① در الگوریتم Hunt، splitting نودها تا جایی ادامه دارد که
همه نمونه‌هایی که در آن قرار می‌گیرند متعلق به یک کلاس باشند، در واقع نمونه‌های هم کلاس قرار
گیرند.

③ a) $1 - P_1^2(\text{train\_set}) - P_0^2(\text{train\_set})$

$$= 1 - \tfrac{1}{4} - \tfrac{1}{4} = \tfrac{1}{2}$$

b) چون ID، nominal است، می‌توانیم از استفاده Multiway split کنیم.
و $I(\text{children})$ را بر اساس ممیس محاسبه کنیم.

$$I = \sum \frac{N_{child}}{N} I_{child} = \sum_{i=1}^{20} \frac{1}{20} \times (1 - 1^2) = 0$$

c)

female — male

$$\bar{I} = \frac{10}{20} \times \left(1 - \left(\tfrac{4}{10}\right)^2 - \left(\tfrac{6}{10}\right)^2\right) + \frac{10}{20} \times \left(1 - \left(\tfrac{6}{10}\right)^2 - \left(\tfrac{4}{10}\right)^2\right)$$

$$= \frac{48}{100}$$

d) $\frac{8}{20} \times (1 - 1^2) + \frac{8}{20} \times \left(1 - \left(\tfrac{7}{8}\right)^2 - \left(\tfrac{1}{8}\right)^2\right) + \frac{4}{20} \times \left(1 - \left(\tfrac{1}{4}\right)^2 - \left(\tfrac{3}{4}\right)^2\right)$

Sports  2  luxury  family

$$= \frac{8}{20} \times \frac{14}{64} + \frac{4}{20} \times \frac{6}{16} = \frac{52}{320} = 0.1625$$

(under $\frac{14}{64}$: $16$)

e) small  $\frac{12}{25}$   large  $\frac{1}{2}$   extra  $\frac{1}{2}$

$\frac{5}{20} \times \left(1 - \left(\tfrac{3}{5}\right)^2 - \left(\tfrac{2}{5}\right)^2\right) + \frac{4}{20} \times \left(1 - \tfrac{1}{4} - \tfrac{1}{4}\right) + \frac{4}{20} \times \left(1 - \tfrac{1}{4} - \tfrac{1}{4}\right)$

$+ \frac{7}{20} \times \left(1 - \left(\tfrac{3}{7}\right)^2 - \left(\tfrac{4}{7}\right)^2\right) = \frac{5}{20} \times \frac{12}{25} + \frac{4}{20} + \frac{7}{20} \times \frac{24}{49}$

medium  $\frac{24}{49}$

$$= \frac{3}{25} + \frac{4}{20} + \frac{1}{35} \simeq 0.348 \cdots$$

f) Car Type بشرطیکه جوی $\sum$ Impurity کمترین راکه دارند و نقطه ها سریعتر خاص

می گوند و عمق درخت کم می شود.

g) با در نظر گرفتن $Z_{child} (ID) = 0$ جوی overfitting می شود زیرا مدل

حفظی ـ یادنی کرد و کلمه حفظ ای کلمه و هیچ داده جدیدی را نمی تواند predict کند.

a)

$$Entropy = - \sum P_i \lg P_i \longrightarrow Z_{parent} \qquad \text{(F)}$$

$$= - \left( P_+ \lg P_+ + P_- \lg P_- \right) = - \frac{4}{9} \lg \frac{4}{9} - \frac{5}{9} \lg \frac{5}{9}$$

b) $I(a_1) = \underbrace{\frac{4}{9} \times \left( -\frac{3}{4} \lg \frac{3}{4} - \frac{1}{4} \lg \frac{1}{4} \right)}_{a_1 = T} + \underbrace{\frac{5}{9} \times \left( -\frac{4}{5} \lg \frac{4}{5} - \frac{1}{5} \lg \frac{1}{5} \right)}_{a_1 = F}$

$Z(a_2) = \underbrace{\frac{5}{9} \times \left( -\frac{2}{5} \lg \frac{2}{5} - \frac{3}{5} \lg \frac{3}{5} \right)}_{a_2 = T} + \underbrace{\frac{4}{9} \times \left( -\frac{1}{2} \lg \frac{1}{2} - \frac{1}{2} \lg \frac{1}{2} \right)}_{a_2 = F}$

c) $\Delta_{info} = Z_{parent} - Z_{children}$

سس $Z_{children}$ را حساب کنیم. اینجا ویژگی $a_3$ را سورت را کرده ای حق:



$X < 4.5$ , $X \geqslant 4.5$ میانگین نقطه . محاسبه $\Delta_{info}$

برای هرکت از مقادیر قرمز بالا:

$$\Delta_{info} = Z_p - Z_c = \underbrace{-\frac{4}{9} \lg \frac{4}{9} - \frac{5}{9} \lg \frac{5}{9}}_{Z_p} - \underbrace{\left( \frac{6}{9} \times Z_{\geqslant 4.5} + \frac{3}{9} \times Z_{<4.5} \right)}_{Z_c}$$

کلا $Z_{<4.5} = -\frac{2}{3} \lg \frac{2}{3} - \frac{1}{3} \lg \frac{1}{3}$ و $Z_{\geqslant 4.5} = -\frac{3}{6} \lg \frac{3}{6} - \frac{3}{6} \lg \frac{3}{6} = -\lg \frac{1}{2}$

