

AS

$$i(N) = - \sum_j P(w_j) \log_2(P(w_j))$$

Date \_\_\_/\_\_\_/\_\_\_

Saat

Total no. of samples = 10

let num-p = no. of positive reviews in the node

~~At root~~ let num-n = no. of negative reviews in node

Deciding split

Initial: num-p = 5  
num-n = 10

$$\therefore i(\text{root}) = - \left( \frac{5}{10} \log_2\left(\frac{5}{10}\right) + \frac{5}{10} \log_2\left(\frac{5}{10}\right) \right) = 1$$

if we split by smell,

left child (woody)

num-p = 2

num-n = 3

$$i(\text{lc}) = - \left( \frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right) \right) = 0.9709$$

right child (sruity)

num-p = 3

num-n = 2

$$i(\text{rc}) = - \left( \frac{3}{5} \log_2\left(\frac{3}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right) \right) = 0.9709$$

$$\therefore \Delta i(N) = 1 - \frac{5}{10} \times 0.9709 - \frac{5}{10} \times 0.9709$$

on splitting based on smell

$$= 0.02904$$

Also,

$$\lim_{x \rightarrow 0} x \log(x) = \lim_{x \rightarrow 0} \frac{\log(x)}{1/x} = \lim_{x \rightarrow 0} \frac{1/x}{-1/x^2} = \lim_{x \rightarrow 0} -x = 0$$

[By L'Hopital's rule]



we split based on taste

left child (sweet)

$$\text{num-p} = 0$$

$$\text{num-n} = 3$$

$$\bar{i}(lc) = -1 \log_2(1) - 0 - 0 = 0$$

right child (salty)

$$\text{num-p} = 2$$

$$\text{num-n} = 2$$

$$\bar{i}(rc) = -\left(\frac{1}{2} \log_2 \frac{1}{2}\right) + \frac{1}{2} \times \log_2 \left(\frac{1}{2}\right) = 1$$

middle child (sour)

$$\text{num-p} = 3$$

$$\text{num-n} = 0$$

$$\bar{i}(\text{middle}) = -1 \log_2(1) - 0 - 0 = 0$$

$$\Delta \bar{i}(N) = \bar{i}(\text{root}) - \frac{3}{10} \times \bar{i}(\text{middle}) - \frac{3}{10} \bar{i}(lc) - \frac{4}{10} \bar{i}(rc)$$

on splitting based on taste

$$= 1 - \frac{3}{10} \times 0 - \frac{3}{10} \times 0 - \frac{4}{10} \times 1$$

Let a node  $N$  be split into  $L$  nodes. Then,

$$\Delta \bar{i}(N) = \bar{i}(N) - \sum_L p_L \bar{i}(N_L)$$

where  $N_L = i^{\text{th}}$  node (child)

$p_L =$  fraction of nodes at  $N$  that were sent to  $N_L$



Date / /

If we split based on portion

Left child (small)

$$\text{num-p} = 4$$

$$\text{num-n} = 1$$

$$i(\text{lc}) = - \left( \frac{1}{5} \log_2 \left( \frac{1}{5} \right) + \frac{4}{5} \log_2 \left( \frac{4}{5} \right) \right) = 0.7219$$

right child (large)

$$\text{num-p} = 1$$

$$\text{num-n} = 4$$

$$i(\text{rc}) = - \left( \frac{1}{5} \log_2 \left( \frac{1}{5} \right) + \frac{4}{5} \log_2 \left( \frac{4}{5} \right) \right) = 0.7219$$

$$\therefore \Delta i(N) = 1 - \frac{5}{10} \times 0.7219 - \frac{5}{10} \times 0.7219 = 0.27807$$

on splitting  
based on  
portion

$$= 0.27807$$

∴ We see that while exploring the information gain (impurity reduction) based on each of the 3 features, the maximum reduction of impurity happens when we split based on taste i.e.,  $\Delta = 0.6$ .

So, assuming that the question asks to greedily choose the split, I decide to split based on taste.

