

Quiz 3 sample solution

Note Title

10/7/2007

Qn 1

- train on 1-6, test on 7-9:
a best rule is $N \mapsto A, D \mapsto B$
(with 1 error on 1-6), gives 1 error on 7-9
 - train on 1,3,7-9, test on 4-6:
a best rule is $N \mapsto A, D \mapsto B$
(with 1 error on 1,3,7-9), gives 1 error on 4-6
 - train on 4-9, test on 1-3:
a best rule is $Q \mapsto A, L \mapsto B$
(with 1 error on 4-9), gives 1 error on 1-3
- \Rightarrow 3/9 errors, error rate is 33.3%.

Qn 2

We choose to use the C4.5 recommended one-sided confidence interval of 25% for estimating conservative error rates. Therefore, cons. err. rates will be computed by adding 0.69 std devs to the regular error rate.

① Node D error rate

$$\text{err rate is } \frac{6}{20} = 0.3$$

$$\text{std dev is } \sqrt{\frac{0.3 \times 0.7}{20}} = 0.10$$

$$\therefore \text{cons. err. rate} = 0.3 + 0.69 \times 0.10 = 0.37$$

② Node E error rate

$$\text{err rate is } \frac{2}{10} = 0.2$$

$$\text{std dev is } \sqrt{\frac{0.2 \times 0.8}{10}} = 0.13$$

$$\therefore \text{cons. err. rate} = 0.2 + 0.69 \times 0.13 = 0.29$$

③ Expected conservative error rate of D & E combined:

$\frac{20}{30}$ instances at D, $\frac{10}{30}$ at E, so apply these

weights to get:

$$\text{expected cons. err. rate} = \frac{20}{30} \times 0.37 + \frac{10}{30} \times 0.29$$

$$= 0.34 \quad \text{---} \quad (\star)$$

④ Error rate at C if pruned

If pruned, would have 16 A and 14 B at C.

$$\text{err rate is } \frac{14}{30} = 0.47$$

$$\text{std dev is } \sqrt{\frac{0.47 \times 0.53}{30}} = 0.09$$

$$\therefore \text{cons. err. rate} = 0.47 + 0.69 \times 0.09 = 0.53 \quad \text{---} \quad (t)$$

⑤ Conclusion:

Since (t) is greater than (\star), do not prune.