(a) ID3

total = 14,    +: 9    −: 5

$S(9+,5-)$    $P(C_1) = 9/14$    $P(C_2) = 5/14$

Entropy $(S) = -\sum_{i}^{2} P(W_i) \log_2 P(W_i)$

$= -\left(\frac{9}{14} \times \log_2 \frac{9}{14} + \frac{5}{14} \times \log_2 \frac{5}{14}\right) \approx 0.94$

(I) root:

① outlook:    Sunny: $5 \frac{2+}{3-}$    overcase: $4+$    Rain: $5 \frac{3+}{2-}$

$E(S) = -(2/5 \times \log_2 2/5 + \frac{3}{5} \log_2 3/5) = 0.971$

$E(O) = -(1 \times \log 1) = 0$

$E(R) = -(3/5 \times \log_2 3/5 + 2/5 \times \log_2 2/5) = 0.971$

$G(S, outlook) = E(S) - \frac{5}{14} E(sunny) - \frac{5}{14} E(overcase) -$

$5/14 \times E(Rain) = 0.94 - \frac{5}{7} \times 0.971 = 0.246$

② temp:    hot: $4 \frac{2+}{2-}$    mid: $6 \frac{4+}{2-}$    cool: $4 \frac{3+}{1-}$

❶ $E(h) = -(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4}) = 1$

$E(m) = -(\frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \times \log_2 \frac{2}{6}) = 0.918$

$E(c) = -(\frac{3}{4} \times \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}) = 0.811$

$G(S, temp) = 0.84 - \frac{4}{14} \times 1 - \frac{6}{14} \times 0.918 - \frac{4}{14} \times 0.81$

$= 0.028$

② Humidity:    H: $7 \frac{3+}{4-}$    N: $7 \frac{6+}{1-}$

$E(H) = -(3/7 \log_2 3/7 + 4/7 \log_2 4/7) = 0.985$

$E(N) = -(6/7 \log_2 6/7 + 1/7 \log_2 1/7) = 0.592$

$G(S, H) = 0.94 - 7/14 \times 0.985 - 7/14 \times 0.592 = 0.151$

④ Wind : $W:8\frac{6+}{2-}$     $S:6\frac{3+}{3-}$
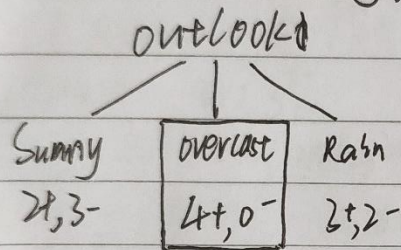
$E(w) = -(6/8 \log 6/8 + 2/8 \log 2/8) = 0.811$

$E(s) = -(3/6 \log 3/6 + 3/6 \log 3/6) = 1$

$G(s,w) = 0.94 - 8/14 \times 0.811 - 6/14 \times 1 = 0.048$

$G(s,0)$ is the biggest, outlook is root

outlook



| Sunny | overcast | Rain |
|-------|----------|------|
| $2+,3-$ | $4+,0-$ | $2+,2-$ |

(II) Next Node in sunny

$E(sunny) = 0.971$

① Temp: $H:2-$ , $M:2\frac{1+}{1-}$ , $C:1+$

$E(H) = 0$ , $E(C) = 0$ , $E(M) = 1$

$G(s, temp) = 0.971 - 2/5 \times 1 = 0.571$
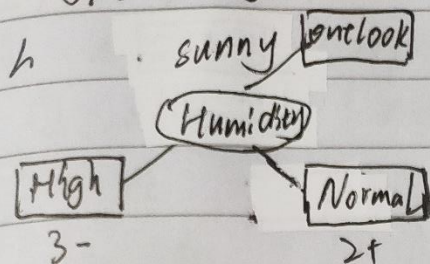
② Humidity : $H:3-$ , $N:2+$

$E(H) = 0$ , $E(N) = 0$

$G(s,H) = 0.971 - 0 = 0.971$

③ Wind: $W:3\frac{1+}{2-}$ , $S:2\frac{1+}{1-}$

$E(w) = -(1/3 \log_2 1/3 + \frac{2}{3} \log_2 2/3) = 0.918$ , $E(s) = 1$

$G(s,w) = 0.971 - 3/5 \times 0.918 - 2/5 \times 1 = 0.0202$

$\therefore$ G(sunny, Humidity) is the larggest,



```
↳      . sunny [outlook]
           (Humidity)
    [High]          [Normal]
     3-                2+
```

