## Option 1

```
        ┌──────────────┐
        │ 20 C | 10 D  │
        │   Gender     │
        └──────┬───────┘
          M ╱      ╲ F
    ┌──────────┐   ┌──────────┐
    │ 19 C| 1 D│   │ 1 C | 9 D│
    └──────────┘   └──────────┘
```

Lesser entropy

Option 1

## Option 2

```
        ┌──────────────┐
        │ 20 C | 10 D  │
        │   Height     │
        └──────┬───────┘
      <5'7" ╱      ╲ ≥ 5'7"
    ┌──────────┐   ┌──────────┐
    │ 10 C| 5 D│   │ 10 C| 5 D│
    └──────────┘   └──────────┘
```

More entropy

Option 2

Decision Trees are concerned with increase overall homogenity.

**Entropy** → way to measure non-uniformness of collection of items of various types

$$En = -\sum p_i \log p_i$$

1.
→ is more mixed; hence more entropy

$p_B = \frac{4}{7}$  $p_B = \frac{3}{7}$

$$En = -\left[\frac{4}{7}\log\frac{4}{7} + \frac{3}{7}\log\frac{3}{7}\right]$$

$\underbrace{\phantom{xxx}}_{0.1388}$  $\underbrace{\phantom{xxx}}_{0.1577}$

$En = 0.2965.$

2.
→ is less mixed; hence lesser entropy

$p_B = \frac{6}{7}$  $p_B = \frac{1}{7}$

$$En = -\left[\frac{6}{7}\log\frac{6}{7} + \frac{1}{7}\log\frac{1}{7}\right] = 0.1781$$

3.
→ is not mixed; minimum possible entropy

$p_B = 1, p_B = 0$

$$En = -\left[1\log 1 + 0\log 0\right] = 0$$

$$X$$

$$\underbrace{x_1 \quad x_2 \quad \cdots \quad x_n}_{\text{classification}} \qquad y$$

- discrete features
- discrete response
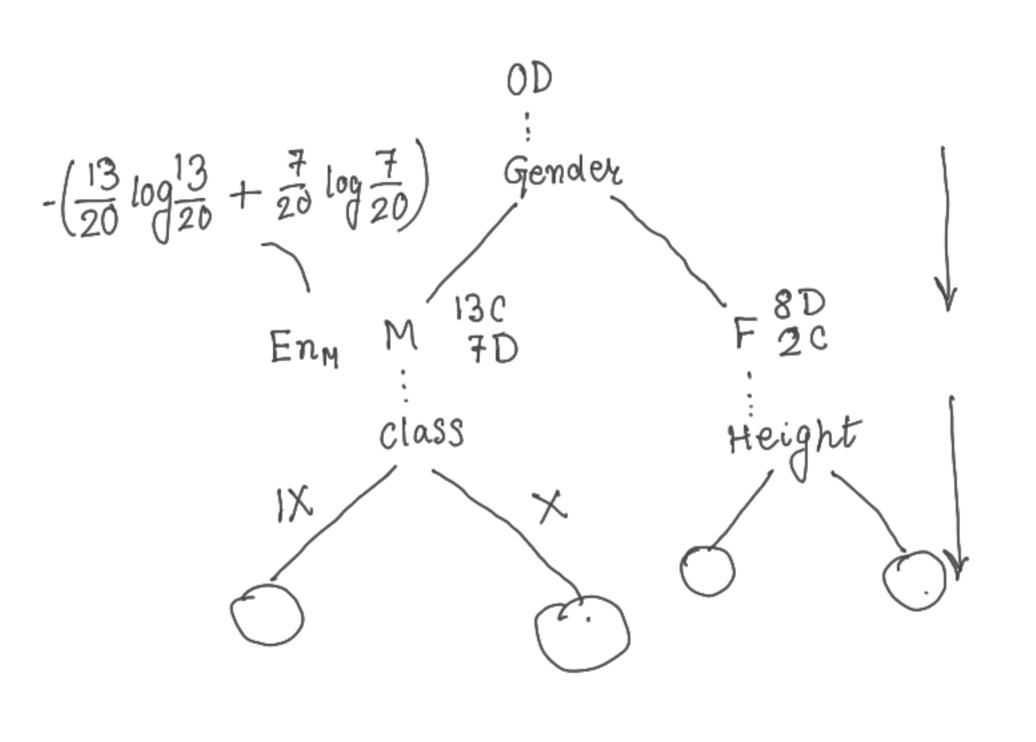
$$30 \begin{cases} 15 & 1/2 \\ 15 & 1/2 \end{cases}$$