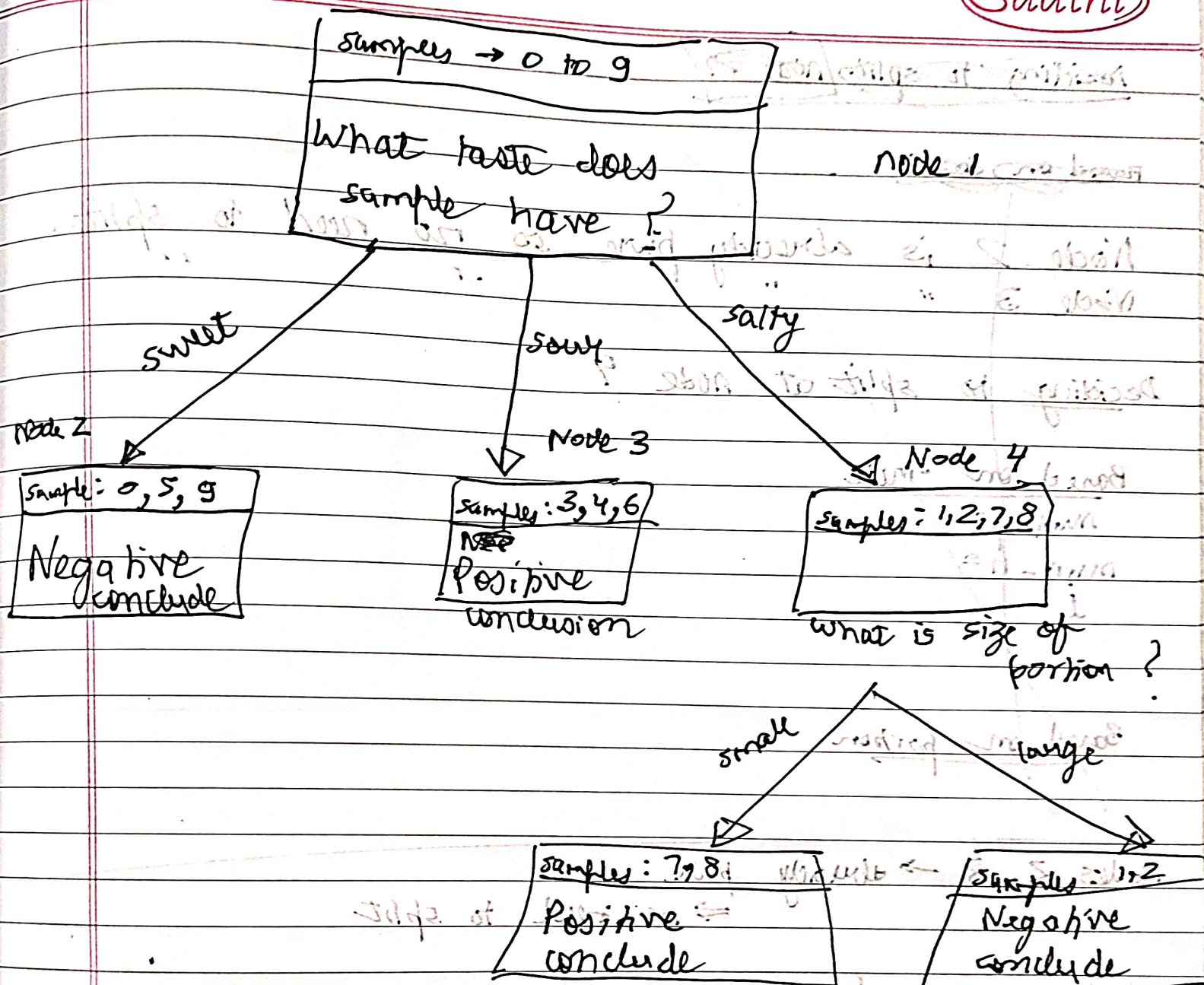
Let samples be indexed as

Let samples have 0 based indexing based on sequence in question.

Date \_\_\_\_ / \_\_\_\_ / \_\_\_\_

107

Saathi



Working for 2<sup>nd</sup> split on negative page



nodes 2, 3  $\rightarrow$  already pure:  $\Rightarrow$  no need to split

### Splitting node 4

Based on small

left child (Woody)

$$\text{num-p} = 0$$

$$\text{num-n} = 0$$

$$\bar{i}(lc) = 0$$

right child (Stuart)

$$\text{num-p} = 2$$

$$\text{num-n} = 2$$

$$\bar{i}(rc) = - \left( \frac{2}{4} \log_2 \left( \frac{1}{2} \right) + \frac{2}{4} \log_2 \left( \frac{1}{2} \right) \right) = 1$$

$$\Delta \bar{i} = 1 - 0 - \frac{1 \times 1}{4} = 0$$

Date \_\_\_/\_\_\_/\_\_\_

109

Saathi

Based on portion

small (left child)

$$\text{num-p} = 2$$

$$\text{num-n} = 0$$

$$i(\text{lc}) = - \left( 0 + \frac{2}{2} \log_2 \left( \frac{2}{2} \right) \right) = 0$$

right child (large)

$$\text{num-p} = 0$$

$$\text{num-n} = 2$$

$$i(\text{rc}) = - \left( \frac{2}{2} \log_2 \left( \frac{2}{2} \right) + 0 \right) = 0$$

$$\Delta i = 1 - 0 - 0 = 1$$

~~✱~~

∴  $\Delta i$  is greatest on splitting based on portion and hence, this is the feature we choose for split on node 4.