(III), next node in Rain

$E(Rain) = 0.971$

① temp = M: $3^{2+}_{1-}$   L: $2^{1+}_{1-}$

$E(mild) = -(2/3 \log_2 2/3 + 1/3 \log_2 1/3) = 0.918$, $E(L) = 1$

$G = 0.971 - \frac{2}{5} - \frac{3}{5} \times 0.918 = 0.02D2$
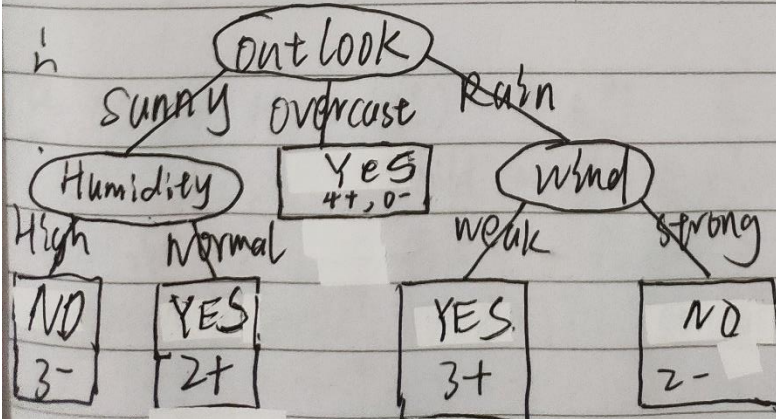
② Humidity = H$2^{1+}_{1-}$  N:$3^{2+}_{1-}$

$E(H) = 1$,   $E(W) = 0.918$

$G = 0.971 - 2/5 - 3/5 \times 0.918 = 0.0202$

③ wind = W:3+,  S:2-

$E(W) = 0$,   $E(S) = 0$,   $G = 0.971$, (biggest)

```
i
↳        (outlook)
   sunny   overcast   Rain
 (Humidity)   [Yes]     (wind)
High  Normal   4+,0-   weak   strong
[NO]  [YES]        [YES]    [NO]
 3-    2+           3+       2-
```

(b) CART.

(I) root, D as the sample

A. outlook     sunny: $5\frac{2+}{2-}$  overcast: $4+$  Rain: $5\frac{3+}{2-}$

① {sunny} | {overcast, Rain}

A₁ (2+,3-)         A₂ (7+,2-)

$Gini(A_1) = 1 - (\frac{2}{5})^2 - (\frac{3}{5})^2 = \frac{12}{25}$

$Gini(A_2) = 1 - (\frac{7}{9})^2 - (\frac{2}{9})^2 = \frac{28}{81}$

$Gini(D, \{A_1, A_2\}) = \frac{5}{14} \times Gini(A_1) + \frac{9}{14} \times Gini(A_2)$

$= \frac{5}{14} \times \frac{12}{25} + \frac{9}{14} \times \frac{28}{81} = 0.394$

② {overcast} | {sunny, Rain},  A₃(4+) | A₄(5+5-)

$G(A_3) = 0$,      $G(A_4) = 1 - 0.5^2 - 0.5^2 = 0.5$

$G(D, \{A_3, A_4\}) = \frac{1}{2} \times \frac{10}{14} = 0.357$

③ {Rain} | {sunny, overcast},  A₅|A₆

$G(A_5) = 1 - (\frac{3}{5})^2 - (\frac{2}{5})^2 = 12/25$

$G(A_6) = 1 - (\frac{6}{8})^2 - (\frac{3}{9})^2 = \frac{4}{9}$

$G(D, \{A_5, A_6\}) = 5/14 \times \frac{12}{25} + \frac{9}{14} \times \frac{4}{9} = \frac{16}{35} = 0.457$


B. temp: H = 4 $(\frac{2+}{2-})$  M = 6 $\frac{4+}{2-}$  C = 4 $\frac{3+}{1-}$

① $Gini(D, \{Hot\}|\{mild, cool\}) = \frac{4}{14} \times (1 - \frac{1}{2}^2 - \frac{1}{2}^2) + \frac{10}{14} \times$

