? For a Sunburn dataset given below, find the first splitting attribute for the decision tree by using the ID3 algorithm.

| Name | Hair | Height | Weight | Lotion | Class |
|------|------|--------|--------|--------|-------|
| Sarah | Blonde | Average | Light | No | Sunburn |
| Dana | Blonde | Tall | Average | Yes | None |
| Alex | Brown | Tall | Average | Yes | None |
| Annie | Blonde | Short | Average | No | Sunburn |
| Emily | Red | Average | Heavy | No | Sunburn |
| Pete | Brown | Tall | Heavy | No | None |
| john | Brown | Average | Heavy | No | None |
| katie | Blonde | Short | light | Yes | None |

Step 1 :- Calculate the entropy of data set

$T$ = data set

$$Info(T) = Entropy(p) = \sum_{j} \log_2 \frac{}{}$$

$$- P_{sunburn} \log_2 P_{sunburn}$$

$$- P_{none} \log_2 P_{none}$$

$$\text{Entropy}\left[\frac{3}{8}, \frac{5}{8}\right] = -\sum_{i=1}^{n} P_i \log_2 P_i$$

$$= -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8}$$

$$= 0.954$$

Step 2 :- For each attribute Hair, Height, weight, Lotion find entropy for all categorical values and also then fend the information gain for the features.

First attribute :- Hair

Categorical values are :- Blonde, Brown, Red

$$\text{Info (Hair, T)} = \sum_{i=1}^{3} \frac{|T_i|}{|T|} \text{info}(T_i)$$

$$= \frac{4}{8} \text{info (Blonde)} + \frac{3}{8} \text{info (Brown)}$$

$$\quad + \frac{1}{8} \text{info (Red)}$$

$$= \frac{4}{8}\left[-\frac{2}{5}\log_2 \frac{2}{5} - \right.$$

$$= \frac{4}{8}\left[-\frac{2}{4}\log_2 \frac{2}{4} - \frac{2}{4}\log_2 \frac{2}{4}\right] +$$

$$\frac{3}{8}\left[-\frac{2}{3}\log_2 \frac{2}{3} - \frac{3}{3}\log_2 \frac{3}{3}\right]$$

$$\frac{1}{8}\left[-\frac{1}{1}\log_2\frac{1}{1} - 0\right]$$

$$= 0.5 + 0 + 0$$

$$= \underline{0.5}$$

Gain (Hair, T) = info(T) - info (Hair)

$$= 0.954 - 0.5$$

$$= \underline{0.454}$$

Second attribute = Height

Categorical values are = Average, tall, short

info ( Height ,T) $= \frac{3}{8}$ info (Average) $+ \frac{3}{8}$ info (Tall)

$+ \frac{2}{8}$ info (short)

$$= \frac{3}{8}\left[-\frac{2}{3}\log_2\frac{2}{3} - \frac{1}{3}\log_2\frac{1}{3}\right] +$$

$$\frac{3}{8}\left[-\frac{0}{3}\log_2\frac{0}{3} - \frac{3}{3}\log_2\frac{3}{3}\right] +$$

$$\frac{2}{8}\left[-\frac{1}{2}\log_2\frac{1}{2} - \frac{1}{2}\log_2\frac{1}{2}\right]$$

$$= 0.344 + 0 + 0.25$$

$$= 0.594$$

$$\text{Gain (Height, T)} = \text{info (T)} - \text{info (Height)}$$

$$= 0.954 - 0.594$$

$$= 0.36$$

Third attribute = weight

Categorical values are = Light, Average, Heavy

$$\text{Info (weight, T)} = \frac{2}{8} \text{ info (Light)} + \frac{3}{8} \text{ info (Average)} + \frac{3}{8} \text{ info (Heavy)}$$

$$= \frac{2}{8} \left[ -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] +$$

$$\frac{3}{8} \left[ -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right] +$$

$$\frac{3}{8} \left[ -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right]$$

$$= 0.25 + 0.344 + 0.344$$

$$= 0.938$$

$$\text{Gain (weight, T)} = \text{info (T)} - \text{info (weight)}$$

$$= 0.954 - 0.938$$

$$= 0.016$$

Fourth attribute = Lotion

Categorical values are = No, Yes

$$\text{info}(\text{Lotion}, T) = \frac{5}{8} \text{info}(No) + \frac{3}{8} \text{info}(Yes)$$

$$= \frac{5}{8}\left[-\frac{3}{5}\log_2\frac{3}{5} - \frac{2}{5}\log_2\frac{2}{5}\right]$$

$$+ \frac{3}{8}\left[-\frac{0}{3}\log_2\frac{0}{3} - \frac{3}{3}\log_2\frac{3}{3}\right]$$

$$= 0.606 \quad 0.607 + 0$$

$$= 0.607$$

$$\text{Gain}(\text{Lotion}, T) = \text{info}(T) - \text{info}(\text{Lotion})$$

$$= 0.954 - 0.607$$

$$= 0.347$$

step 3 :- Here the attribute with maximum inform^n gain is Hair. So Hair is the root of decision tree.

| attribute | inform$^n$ gain |
|-----------|-----------------|
| Hair      | 0.454           |
| Height    | 0.36            |
| weight    | 0.016           |
| Lotion    | 0.347           |