

ID3 范例：

题目：

Target : 属性值分布: $\{T: 9, N: 5\}$ Appearance : $\{A_h: 5 = 3T + 2N, Good: 5 = 2T + 3N,$

Great: 4 = 4T}

Attributed Income : $\{Low: 4 = 2T + 2N, Good: 6 = 4T + 2N,$

Great: 4 = 3T + 1N}

Age : $\{Younger: 7 = 6T + 1N, Older: 7 = 3T + 4N\}$ Profession : $\{Unstable: 6 = 3T + 3N, Steady: 8 = 6T + 2N\}$

第一步：系统熵计算

$$H(D) = - \sum_{k=1}^K P_k \log_2 P_k = - \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

第二步：计算每个特征条件下的熵

$$\text{Appearance : } H(F_{Great}) = - \frac{4}{4} \log_2 \left(\frac{4}{4} \right) = 0$$

$$H(F_{Good}) = - \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$H(F_{Age}) = - \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

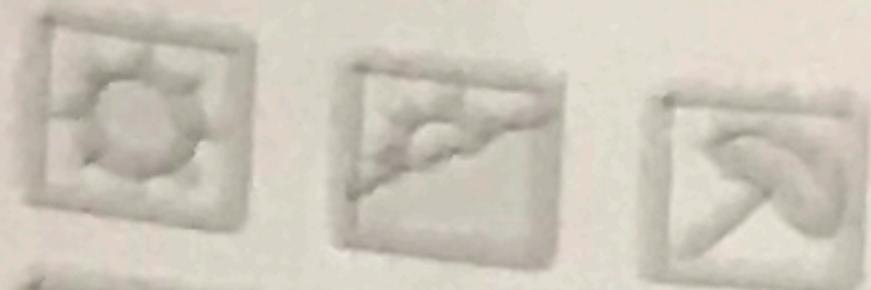
$$H(D|F_{App}) = \frac{4}{14} H(F_{Great}) + \frac{5}{14} H(F_{Good}) + \frac{5}{14} H(F_{Age}) \\ = 0.693$$

$$\text{Income : } H(F_{Great}) = - \frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$$

$$H(F_{Good}) = - \frac{6}{10} \log_2 \frac{6}{10} - \frac{2}{10} \log_2 \frac{2}{10} = 0.918$$

$$H(F_{Low}) = - \frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 0.1$$

$$H(P|F_{Inc}) = \frac{4}{14} H(F_{Great}) + \frac{6}{14} H(F_{Good}) + \frac{4}{14} H(F_{Low}) \\ = 0.911$$



Mo Tu We Th Fr Sa Su

Memo No. _____

Date / /

$$\text{Age: } H(F_{\text{Age}}) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.985$$

$$H(F_{\text{Younger}}) = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.592$$

$$H(D|F_{\text{Age}}) = \frac{3}{7} H(F_{\text{Older}}) + \frac{4}{7} H(F_{\text{Younger}}) = 0 \\ = \cancel{0.789}$$

$$\text{Profession: } H(F_{\text{steady}}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.84$$

$$H(V_{\text{Unstable}}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

$$H(D|F_{\text{pro}}) = \frac{6}{7} H(\cancel{F_{\text{Funstable}}}) + \frac{1}{7} H(F_{\text{steady}}) \\ = 0.892$$