$(1 - \frac{7^2}{10} - \frac{3^2}{10}) = \frac{4}{14} \times \frac{1}{2} + \frac{10}{14} \times \frac{21}{50} = 0.443$

② $Gini(D, \{mild\}|\{Hot, cool\}) = (1 - (\frac{4}{6})^2 - (\frac{2}{6})^2) \times \frac{6}{14} + \frac{8}{14} \times (1 -$

$(\frac{5}{8})^2 - (\frac{3}{8})^2) = \frac{6}{14} \times \frac{4}{9} + \frac{8}{14} \times \frac{15}{32} = \frac{11}{24} = 0.458$

③ $Gini(D, \{cool\}|\{Hot, mild\}) = \frac{4}{14} \times (1 - (\frac{3}{4})^2 - (\frac{1}{4})^2) + \frac{10}{14} \times (1 - (\frac{6}{10})^2 - (\frac{4}{10})^2)$

$= \frac{4}{14} \times \frac{3}{8} + \frac{10}{14} \times \frac{12}{25} = \frac{9}{20} = 0.45$

C. Humidity : $H:7(\frac{3+}{4-})$  $N:7(\frac{6+}{1-})$
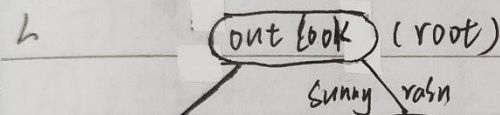
$Gi(D, \{Gini(H), Gini(N)\}) = \frac{1}{2} \times (1-(\frac{3}{7})^2-(\frac{4}{7})^2) + \frac{1}{2} \times$

$(1-(\frac{6}{7})^2-(\frac{1}{7})^2) = \circ \frac{24+12}{49\times2} = \frac{18}{49} = 0.367$

D. Wind : $W:8\frac{6+}{2-}$  $S:6\frac{3+}{3-}$

$Gini(D, \{Gini(W), Gini(S)\}) = \frac{8}{14}(1-(\frac{6}{8})^2-(\frac{2}{8})^2) + \frac{6}{14}(1-(\frac{3}{6})^2-(\frac{3}{6})^2)$

$= \frac{8}{14}\times \frac{3}{8} + \frac{6}{14}\times \frac{1}{2} = 0.2429$

∵ $Gini(D; \{overcase\}|\{Sunny; Rain\})$ has the smallest

∟                                    $Gini : 0.357$

(outlook) (root)

Sunny / rain

$4^+$ [overcast] Yes (    )    $D = D\{1,2, 4,5, 6, 8,9,10, 13,14\}$

(ZZ) second node

Ⓐ temp: Hot:$2-$, Mid: $8\frac{3+}{2-}$, cool : $3\frac{2+}{1-}$

① $Gini(D, \{G(H)|G(mild,coo)\}) = \frac{2}{10}\times 6 + \frac{8}{10}\times\frac{15}{32} = 0.375$

② $Gini(D, \{G(m)|G(c,H)\}) = \frac{5}{10}\times(1-\frac{12}{25}) + \frac{5}{10}\times\frac{12}{25} = 0.48$

③ $Gini(D, \{G(c)|G(m,H)\}) = \frac{3}{10}\times\frac{4}{9} + \frac{7}{10}\times\frac{24}{49} = 0.476$

Ⓑ Wind  $W:6\frac{4+}{2-}$  $S:4\frac{1+}{3-}$

$G(D,\{w,s\}) = \frac{6}{10}\times\frac{16}{36} + \frac{4}{10}\times\frac{6}{16} = 0.417$

Ⓒ Humidity  $H:5\frac{1+}{4-}$  $N:5\frac{4+}{1-}$

$G(D, \{H,N\}) = \frac{1}{2}\frac{8}{25} + \frac{1}{2}\frac{8}{25} = \frac{8}{25} = 0.32$

