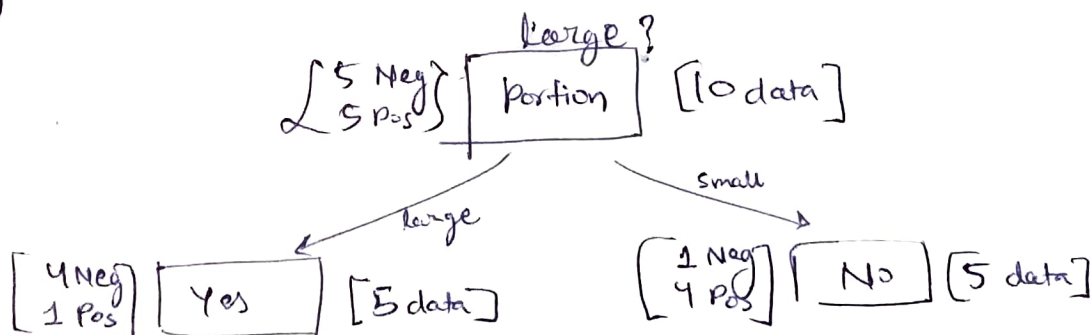
$$\Delta i = \quad 1 - 0.971 = \boxed{0.029} \quad\text{————①}$$

- Entropy " if data are splitted on the basis of Taste "



$$\text{Total } H = \frac{3}{10}\left[-\frac{3}{3}\log_2\left(\frac{3}{3}\right) - \frac{0}{3}\log_2\left(\frac{0}{3}\right)\right] +$$

$$\frac{4}{10}\left[-\frac{2}{4}\log\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right)\right] +$$

$$\frac{3}{10}\left[-\frac{0}{3}\log_2\left(\frac{0}{3}\right) - \frac{3}{3}\log_2\left(\frac{3}{3}\right)\right]$$

$$= \quad 0 + 0.4 + 0. = \quad 0.4$$

Then $\Delta i = \quad 1 - 0.4 = \boxed{0.6} \quad\text{————②}$

- Entropy "if data are splitted on the basis of "postion". )

Total $H$ = Total entropy

$$= \frac{5}{10}\underbrace{\left[-\frac{4}{5}\log_2\left(\frac{4}{5}\right) - \frac{1}{5}\log_2\left(\frac{1}{5}\right)\right]}_{\text{for } \textcircled{Yes}\ =\text{Large}} + \frac{5}{10}\underbrace{\left[-\frac{1}{5}\log_2\left(\frac{1}{5}\right) - \frac{4}{5}\log_2\left(\frac{4}{5}\right)\right]}_{\substack{\text{for } \textcircled{No}\\ =\text{Small}}}$$

$$= \quad 0.72$$

$$\Delta i = \quad 1 - 0.72 = \boxed{0.28} \qquad \text{——} \quad \textcircled{3}$$

Comparing $\textcircled{1}, \textcircled{2} \ \& \ \textcircled{3}$ ; we have $\textcircled{2} > \textcircled{3} > \textcircled{1}$.

$\Rightarrow$ 'Taste' should be at the root node.

$\rightarrow$ Now, once "taste" has took the root node, it's left & right node (i.e., sweet & sour are "pure" nodes & hence the branch won't be extended from them & they will be regarded as leaf nodes).

$\rightarrow$ from "salty" again splitting would be done (with the same process & steps as done for root node) but with other features.



$$\text{Entropy}[\text{salty -node}] = \quad -\frac{2}{4}\log_2\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right)$$

$$= \quad 1.0$$

• taking smell = woody ? question :

Total-H =
$$\frac{0}{4}\left[\underbrace{\phantom{XXXX}}_{\text{for Woody}}\right] + \frac{4}{4}\left[\underbrace{-\frac{2}{4}\log_2\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right)}_{\text{(for fruity)}}\right]$$

$$= 1$$

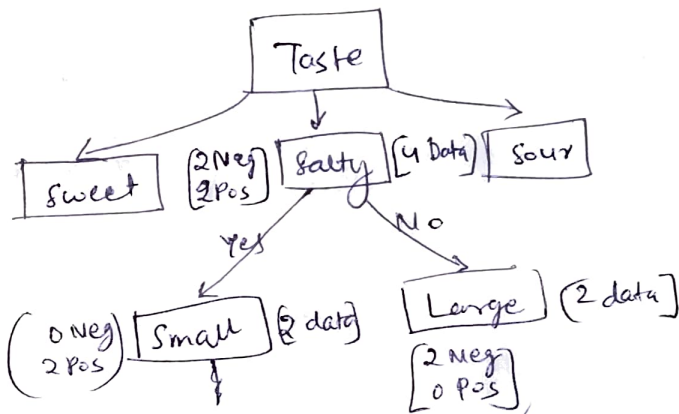$$\Rightarrow \Delta i = 1 - 1 = 0 \quad (\text{Not preferable}) \quad\text{——}④$$

• Taking **portion = Small** ?

Total-H
$$= \frac{2}{4}\left[-\frac{0}{2}\log_2\left(\frac{0}{2}\right) - \frac{2}{2}\log_2\left(\frac{2}{2}\right)\right]$$

$$+ \frac{2}{4}\left[-\frac{2}{2}\log_2\left(\frac{2}{2}\right)\right.$$
$$\left. - \frac{0}{2}\log_2\left(\frac{0}{2}\right)\right]$$

$$= 0$$

$$\Rightarrow \quad \Delta i = \underset{\underset{(\text{at salty})}{\uparrow}}{1 - \overset{\checkmark}{0}} = 1 \quad\text{——}⑤$$



Taste → sweet | (2 Neg / 2 Pos) | salty | (4 Data) | sour

salty: Yes → (0 Neg / 2 Pos) | Small (2 data) ; No → Large (2 data) [2 Neg / 0 Pos]

Both are "pure" (hence "leaf") nodes → indicates termination of tree.

★ Since ⑤ > ④ ⟹ At "salty" node, the question which should be asked = (portion = small)? (i.e. splitting would be done based on feature 'portion')

→ Now, as we have got 'leaf' nodes (i.e pure nodes) our tree has fully grown. So final tree would look like something ~~at back~~ in next page.

$\left\{ \begin{array}{l} 5 \text{ Negative} \\ 5 \text{ positive} \end{array} \right\}$ | Taste | [10 data]

Sweet

$\left\{ \begin{array}{l} 3 \text{ Neg} \\ 0 \text{ Pos} \end{array} \right\}$

Salty

Sour

$\left\{ \begin{array}{l} 0 \text{ Neg} \\ 3 \text{ Pos} \end{array} \right\}$ | +ve | [3 data]

| .... -ve | [3 data]

‖
−ve = Negative
Class
(Pure Node)

$\left\{ \begin{array}{l} 2 \text{ Neg} \\ 2 \text{ Pos} \end{array} \right\}$ | Portion | [4 data]

‖
+ve = Positive
class
(Pure Mode

large

Small

$\left\{ \begin{array}{l} 2 \text{ Neg} \\ 0 \text{ Pos} \end{array} \right\}$ | −ve | [? data]

$\left\{ \begin{array}{l} 0 \text{ Neg} \\ 2 \text{ Pos} \end{array} \right\}$ | +ve | [2 data]

Note : 'Smell' feature remain useless (or doesn't have any impact on final decision label i.e.. Negative or Positive ). ✓