第三步：计算不同特征值的熵增：

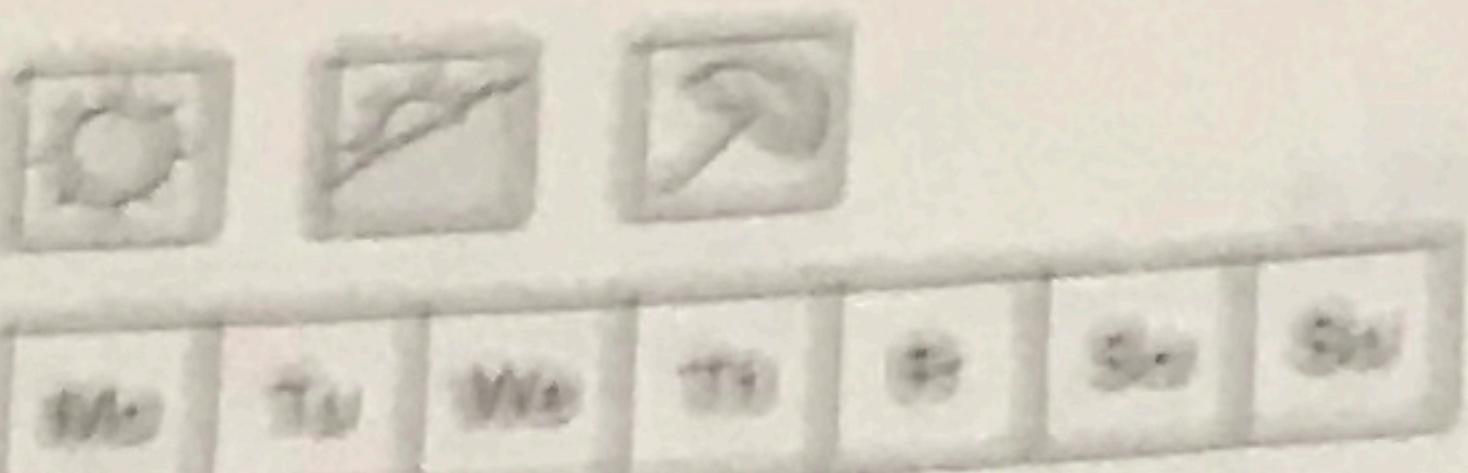
$$G(D|F_{\text{app}}) = H(D) - H(D|F_{\text{app}}) = 0.246$$

