

## Business Data Mining (IDS 572)

### Homework 5 (Question 1)-Solution

#### Question 1

(a)  $P(\text{senior}) = 5/11$  and  $P(\text{junior}) = 6/11$ .

(b)  $P(\text{department} | \text{status})$

Class	Sales	Systems	Marketing	Secretary
Senior	1/5	2/5	1/5	1/5
Junior	2/6 = 1/3	2/6 = 1/3	1/6	1/6

$P(\text{age} | \text{status})$

Class	21...25	26...30	31...35	36...40	41...45	46...50
Senior	0/5 = 0	0/5 = 0	1/5	2/5	1/5	1/5
Junior	1/6	3/6 = 1/2	2/6 = 1/3	0/6 = 0	0/6 = 0	0/6 = 0

$P(\text{salary} | \text{status})$

Class	26K-30K	31K-35K	36K-40K	41K-45K	46K-50K	66K-70K
Senior	0/5 = 0	0/5 = 0	1/5	0/5 = 0	2/5	2/5
Junior	2/6 = 1/3	1/6	0/6 = 0	1/6	1/6	1/6

(c) For the test instance  $A = \{\text{marketing}, 31...35, 46K-50K\}$  we have,

$$P(\text{senior} | A) = P(\text{senior})P(\text{marketing} | \text{senior})P(31...35 | \text{senior})P(46K-50K | \text{senior}) \\ = 5/11 \times 1/5 \times 1/5 \times 2/5 = 0.007$$

$$P(\text{junior} | A) = P(\text{junior})P(\text{marketing} | \text{junior})P(31...35 | \text{junior})P(46K-50K | \text{junior}) \\ = 6/11 \times 1/6 \times 1/3 \times 1/6 = 0.005$$

Therefore the label for this instance is “senior” since  $P(\text{senior} | A) > P(\text{junior} | A)$

For the instance  $B = \{\text{sale}, 31...35, 66K-70K\}$ ,

$$P(\text{senior} | B) = P(\text{senior})P(\text{sale} | \text{senior})P(31...35 | \text{senior})P(66K-70K | \text{senior}) \\ = 5/11 \times 1/5 \times 1/5 \times 2/5 = 0.007$$

$$P(\text{junior} | B) = P(\text{junior})P(\text{sale} | \text{junior})P(31...35 | \text{junior})P(66K-70K | \text{junior}) \\ = 6/11 \times 1/3 \times 1/3 \times 1/6 = 0.01$$

Therefore the label for this instance is “junior” since  $P(\text{senior} | B) < P(\text{junior} | B)$ .

(d) Now suppose  $A = \{\text{marketing}, 31\ldots 35, 46\text{K}-50\text{K}, 46\text{K}-50\text{K}\}$  and  $B = \{\text{sale}, 31\ldots 35, 66\text{K}-70\text{K}, 66\text{K}-70\text{K}\}$ . Therefore we have,

$$P(\text{senior} | A) = P(\text{senior})P(\text{marketing} | \text{senior})P(31\ldots 35 | \text{senior})P(46\text{K}-50\text{K} | \text{senior}) \\ P(46\text{K}-50\text{K} | \text{senior}) = 5/11 \times 1/5 \times 1/5 \times 2/5 \times 2/5 = 0.003$$

$$P(\text{junior} | A) = P(\text{junior})P(\text{marketing} | \text{junior})P(31\ldots 35 | \text{junior})P(46\text{K}-50\text{K} | \text{junior}) \\ P(46\text{K}-50\text{K} | \text{junior}) = 6/11 \times 1/6 \times 1/3 \times 1/6 \times 1/6 = 0.0008$$

Therefore the label for the instance A is “senior” since  $P(\text{senior} | A) > P(\text{junior} | A)$

$$P(\text{senior} | B) = P(\text{senior})P(\text{sale} | \text{senior})P(31\ldots 35 | \text{senior})P(66\text{K}-70\text{K} | \text{senior}) P(66\text{K}-70\text{K} | \text{senior}) = 5/11 \times 1/5 \times 1/5 \times 2/5 \times 2/5 = 0.003$$

$$P(\text{junior} | B) = P(\text{junior})P(\text{sale} | \text{junior})P(31\ldots 35 | \text{junior})P(66\text{K}-70\text{K} | \text{junior}) P(66\text{K}-70\text{K} | \text{junior}) = 6/11 \times 1/3 \times 1/3 \times 1/6 \times 1/6 = 0.02$$

Therefore the label for the instance B is “junior” since  $P(\text{senior} | B) < P(\text{junior} | B)$ .

(e) The assumption of independent inputs is clearly violated in part (d). But Naïve Bayes model does not consider the dependency of the input variables. In addition, the occurrence of A and B is zero. In general the Naïve Bayes model could compute large probabilities even for the cases that have very low occurrence.