

## Attribute selection technique

data set

| Age           | Income | student | CR        | class<br>Buy computer |
|---------------|--------|---------|-----------|-----------------------|
| youth         | High   | no      | Fair      | N                     |
| y...          | H...   | no      | Excellent | N                     |
| middle<br>age | H...   | no      | F         | Y                     |
|               | medium | no      | F         | Y                     |
| Sen           | low    | yes     | F         | Y                     |
| Sen           | low    | yes     | F         | N                     |
| m...          | low    | yes     | F         | N                     |
| y...          | m...   | yes     | E         | Y                     |
| y...          | low    | no      | F         | N                     |
| s...          | m...   | yes     | E         | Y                     |
| y...          | m...   | yes     | E         | Y                     |
| m             | m...   | no      | E         | Y                     |
| m             | H...   | yes     | F         | Y                     |
| m             | m...   | no      | E         | N                     |
| S             | Low    | yes     | E         | Y                     |

1) Information gain / Entropy measure

$$\text{Info}(\mathcal{D}) = - \sum_{i=1}^m P_i \log_2(P_i)$$

$P_i$  - probability that arbitrary tuple in ' $\mathcal{D}$ '

$$\text{Info}_A(\mathcal{D}) = \sum_{j=1}^J \frac{|\mathcal{D}_j|}{|\mathcal{D}|} \text{info}(\mathcal{D}_j)$$

$$\text{Gain} = \text{info}(\mathcal{D}) - \text{info}_A(\mathcal{D})$$

# of "no" in class = 5  $\hookrightarrow$  5 + 9 = 14

# of 'yes' in class = 9

$$\begin{aligned}\text{Info(D)} &= \frac{-9}{14} \log_2(9/14) - (5/14) \log_2(5/14) \\ &= -(-0.4098) - (-0.5305)\end{aligned}$$

$$\text{Info(D)} = 0.9403$$

say now we select "Age" as the factor to decide

$$\begin{aligned}\text{Info(D)}_{\text{Age}} &= \frac{5}{14} \left[ -\frac{2}{5} \log_2(2/5) - \frac{3}{5} \log_2(3/5) \right] + \\ &\quad \frac{4}{14} \left[ -\frac{4}{4} \log_2(4/4) - \frac{0}{4} \log_2(0/4) \right] + \\ &\quad \frac{5}{14} \left[ -\frac{3}{5} \log_2(3/5) - \frac{2}{5} \log_2(2/5) \right] \\ &= 0.694\end{aligned}$$

$$\begin{aligned}\text{gain}_{\text{(age)}} &= 0.9403 - 0.694 \\ &= 0.2463\end{aligned}$$

by

$$\text{gain}(\text{income}) = 0.029$$

$$\text{gain}(\text{student}) = 0.151$$

$$\text{gain}(\text{credit rating}) = 0.048$$

$$0.2463_{\text{gain}} > 0.151_{\text{student}} > 0.048_{\text{credit}} > 0.029_{\text{income}}$$