$$G(D|F_{\text{Inc}}) = H(D) - H(D|F_{\text{Inc}}) = 0.029$$

$$G(D|F_{\text{age}}) = H(D) - H(D|F_{\text{age}}) = 0.151$$

$$G(D|F_{\text{pro}}) = H(D) - H(D|F_{\text{pro}}) = 0.048$$

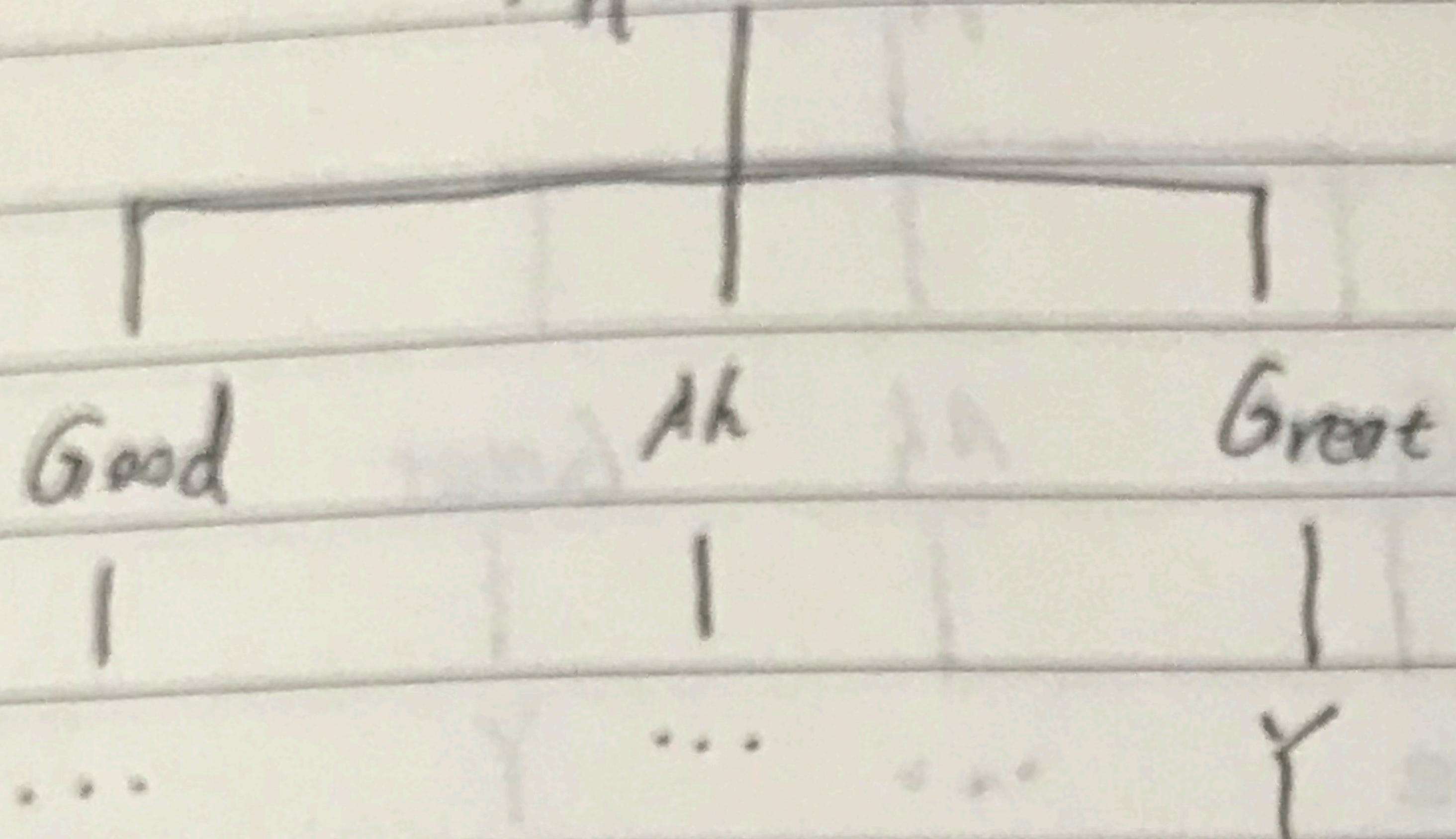
显然，上式中 Appearance 对应的熵增最大，故选 Appearance 作为决策树的根节点，又 Appearance 分为了 3 种情况，即含有 3 个分支，其中 Appearance 为 Great 的熵增为 0，故不再对 Appearance 的 Great 分支进行划分。



Memo No. _____
Date / /

截止目前决策树图为

Appearance



接下来观察 Appearance 的 Good 分支，包含情况如下：

~~Appearance - Good~~: $\{ \text{Low: } 2 = 2N, \text{Good: } 2 = 1Y + 1N, \text{Great: } 1 = 1Y \}$

$\{ \text{Income: } \{ \text{Low: } 2 = 2N, \text{Good: } 2 = 1Y + 1N, \text{Great: } 1 = 1Y \} \}$

$\{ \text{Age: } \{ \text{Younger: } 2 = 2Y, \text{Older: } 3 = 3N \} \}$

$\{ \text{Profession: } \{ \text{Unstable: } 2 = 1Y + 1N, \text{Steady: } 3 = 1Y + 2N \} \}$