Ⓓ Outlook  sunny $5\frac{2+}{3-}$  Rain $5\frac{3+}{2-}$

$G(D, \{S,R\}) = \frac{1}{2}\frac{12}{25}\times 2 = 0.48$

the smallest is $Gini(D, \{High, Normal\}) = 0.32$



the left, $D = D_2$

(III) the left, $D = D_2$

(A) outlook: Sunny: $3^-$, Rain: $2\,^{1+}_{1-}$

$Gini(D, \{S, R\}) = \frac{3}{5} \times 0 + \frac{2}{5} \times \frac{1}{2} = 0.2$
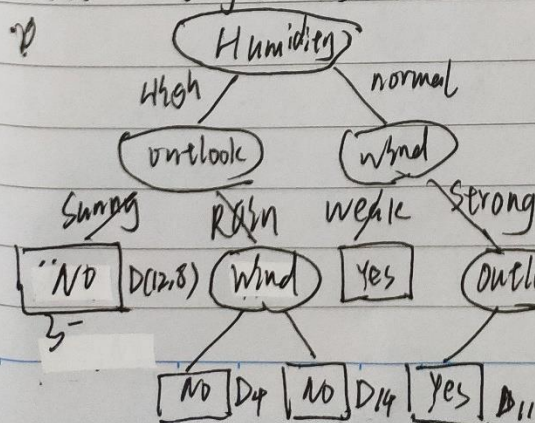
(B) temp: $H: 2^-$  mild: $3\,^{1+}_{2-}$

$Gini(D, \{H, m\}) = \frac{2}{5} \times 0 + \frac{3}{5} \times (1 - (\frac{1}{3})^2 - (\frac{2}{3})^2) = \frac{4}{15} = 0.267$

(C) wind: $W: 3\,^{1+}_{2-}$   $S: 2^-$
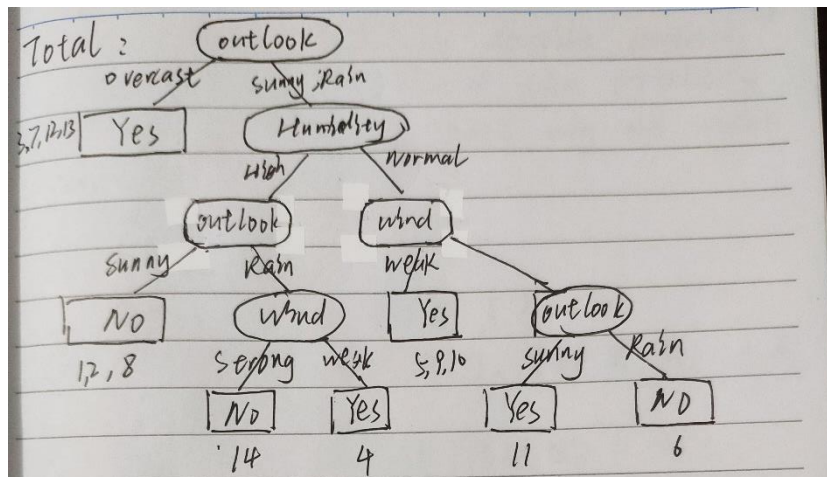
$Gini(D, \{W, S\}) = \frac{3}{5} \times \frac{4}{8} + 0 = 0.267$

$\therefore Gini(D, \{S, R\}) = 0.2$ is the smallest

$\therefore$ Left is outlook, right is wind

RND, Because the last two node have only 1+ sample and 1- sample, which creat 4 nodes.

Total :



(c) ① CART use minimization of squared error, ID3 use log, when calculate, CART is faster than ID3

② CART tree is a binary tree, ID3 is not. The branching factor of ID3 equals to the number of decision category

③ CART has a deeper tree than ID3, so when decision, ID3 is faster.

④ ID3 use information gain to describe choosing, so choose the biggest one which can make the remaining samples have more purity

CART use Gini index to describe decision, the goal is to minimize the probability of miss classification. That's why we choose the smallest one.

HW4 附加题页
K 聚类如何选初始点?
① 随机选取, 多次平均。重效果一般, 也是最原始的选择法
② K-means ++. 选尽可能远的 K 个点。对于第 n (1<n≤k) 个点, 选离前 n-1 个点的中心最远的点
③ ISODATA: 当某一类点过少, 去除之。当某一类点过多、分散较大, 分裂的两个聚类