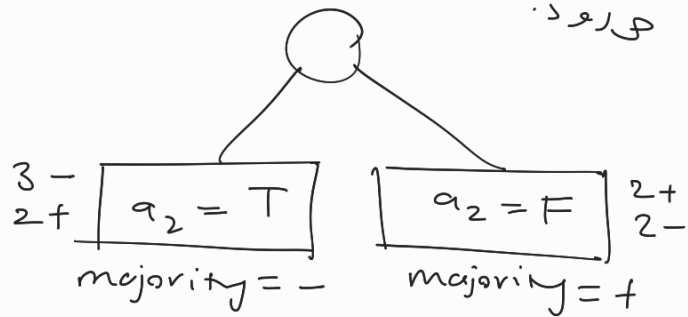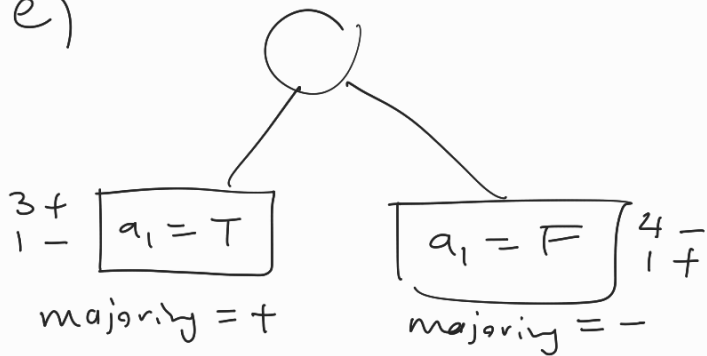به طور مشابه برای مابقی موارد قرمز محاسبه می شود.

d) برای $a_3$ به این علت که split ها طبیعی و سودمند نیست و فقط ۱

دو برای $a_1$:

$$\Delta_{M.}(a_2) = I_{parent} - \frac{5}{9} \times \left(-\frac{2}{5}\lg\frac{2}{5} - \frac{3}{5}\lg\frac{3}{5}\right) - \frac{4}{9} \times \left(-\frac{1}{2}\lg\frac{1}{2} - \frac{1}{2}\lg\frac{1}{2}\right)$$

$$\uparrow$$
$$a_2 = T \qquad\qquad\qquad\qquad a_2 = F$$

$$\Delta(a_1) = I_{parent} - \frac{4}{9} \times \left(-\frac{3}{4}\lg\frac{3}{4} - \frac{1}{4}\lg\frac{1}{4}\right) - \frac{5}{9} \times \left(-\frac{1}{5}\lg\frac{1}{5} - \frac{4}{5}\lg\frac{4}{5}\right)$$

$$\uparrow \qquad\qquad\qquad\qquad\qquad \nearrow$$
$$a_1 = T \qquad\qquad\qquad\qquad a_1 = F$$

هم کدام را انتخاب کنیم؟ ببین $\Delta(a_1)$، $\Delta(a_2)$، برای split $a_3$ که به این نمی‌کنیم، بزرگ‌ترین، بزرگ‌ترین انتخاب می‌شود و به طور کلی برای splitting به این می‌آوریم.

e)

3 t
1 − [ $a_1 = T$ ] [ $a_1 = F$ ] 4 −
                                    1 t

majority $= t$     majority $= -$

3 −
2 f [ $a_2 = T$ ] [ $a_2 = F$ ] 2 +
                                      2 −

majority $= -$     majority $= t$

$$\text{miss} = \frac{f_{+-} + f_{-+}}{9} = \frac{2}{9}$$
$$\downarrow$$
on traing set

$$\frac{2}{9} < \quad \text{miss} = \frac{2 + 2}{9} = \frac{4}{9}$$

Hence $a_1$ is best split.

f) $\Delta(a_1) = I_p -$

$$\Delta(a_2) = I_p - \frac{5}{9} \times \left(1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2\right) - \frac{4}{9} \times \left(1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2\right)$$

The bigger is the best split because we would obtain much purer children nodes.

a) optimistic approach → $err_{gen} = err_{train}$

pessimistic approach → $err_{gen} = err_{train} + 0.5 \times \dfrac{K}{N_{train}}$

$err_{train} = \dfrac{\overset{A=0,B=1}{1} + \overset{A=1,C=0}{3} + \overset{A=1,C=1}{1}}{10} = \dfrac{1}{2}$

$K = 4$, $N_{train} = 10$

Hence optimistic $= 0.5$ , pessimistic $= \dfrac{1}{2} + \dfrac{1}{2} \times \dfrac{4}{10}$

$= \dfrac{7}{10} = 0.7$

validation $= \overset{A=0,B=1}{\dfrac{1}{5}} = 0.2$

underfitting زمانی اتفاق می‌افتد که مدل الگوهای موجود در train set را یاد نمی‌گیرد و خطای بالایی در هر دو دیتاست هم روی train set و هم test set می‌دهد. در این موارد مدل بسیار ساده است و به اندازه کافی complex نیست که بتواند to learn کند. مثلا می‌تواند درختی باشد که عمق کمی داشته باشد.

to pattern نتواند to leaf سریع رسیده و نتواند درست کار کند.