

## Individual Assignment 4

### Topics: K-NN and Naïve Bayes, Classification vs. Ranking

105078702 Wendy Huang

#### Submission instruction:

1. Upload a single PDF to Canvas. It should include all the screenshots and answers 1-8.
2. Hand in the printed copy of your document in class.

#### Data Details and Goal:

The file “user\_changes.xlsx” contains information on EZTABLE restaurant reservations by members and an analysis. Worksheet “train” (also available as file train.txt) is the data from the Kaggle contest: it includes the first reservation made by EZTABLE members. For each reservation, there is information that the user entered at the time of booking, and column **status** which is entered by EZTABLE. **status** has two values that indicate a change that a user made on the EZTABLE website after his/her initial booking: “change” means that the user changed the reservation in some way. “canceled” means that the user canceled it online.

We want to build a model that predicts, right after a booking, whether the user will change/cancel that booking in the future (**status** = *changed* or *canceled*).

#### Data Preparation:

- Set random seed to 12345. (R: `set.seed(12345)`)
- Take a random sample of 10,000 records (from worksheet *train*).
- Create column **cancel\_change** that takes 1 if *status*={canceled or changed}, 0 otherwise.
- Create column **time\_to\_booking** from *datetime* (restaurant reservation date) and *cdate*
- Partition the data randomly into training/validation/test sets using 50%-30%-20%. Keep *booking\_id*, *time\_to\_booking*, *people*, and *cancel\_change*

#### Naïve Benchmark

1. Using the information in the training dataset, if a member just logged in but did not complete their booking (we have no information), would we classify the booking as **cancel\_change** =1 or =0? What is the probability that this record will be canceled/change?

Classification: \_\_0\_\_

Probability of cancellation/change = \_\_21.48%\_\_

**K-NN**

Run a K-NN classifier on the partitioned datasets: Use predictors *people* and *time\_to\_booking* and outcome *cancel\_change* (with 1 as the success class). Cutoff = 0.5. Normalize the predictor variables (XLMiner: “normalize”; R use `preProcess()` - see Table 7.2 in book); Let the validation set find the best k between 1-20 (XLMiner: “score on be k between 1 and specified value”, choose 20 or the largest value allowed; R - see Table 7.3 in book)

2. From the output (XLMiner: worksheets KNNC\_Output, KNNC\_TestLiftChart) fill in:

Best k = \_\_\_\_\_ 19 \_\_\_\_\_

Training overall error = \_\_\_\_\_ 21.88 \_\_\_\_\_ %

Validation overall error = \_\_\_\_\_ 21.36 \_\_\_\_\_ %

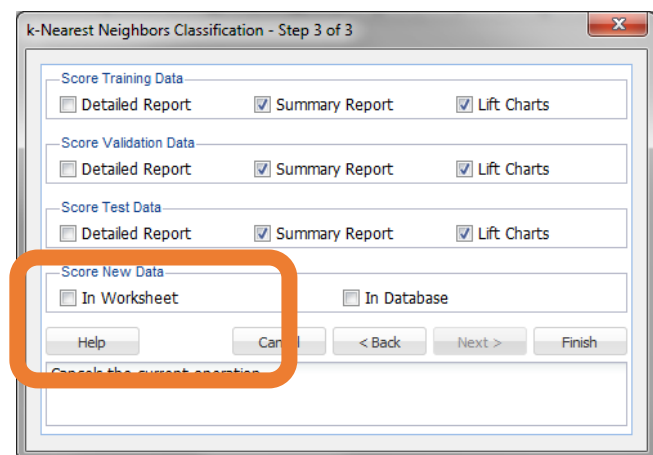
Test overall error = \_\_\_\_\_ 21.05 \_\_\_\_\_ %

Test Lift for first decile = \_\_\_\_\_ 1.39 \_\_\_\_\_ (R: use `gains()`, see Fig 5.7 in book)

3. Fill in the predicted probabilities and classifications **for the 10 records in the table below.**

XLMiner: copy the three left columns into a new worksheet; re-run KNN, in Step 3 click “Score New Data” in Worksheet and choose this worksheet.

