

CS 4811

Feb 15, 2017

Wednesday (1)

week 6 now hw 2 due Friday search

week 7 M ✓
W career fair no class

week 8 M
W exam 1 4⁰⁰-5³⁰

Spring Break
↓ search ch 3
ch 4
ch 5
neural network

UPE resume clinic

Fri 4⁰⁰

career services

Google workshop

Previous class

Constructing decision trees

'learning'

classification

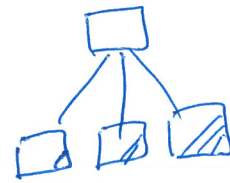
multi-valued concepts

↓
problem

data

→ concise representation

	A	B	C	D
				○
				○
				○
				○



↓ series of questions

minimum depth tree

look for the minimal tree that accurately represents this table

search: for this tree

optimization problem

(goal is not explicitly defined)

goal: to classify objects
"efficiently"

(2)

reduce uncertainty

13 examples ~~3~~ = false = 7 true = 6

uncertainty: $\frac{6}{13}, \frac{7}{13}$

One question: which one would reduce the uncertainty the most?

heuristic:

A is true

A is false

originally
probability of
error,
without asking
any questions

class
↓
classification
 $\left\{ \begin{array}{l} \text{true} = 6 \\ \text{false} = 2 \end{array} \right.$
 $\frac{2}{8}$

$\frac{8}{13} \frac{2}{8}$

true: 0
false: 5

$\frac{0}{5}$

$\frac{5}{13} \frac{0}{5}$

$\frac{6}{13}$

$\frac{6}{13}$

$\frac{2}{13}$

B: $\frac{5}{13}$

C: $\frac{4}{13}$

D: $\frac{3}{13}$

E: $\frac{6}{13}$

A

Entropy: measure of uncertainty

(3)

how many bits are needed to describe the possible outcomes

one value: 0 bits
2 values: 1 bit
4 values: 2 bits

Random variable V with values v_k each of which with probability $P(v_k)$
a variable that can take possible some values each with a probability

$$\text{entropy} = H(V) = \sum_k P(v_k) \log_2 \frac{1}{P(v_k)}$$

$$= - \sum_k P(v_k) \log_2 P(v_k)$$

Example fair coin flip

$$H(\text{fair coin}) = - \left[\underbrace{\frac{1}{2}}_{0.5} \underbrace{\log_2 \frac{1}{2}}_{-1} + \underbrace{\frac{1}{2}}_{-0.5} \underbrace{\log_2 \frac{1}{2}}_{-1} \right]$$
$$= 1$$

$$H(\text{loaded}) = - \left[\underbrace{(0.99 \log_2 0.99)}_{-0.01} + \underbrace{(0.01 \log_2 0.01)}_{-6.64} \right] \quad (4)$$

$$= 0.08$$

$B(q)$ = entropy of a Boolean variable that is true with probability q .

$$B(q) = - (q \log_2 q + (1-q) \log_2 (1-q))$$

p : # of positive examples 6 total = 13
 n : # of negative examples 7 all = $p+n$

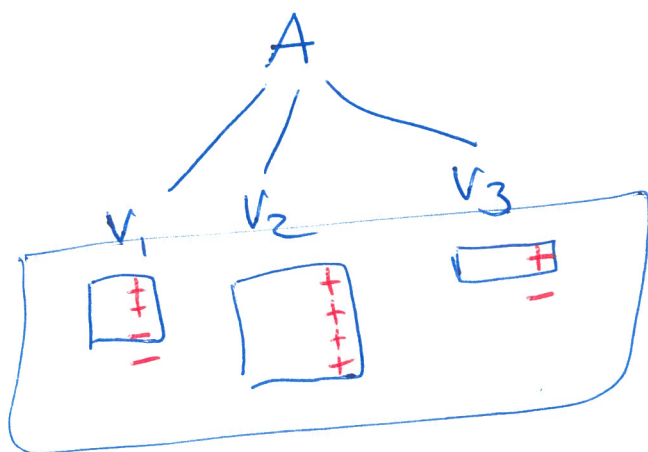
entropy before asking any questions
 $B\left(\frac{p}{p+n}\right)$ $B\left(\frac{6}{13}\right)$

After picking position A
 Remainder (A)

A has d answers
 for each: how much uncertainty do all I have?

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

p_k : # of positive examples when A 's value is v_k
 n_k : # of negative examples when A 's value is v_k



"expectation"

lottery	0.1	1000
	0.2	10
	0.7	0

expected value of this lottery is

$$\frac{0.1 \times 1000 + 0.2 \times 10 + 0.7 \times 0}{100} = 102$$

using the entropy formulas gain a reduction in uncertainty when we ask a question

$$\text{gain}(A) = B \left(\frac{P}{p+n} \right) - \text{remainder}(A)$$

The question (attribute) that has the highest gain will win a spot at the root of the tree.

$$\text{gain}(A) = B \left(\frac{6}{13} \right) - \left[\frac{2}{13} B \left(\frac{6}{8} \right) + \frac{5}{13} B \left(\frac{0}{5} \right) \right]$$

1.00
0.62
0.81
0

0.50

$$= 0.50 \text{ for } A$$

B, C, D, E will have lower values for gain

Summary

(7)

entropy = measure of uncertainty

random variable V with values v_k
each with probability $P(v_k)$

$$H(V) = \sum_k P(v_k) \log_2 \frac{1}{P(v_k)}$$

$$= - \sum_k P(v_k) \log_2 P(v_k)$$

$$B(q) = - (q \log_2 q + (1-q) \log_2 (1-q))$$

$$H(\text{"goal"}) = B\left(\frac{p}{p+n}\right)$$

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

$$\text{gain}(A) = H(\text{goal}) - \text{remainder}(A)$$