

18 decision trees

问题: Choose which feature to split (according to importance)

Note: 选择 feature 的顺序不影响结果
但会影响 decision tree 的大小

↓
reduces uncertainty
the most

entropy: $H(x) = - \sum_k p(x=x_k) \log_2 p(x=x_k)$

x_k 代表 x 所有可能的取值

info gain: $B(q) = - (q \log_2 q + (1-q) \log_2 (1-q))$

实际计算:

- n - #negative samples
- p - #positive samples
- d - #outcomes of feature A

$$\textcircled{1} \text{Rem}(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right)$$

$$\textcircled{2} \text{Gain}(A) = B\left(\frac{p}{p+n}\right) - \text{Rem}(A)$$

pruning: start with a full tree

from top to down:

if a test is uninformative, replace the node with a leaf

↓
statistical significance testing

compact representation of decision process

explanatory (可解释的)