~~APP-Good~~

此时，系统熵即为之前的 $H(F_{APP}) = 0.971$

类比之前的计算方法，可得：

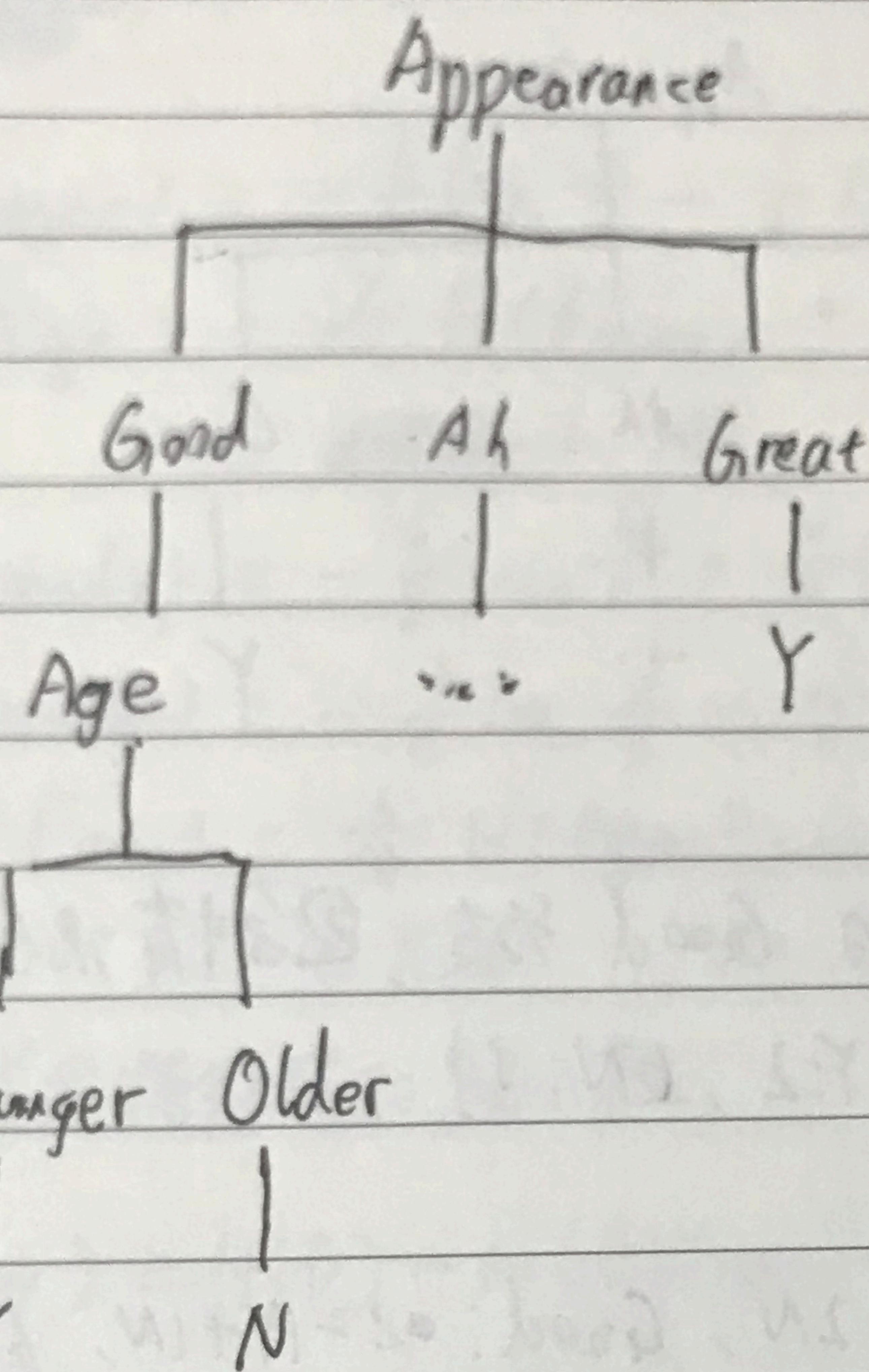
$$G(F_{APP} | F_{Inc}) = H(F_{APP}) - H(F_{APP} | F_{Inc}) = 0.571$$

$$G(F_{APP} | F_{Age}) = H(F_{APP}) - H(F_{APP} | F_{Age}) = 0.971$$

$$G(F_{APP} | F_{Pro}) = H(F_{APP}) - H(F_{APP} | F_{Pro}) = 0.020$$

故 Appearance - Good 的分支选取 Age 作为根节点，其下属 2 支分
支的 K 值均 > 0，故不再划分，~~再划~~

此时决策树对更新如下：



同理对 ~~Age~~ Appearance 和 ~~Age~~ 以及 ~~Profession~~ 进行计算，可知 ~~Profession~~ 的熵增最大，故将 ~~Profession~~ 作为 ~~Appearance - Ah~~ 的根节点，~~因~~ 此时 ~~Profession~~ 的两分支熵均为 0，故不再划分。最终整理决策树如下：

