

Decision Trees and Naive Bayes

Feyyaz Demir

October 11, 2013

Decision Trees

- ▶ $H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$
- ▶ $H(X|Y = y) = - \sum_{x \in X} p(x|y) \log_2 p(x|y)$
- ▶ $H(X|Y) = - \sum_{y \in Y} p(y) \sum_{x \in X} p(x|y) \log_2 p(x|y)$
- ▶ At each node, we select the attribute Y that minimizes $H(X|Y)$ and split.
- ▶ We stop when entropy of a node becomes 0 (or less than some threshold).

Tennis Playing Data

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Decision Trees

Outlook	PlayTennis = No	PlayTennis = Yes	Total
Sunny	3	2	5
Rain	2	3	5
Overcast	0	4	4
Total	5	9	14

$$P(\text{Outlook} = \text{Sunny}) = \frac{5}{14}$$

$$P(\text{PlayTennis} = \text{Yes} | \text{Outlook} = \text{Sunny}) = \frac{2}{5}$$

$$H(\text{PlayTennis} | \text{Outlook} = \text{Sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$H(\text{PlayTennis} | \text{Outlook}) = -\frac{5}{14} \left(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right) - \frac{5}{14} \left(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) - \frac{4}{14} (0 \log_2 0 + 1 \log_2 1)$$

$$H(\text{PlayTennis} | \text{Outlook}) = 0.69$$

Decision Trees

$$H(\text{PlayTennis}|\text{Temperature}) = -\frac{4}{14}\left(\frac{1}{2}\log_2\frac{1}{2} + \frac{1}{2}\log_2\frac{1}{2}\right) - \frac{6}{14}\left(\frac{1}{3}\log_2\frac{1}{3} + \frac{1}{3}\log_2\frac{1}{3}\right) - \frac{4}{14}\left(\frac{1}{4}\log_2\frac{1}{4} + \frac{3}{4}\log_2\frac{3}{4}\right)$$

$$H(\text{PlayTennis}|\text{Temperature}) = 0.91$$

$$H(\text{PlayTennis}|\text{Humidity}) = -\frac{1}{2}\left(\frac{4}{7}\log_2\frac{4}{7} + \frac{3}{7}\log_2\frac{3}{7}\right) - \frac{1}{2}\left(\frac{1}{7}\log_2\frac{1}{7} + \frac{6}{7}\log_2\frac{6}{7}\right)$$

$$H(\text{PlayTennis}|\text{Humidity}) = 0.78$$

$$H(\text{PlayTennis}|\text{Wind}) = -\frac{8}{14}\left(\frac{2}{8}\log_2\frac{2}{8} + \frac{6}{8}\log_2\frac{6}{8}\right) - \frac{6}{14}\left(\frac{1}{2}\log_2\frac{1}{2} + \frac{1}{2}\log_2\frac{1}{2}\right)$$

$$H(\text{PlayTennis}|\text{Wind}) = 0.89$$

Outlook=Sunny

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Decision Trees

