

DSAA 5002 - Data Mining and Knowledge Discovery in Data Science

(Fall Semester 2023)

Homework 2

Deadline: 7 Nov 2023 11:59pm

Full Mark: 100 Marks

Please hand in via Canvas. You may submit your solutions several times to correct some mistakes, but please make sure that each submission is a whole submission.

Please start early. For each day overdue, the grade for this assignment will decrease by 10%.

In all the computations, if you need to round a number, please keep four decimal places, like

$$\frac{1000}{7} = 142.8571, \quad \frac{0.02}{7} = 0.0029$$

1. (25 marks) Consider the following training data with labels 0 and 1, and three attributes A, B, and C.

id	A	B	C	class
1	0.62	yes	yes	0
2	3.84	no	no	0
3	6.61	yes	no	0
4	6.87	yes	no	0
5	7.71	no	yes	0
6	8.98	no	yes	0
7	1.77	yes	no	0
8	2.02	yes	no	1
9	2.06	no	yes	1
10	2.66	no	yes	1
11	3.72	no	yes	1
12	4.98	yes	yes	1
13	5.73	yes	yes	1
14	6.29	yes	yes	1
15	9.08	no	no	1
16	9.45	no	no	1

- (a) (10 marks) Try threshold 2, 5, and 8 for attributes A (that is, use the “ $A > 2$, $A < 2$ ”, “ $A > 5$, $A < 5$ ”, and “ $A > 8$, $A < 8$ ” respectively). Use the Gini score to determine the best one θ_a among them. Recall

$$Gini(t) = 1 - \sum_{i=1}^c [p(i|t)]^2$$

(b) (15 marks) Use θ_a obtained above, and the Gini score, determine which attributes should firstly be used for developing a decision tree.

2. (30 marks) The table below is a small part of the Acute Inflammations Data Set.

- a1 Temperature of patient (35C-42C)
- a2 Occurrence of nausea (yes, no)
- a3 Lumbar pain (yes, no)
- a4 Urine pushing (continuous need for urination) (yes, no)
- a5 Micturition pains (yes, no)
- a6 Burning of urethra, itch, swelling of urethra outlet (yes, no)
- d1 Decision: Inflammation of urinary bladder (yes, no)
- d2 Decision: Nephritis of renal pelvis origin (yes, no)

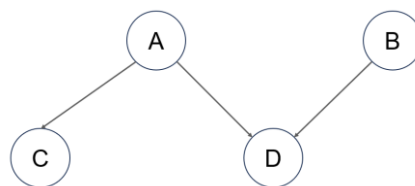
Here the attributes a1-a6 are observations, and the decisions d1 and d2 are made by a medical expert. The purpose of studying this data set is to predict presumptive diagnosis of two disease of the urinary system, namely, “Inflammation of urinary bladder” and “Nephritis of renal pelvis origin”.

a1	a2	a3	a4	a5	a6	d1	d2
37.3	no	yes	no	no	no	no	no
37.4	no	no	yes	no	no	yes	no
37.5	yes	yes	no	no	no	no	no
37.6	no	no	yes	yes	yes	yes	yes
37.7	no	no	yes	no	no	yes	no
37.7	no	no	yes	yes	no	yes	no
37.7	no	no	yes	yes	no	yes	no
37.8	no	yes	no	no	no	no	no
37.9	no	no	yes	yes	yes	yes	no
37.9	no	no	yes	no	no	yes	no
38.0	no	yes	yes	no	yes	no	yes
38.0	no	yes	yes	no	yes	no	yes
38.1	no	yes	yes	no	yes	yes	yes
38.3	no	yes	yes	no	yes	no	yes
38.5	no	yes	yes	no	yes	no	no
38.7	no	yes	yes	no	yes	no	yes
38.9	no	yes	yes	no	yes	yes	yes
39.0	no	yes	yes	no	yes	no	yes
39.4	no	yes	yes	no	yes	no	yes
39.5	no	yes	yes	no	yes	no	yes

- (a) (10 marks) Consider the procedures of building a decision tree with Gini score. If we plan only to use the attributes a3 and a5 to predict the decision d2, which attribute should we use first?
- (b) (20 marks) Use the naïve Bayes algorithm, the attributes a1 (with the threshold $\theta_1 = 37.95$), a2, and a3 only, to predict the decision d2 for the following data of a new patient. (For simplicity you do NOT need to use the Laplacian correction.)

a1	a2	a3	a4	a5	a6	d1	d2
40.0	yes	no	no	no	no	?	?

3. (15 marks) There is a BBN below, which comprises four Random Variables(RV). Each RV is a Boolean RV.



$$\begin{aligned}
 P(A) &= 0.1 & P(B) &= 0.5 & P(C|A) &= 0.7 \\
 P(C|\neg A) &= 0.2 & P(D|A, B) &= 0.9 & P(D|\neg A, B) &= 0.6 \\
 P(D|A, \neg B) &= 0.7 & P(D|\neg A, \neg B) &= 0.3 & &
 \end{aligned}$$

- (a) (7 marks) What is $P(\neg A, B, \neg C, D)$?
- (b) (8 marks) What is $P(A | B, C, D)$?
4. (30 marks) Consider a simple neural network with a single hidden layer. The input layer consists of three dimensional $\mathbf{x} = (x_1, x_2, x_3)^T$. The hidden layer includes two dimensional $\mathbf{h} = (h_1, h_2)$. The output layer includes one scalar o . We ignore bias terms for simplicity.

We use linear rectified (ReLU) as activation function **for the hidden and output layer BOTH**.

$$\begin{aligned}
 ReLU(x) &= \max(0, x) \\
 ReLU'(x) &= \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}
 \end{aligned}$$

Moreover, denote the loss function (also called *error* in slides) by $J(o, t) = \frac{1}{2}|o - t|^2$ where t is the associated label (target) value for scalar output o .

Denote by W and V weight matrices connecting input and hidden layer, and hidden layer and output respectively. They are **initialized** (i.e., the initial condition before first updating round) as follows:

$$W = \begin{bmatrix} 1 & 0 & 1 \\ -3 & -1 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} 0 & 1 \end{bmatrix}, \quad \text{Moreover, one training sample is } \mathbf{x} = (-1, 2, -1)^T, \quad t = 0.$$

Now, try to solve the following parts.

- (a) (5 marks) Write out symbolically (thus, no need to plug in the specific values of W and V just yet) the mapping $\mathbf{x} \rightarrow o$ using ReLU, W, V .
- (b) (10 marks) Given the condition $\mathbf{x} = (1, 2, 1)^T$, $t = 1$, compute the numerical output value o , clearly show all intermediate steps. You can reuse the results of the previous question.
- (c) (15 marks) Compute the gradient of the loss function with respect to the V weights, and evaluate the gradients at specific $\mathbf{x} = (1, 2, 1)^T$, $t = 1$.