

6.3) Problem Solving on Paper

A. Research the concept of “Information Gain Ratio” (which was discussed in lectures). In 100 words, discuss what is wrong with regular “Information Gain” and how does “Information Gain Ratio” partially remedy this. Explain with one example.

(total words: 100).

Answer: The issue with information gain is bias towards attributes with all distinct values, as information gain in that scenario is equivalent to the original total entropy of the data provided, thus always selecting that attribute as the first attribute to split tree, this could lead to model learning training data too well, leading to overfitting.

As in the example the distinct valued attribute is the Name, this has highest initial Information gain, which when divided by intrinsic information provided by the feature gives correct result.

EXAMPLE:

| Name | Time Spent on Unit(in hrs) | Grade in Unit |
|----------|----------------------------|---------------|
| Navdeep | 120 (0 Class) | HD |
| Michael | 110 (0 class) | HD |
| Abhishek | 80 (1 class) | C |
| Gurpreet | 90 (1 class) | D |

Split Info(Name)= $(-\frac{1}{4} * \log(\frac{1}{4}) + \frac{1}{4} * \log(\frac{1}{4}) + \frac{1}{4} * \log(\frac{1}{4}) + \frac{1}{4} * \log(\frac{1}{4})) = 0.60$

Split Info(Time Spent On Unit)= $(-\frac{2}{4} * \log(\frac{2}{4}) + \frac{2}{4} * \log(\frac{2}{4})) = 0.15$

Split Info(Grade in Unit)= $(-\frac{2}{4} * \log(\frac{2}{4}) + \frac{1}{4} * \log(\frac{1}{4}) + \frac{1}{4} * \log(\frac{1}{4})) = 0.45$

Class HD = 2 times p = 2
Class C = 1 time q = 1
Class D = 1 time r = 1

Total (n) = 4.

Entropy prior to split:

(Class HD) $q = \frac{2}{4} = \frac{1}{2}$
 $B(q) = -(q \log p + (1-q) \log (1-q))$
 $= -(\frac{1}{2} \log \frac{1}{2} + (\frac{1}{2} \log \frac{1}{2}))$
 $= 0.30$

(Class C) $q = \frac{1}{4}$
 $B(q) = -(\frac{1}{4} \log \frac{1}{4} + \frac{3}{4} \log \frac{3}{4})$
 $= 0.24$

(Class D) $q = \frac{1}{4}$
 $B(q) = 0.24$

Total Entropy $B(Q) = 0.24 + 0.24 + 0.30$
 $= 0.78$

Split on Time Spent:
Gain (Time) = $B(Q) - \sum_{k=1}^d (\frac{p_k + n_k}{p+n} B(\frac{p_k}{p_k + n_k}))$

Class HD: 2 in class 0 time spent 0 in class 1
Class C: 0 in class 0 time spent 1 in class 1
Class D: 0 in class 0 time spent 1 in class 1.

$$\text{Gain (Time)} = B(0) - \left[\frac{2}{2} B\left(\frac{2}{2}\right) + \frac{1}{2} B\left(\frac{1}{2}\right) + \frac{1}{2} B\left(\frac{1}{2}\right) \right]$$

$$= B(0) - \left[B(1) + \frac{1}{2} B\left(\frac{1}{2}\right) + \frac{1}{2} B\left(\frac{1}{2}\right) \right]$$

$$= B(0) - \left[B(1) + 2 \times B\left(\frac{1}{2}\right) \right]$$

$$= B(0) - \left[-\log 1 + 2 \times 0.30 \right]$$

$$= B(0) - \left[0 + 0.60 \right]$$

$$\text{Gain (Time)} = \frac{0.78 - 0.60}{0.48}$$

But now if gain is calculated on Names

$$\text{Gain (Names)} = 0.78 - 0$$

$$= 0.78$$

it would be equal to total entropy since $B(1)$ will always be 0, leading to be attribute 'Name' chosen first.

Since this is when the information gain fails, the ratio can be calculated by dividing the gain found with the Split info of relevant attribute, this gives true and corrected values which can be used in building decision tree.

$$\text{Gain Ratio (Names)} = \frac{0.78}{0.60} = \boxed{1.3}$$

$$\text{Gain Ratio (Time spent on Units)} = \frac{0.48}{0.15} = \boxed{3.2}$$

B. Research the concept of cross validation. In 100 words, discuss what is different between StratifiedKFold and ShuffleSplit. Explain with one example.

(100 words)

Answer:

| Stratified K Fold | Shuffle Split |
|---|--|
| 1)Each test should not overlap with previous split. 2) The data is first shuffled at the start and then split is done depending on the size, each time, testing and training dataset being distinct. | 1)Tests could overlap with previous split. 2) The data is shuffled every time the model is tested and split is done depending on the size, each time, testing and training dataset, having the chance of overlap with last split. |

Here is an example K fold shuffled once, shuffle split, every time data is shuffled.

| KFold | Shuffle Split |
|--|--|
| TRAIN: [0 2 3 4 5 6 7 9] TEST: [1 8] TRAIN: [0 1 2 3 5 7 8 9] TEST: [4 6] TRAIN: [0 1 3 4 5 6 8 9] TEST: [2 7] TRAIN: [1 2 3 4 6 7 8 9] TEST: [0 5] TRAIN: [0 1 2 4 5 6 7 8] TEST: [3 9] | TRAIN: [8 4 1 0 6 5 7 2] TEST: [3 9] TRAIN: [7 0 3 9 4 5 1 6] TEST: [8 2] TRAIN: [1 2 5 6 4 8 9 0] TEST: [3 7] TRAIN: [4 6 7 8 3 5 1 2] TEST: [9 0] TRAIN: [7 2 6 5 4 3 0 9] TEST: [1 8] |