R - see Table 7.4 in book.



Booking #	people	time_to_booking	KNN: Prob of 1	KNN Classification (0/1)
1	2	27	0.3157895	0
2	2	63	0.3000000	0
3	2	1.85	0.4210526	0
4	10	1.7	0.4210526	0
5	2	13	0.3157895	0
6	11	8.7	0.2631579	0
7	2	1.8	0.4210526	0
8	4	7	0.2631579	0
9	2	2.5	0.2105263	0
10	6	10	0.2105263	0

**Naïve Bayes**

From the SampleWS worksheet (or file SampleWS.txt), convert *people* into a categorical variable *Binned\_people*, by splitting it into **two** categories with an equal number of records. Similarly, convert *time\_to\_booking* into *Binned\_time\_to\_booking* with two categories with an equal number of records.

- XLMiner: Transform > Bin Continuous Data
- R: use function `median`: `Binned_people <- ifelse(people < median(people), 1, 2)`

Partition the data again (50-30-20%) and keep *booking\_id*, *cancel\_change*, and the two Binned columns.

Run a naïve Bayes classifier on the partitioned dataset using only the two new categorical predictors you created. Keep all other options as in the KNN. R: see Table 8.4 in textbook.

4. From the output (XLMiner worksheets NB\_Output, NNB\_TestLiftChart), fill in:

Training overall error = \_\_\_\_\_ 21.43 \_\_\_\_\_ %

Validation overall error = \_\_\_\_\_ 20.52 \_\_\_\_\_ %

Test overall error = \_\_\_\_\_ 22.06 \_\_\_\_\_ %

Test Lift for first decile = \_\_\_\_\_ 1.147 \_\_\_\_\_ (R: use `gains()`, see Fig 5.7 in book)

5. Fill in the predicted probabilities and classifications for the 10 records in the table below.

You'll need to add their Binned values before NB can score them. R: see Table 8.6 in book.

Booking #	people	time_to_booking	Binned_people	Binned_time_to_booking	NB: Prob of 1	NB Classification (0/1)
1	2	27	1	2	0.2488304	0
2	2	63	1	2	0.2488304	0
3	2	1.85	1	1	0.1991569	0
4	10	1.7	2	1	0.1839551	0
5	2	13	1	2	0.2488304	0
6	11	8.7	2	2	0.2309303	0
7	2	1.8	1	1	0.1991569	0
8	4	7	2	2	0.2309303	0
9	2	2.5	1	1	0.1991569	0
10	6	10	2	2	0.2309303	0

### Exact Bayes

6. Use a pivot table to compute the **count** of *cancel\_change*=1 and 0 in the **training dataset** as a function of the two categorical predictors.

	cancel_change=1	cancel_change=0
Binned_people = 1, Binned_time_to_booking = 1	652	2473
Binned_people = 1, Binned_time_to_booking = 2	772	3585
Binned_people = 2, Binned_time_to_booking = 1	720	2293
Binned_people = 2, Binned_time_to_booking = 2	1070	3435

7. Compute manually the **exact** Bayes conditional probabilities of a cancel\_change=1, given *Binned\_people* and *Binned\_time\_to\_booking* as the two predictor variables. Show your calculation for the first record.

Then, use cutoff=0.5 to get the classifications.

Booking #	people	time_to_booking	Binned_people	Binned_time_to_booking	Exact Bayes: Prob of 1	Exact Bayes Classification (0/1)
1	2	27	1	2	0.2401991	0
2	2	63	1	2	0.2401991	0
3	2	1.85	1	1	0.2028625	0
4	10	1.7	2	1	0.2240199	0
5	2	13	1	2	0.2401991	0
6	11	8.7	2	2	0.3329185	0
7	2	1.8	1	1	0.2028625	0
8	4	7	2	2	0.3329185	0
9	2	2.5	1	1	0.2028625	0
10	6	10	2	2	0.3329185	0

8. Compare the Naïve Bayes and exact Bayes scores for the 10 records.

Are the resulting probabilities identical? No (choose one)

Are the resulting classifications identical? Yes (choose one)

Is the ranking (= ordering) of observations identical? No (choose one)