1. Entropy of the dataset $-\left[\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}\right] = 1$

2. Identify the feature that leads to a $\underbrace{\text{maximum decrease in entropy}}_{\text{information gain.}}$

3. Using the most useful feature as parent node; we find another feature that leads to a max. decrease in entropy

4. Do (3) until all features are exhausted or you hit a stopping criteria.

Original DS :     15 C , 15 D

$\overline{En} = 1$

Gender :   8D, 2C | 13C, 7D

$$En_g = 10 \times E_1 + 20 \times E_2 / (10 + 20)$$

$$-\left( \frac{8}{10} \log \frac{8}{10} + \frac{2}{10} \log \frac{2}{10} \right) \xrightarrow{=} \left( \frac{13}{20} \log \frac{13}{20} + \frac{7}{20} \log \frac{7}{20} \right)$$

$$1G = En - En_{g/h/c}$$

$\underbrace{\qquad}_{E_1}$        $\underbrace{\qquad}_{E_2}$

Height :    5C, 7D | 10C, 8

$$En_h = \frac{12 \times E_3 + 18 \times E_4}{(12 + 18)}$$

$\underbrace{\qquad}_{E_3}$    $\underbrace{\qquad}_{E_4}$

class :

$En_c$

$$-\left(\frac{13}{20}\log\frac{13}{20} + \frac{7}{20}\log\frac{7}{20}\right)$$

OD

Gender

$En_M$   M   13C 7D

F   8D 2C

class

1X     X

Height

$avg\left(En_H\right)$

| H | sort(H) | Splits |
|-----|---------|--------|
| 165 | 163 | 164 |
| 172 | 165 | 165.5 |
| 179 | 166 | 166.5 |
| 181 | 167 | 168 |
| 171 | 169 | |

$H < 164 \rightarrow En_{H<164}$

$H < 165.5 \rightarrow En_{H<165.5}$

$H < 166.5$

# CART

— it uses Gini index to measure impurity / homogenity .

$$G = 1 - \sum p_i^2$$

$$G = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = \frac{24}{49} \sim 0.5$$

$$G = 1 - (1)^2 = 0$$

Algorithm uses GI instead of entropy; rest similar to ID3.

$x_1 \; x_2 \ldots x_p \quad y \to$ cont

$\begin{aligned} & y_1 \\ & y_2 - \\ & \quad : \\ & y_n \end{aligned}$

DT $\to$ minimizes MSE

O D` : $MSE = \frac{1}{n} \sum (\bar{y} - y_i)^2$

$\quad\quad\quad\quad n_1 \quad\quad n_2$

$x_1 : \quad ( D_1 , \quad D_2 )$

$\overline{y_1} \quad\quad \overline{y_2} \to \frac{1}{n_2} \sum (y_i - \overline{y_2})^2$

$MSE_{x_1} : \quad \underline{\frac{1}{n_1} \sum (y_i - \overline{y_1})^2} + \frac{1}{n_2} \sum (y_i - \overline{y_2})^2$

$n_1 + n_2 = n.$

$MSE_{x_2} :$

$100$

$70$

$\overline{y_1}$     $\dfrac{30}{\overline{y_2}}$

$\rightarrow \overline{y_1}$

$1,2,3$    $2$       $100,101,102$    $101$



$\overline{x_G}=2$    $\overline{x}=51$    $\overline{x_B}=101$

1. $\sum (x_i - \overline{x})^2 = (1-51)^2 + (2-51)^2 + (3-51)^2$
$$+ (100-51)^2 + (101-51)^2$$
$$+ (102-51)^2$$

$\dfrac{(1-2)^2 + (2-2)^2 + (3-2)^2}{}$

$= 2. \begin{cases} \sum (x_{gi} - \overline{x_G})^2 + \approx 15000. = MSE \\ \sum (x_{bi} - \overline{x_B})^2 \end{cases}$

$(100-101)^2 +$
$(101-101)^2 +$
$(102-101)^2$

$= 4$ new MSE

$14496 \quad \triangle MSE$