

CISC3023 Machine Learning

1. Let us apply the naive Bayes classifier to the problem of classifying days according to whether someone will play tennis. Table below provides a set of 14 historical examples of the target concept PlayTennis, where each day is described by the attributes Outlook, Temperature, Humidity, and Wind. Please use naive Bayes classifier (do not use the “imaginary” samples) and the historical data from this table to classify the following novel instance; the target PlayTennis should be Yes or No? (**Show every step of your calculation**)
(Outlook = rain, Temperature = mild, Humidity = normal, Wind = strong)

<i>Day</i>	<i>Outlook</i>	<i>Temp</i>	<i>Humidity</i>	<i>Wind</i>	<i>Play Tennis</i>
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$$P(\text{play=yes})=9/14$$

$$P(\text{play=no})=5/14$$

$$\begin{aligned} P(\text{yes})P(\text{rain}|\text{yes})P(\text{mild}|\text{yes})P(\text{normal}|\text{yes})P(\text{strong}|\text{yes}) \\ =9/14 \times 3/9 \times 4/9 \times 6/9 \times 3/9 = 0.0211640212 \end{aligned}$$

It is larger than:

$$\begin{aligned} P(\text{no})P(\text{rain}|\text{no})P(\text{mild}|\text{no})P(\text{normal}|\text{no})P(\text{strong}|\text{no}) \\ =5/14 \times 2/5 \times 2/5 \times 1/5 \times 3/5 = 0.00685714286 \end{aligned}$$

So the answer is Yes.

2. Suppose that training examples are given in the following table. The last column is used to predict whether a particular type of rice will be good for making Sushi or not.

rice #	length	width	weight	color	Good for Sushi?
1	long	wide	light	white	yes
2	short	wide	heavy	yellow	no
3	short	narrow	light	white	yes
4	long	normal	light	red	no
5	long	normal	heavy	white	no

- a) Use **Find-S** Algorithm to find the most specific hypothesis that satisfies above training examples **Show every step of your derivation.** (20/100)

S1=<l w l w>
 S2=<l w l w>
 S3=<? ? l w>
 S4=<? ? l w>
 S5=<? ? l w>

- b) Use **candidate-elimination algorithm** to derive the final **version space** from the given examples. **Show every step of your derivation** (20/100)

S0=<∅ ∅ ∅ ∅ > G0=<? ? ? ?>
 S1=<l w l w> G1=<? ? ? ?>
 S2=<l w l w> G2=<l ? ? ?><? ? l ?><? ? ? w>
 S3=<? ? l w> G3=<? ? l ?><? ? ? w>
 S4=<? ? l w> G4=<? ? ? w>
 S5=<? ? l w> G5=<? ? l w>
 Version space <? ? l w>

- c) Calculate the entropy of current collection **Entropy(S)**, and Information Gains: **Gain(S, weight)**, i.e., the information gain if we set ‘**weight**’ as the root node of a decision tree (20/100)

$$\text{Entropy}(s) = -2/5 \log_2(2/5) - 3/5 \log_2(3/5) = 0.97$$

$$\text{Entropy}(\text{weight}=\text{light}) = -1/3 \log_2(1/3) - 2/3 \log_2(2/3) = 0.92$$

$$\text{Entropy}(\text{weight}=\text{heavy}) = -0/2 \log_2(0/2) - 2/2 \log_2(2/2) = 0$$

$$\text{Gain}(S, \text{weight}) = \text{Entropy}(s) - 3/5 \text{Entropy}(\text{weight}=\text{light}) - 2/5 \text{Entropy}(\text{weight}=\text{heavy}) = 0.42$$

3. Suppose that 5 cases are given in the case base in following table.

Case ID	x_1	x_2	x_3	Output $f(X)$
X_1	7	5	3	2
X_2	3	9	4	10
X_3	10	3	2	6
X_4	4	9	7	2
X_5	5	2	10	8

- a) For a query $x_q = (4, 8, 4)$, calculate the estimated output $f(x_q)$ using “Distance-weighted 3-Nearest Neighbor algorithm for discrete-valued target function”. (In this case, the last column in above table shows the discrete outputs. Then the outputs like 2 and 10 are discrete class labels, **not the numeric numbers**) (20/100)

$$d_i = |x_i - x_q|^{1/2}$$

So according to d_i^2 , we know 3NN of x_q are x_1 x_2 x_4

$$\text{Weight for } x_1: w_1 = (1/d_1)^2 = 1/[(7-4)^2 + (5-8)^2 + (3-4)^2] = 1/19$$

$$\text{Weight for } x_2: w_2 = (1/d_2)^2 = 1/[(3-4)^2 + (9-8)^2 + (4-4)^2] = 1/2$$

$$\text{Weight for } x_4: w_4 = (1/d_4)^2 = 1/[(4-4)^2 + (9-8)^2 + (7-4)^2] = 1/10$$

Outputs for x_1 , x_4 are label 2, their weights together is $1/19 + 1/10$

Output for x_2 is label 10, its weight is $1/2$

Then the weight for output 10 is larger, output for x_q is 10

- b) For a query $x_q = (4, 8, 4)$, calculate the estimated output $f(x_q)$ using “Distance-weighted 3-Nearest Neighbor algorithm for real-valued target function”. (the last column in above table shows the numeric outputs this time) (20/100)

$$\text{Output}(x_q) = [w_1 f(x_1) + w_2 f(x_2) + w_4 f(x_4)] / (w_1 + w_2 + w_4) = [1/19 * 2 + 1/2 * 10 + 1/10 * 2] / (1/19 + 1/2 + 1/10) = 8.1$$