

# Homework 1

Miriam Wagner  
373045

3. Mai 2018

## Question 1

I used for solving this exercise Matlab. My code you can see here:

```
1 amount_req = [7000, 15000, 23112, 6000, 10000,2500,
    19000,
    5000,13000,35000,5000,7500,6000,6500,25000,10000,7800,8500,17000,3000];

2 case_duration =
    [0.29309,28.43964,0.000532,0.048206,12.95023,0.021134,19.78213,29.51885,0.515

3 total_activities =
    [6,74,3,25,26,7,45,46,10,88,25,11,22,6,41,3,52,3,12,15];

4 c1 = [50000, 1.5, 6];
5 c2 = [310000, 0.7, 10];
6 c3 = [30000, 4, 26];
7
8 clust1x = [];
9 clust1y = [];
10 clust1z = [];
11 clust2x = [];
12 clust2y = [];
13 clust2z = [];
14 clust3x = [];
15 clust3y = [];
16 clust3z = [];
17 change = 100000;
18 k = 0;
19 while change > 10^-15
20     for i = 1:length(amount_req)
21         distc1(i) = distance([amount_req(i),case_duration
            (i),total_activities(i)],c1);
```

```

22     distc2(i) = distance([amount_req(i),case_duration
23                          (i),total_activities(i)],c2);
24     distc3(i) = distance([amount_req(i),case_duration
25                          (i),total_activities(i)],c3);
26
27     if (distc1(i) == min([distc1(i),distc2(i),distc3(
28                          i)]))
29         clust1x = [ clust1x , amount_req(i)];
30         clust1y = [ clust1y , case_duration(i)];
31         clust1z = [ clust1z , total_activities(i)];
32     elseif (distc2(i) == min([distc1(i),distc2(i),
33                          distc3(i)]))
34         clust2x = [ clust2x , amount_req(i)];
35         clust2y = [ clust2y , case_duration(i)];
36         clust2z = [ clust2z , total_activities(i)];
37     else
38         clust3x = [ clust3x , amount_req(i)];
39         clust3y = [ clust3y , case_duration(i)];
40         clust3z = [ clust3z , total_activities(i)];
41     end
42     c1new = centroid([clust1x;clust1y;clust1z]);
43     c2new = centroid([clust2x;clust2y;clust2z]);
44     c3new = centroid([clust3x;clust3y;clust3z]);
45     change = distance(c1,c1new) + distance(c2,c2new) +
46             distance(c3,c3new);
47     c1 = c1new;
48     c2 = c2new;
49     c3 = c3new;
50     k = k + 1;
51     namen = {'distance_to_cluster_1','
52             distance_to_cluster_2','distance_to_cluster_3'};
53     T = table(distc1','distc2','distc3','VariableNames',
54             namen);
55     filename = 'Question1.xlsx';
56     writetable(T,filename,'Sheet',k,'Range','A1');
57     namen2 = {'centroid_1','centroid_2','centroid_3'};
58     T = table(c1',c2',c3','VariableNames',namen2);
59     filename = 'Question1centroids.xlsx';
60     writetable(T,filename,'Sheet',k,'Range','A1')
61 end
62 c1
63 c2
64 c3

```

```

61 k
62
63     namen3 = { 'amount_req', 'case_duration', '
64               total_activities' };
65     T = table(clust1x', clust1y', clust1z', 'VariableNames',
66             namen3);
67     filename = 'Question1cluster1.xlsx';
68     writetable(T, filename, 'Sheet', 1, 'Range', 'A1')
69
70     namen4 = { 'amount_req', 'case_duration', '
71               total_activities' };
72     T = table(clust2x', clust2y', clust2z', 'VariableNames',
73             namen4);
74     filename = 'Question1cluster2.xlsx';
75     writetable(T, filename, 'Sheet', 1, 'Range', 'A1')
76
77     namen5 = { 'amount_req', 'case_duration', '
78               total_activities' };
79     T = table(clust3x', clust3y', clust3z', 'VariableNames',
80             namen5);
81     filename = 'Question1cluster3.xlsx';
82     writetable(T, filename, 'Sheet', 1, 'Range', 'A1')

```

The functions used I does not write down them here. 'distance' calculates the euclidian distance and centroid sums up all vectors and divides by the number of vectors.

My clustering algorithm now calculates for every Datavector the distances to every centroid and then checks which centroid is closed. The vector is add to the cluster. When all datavectors are put in a cluster the new cluster centroids are calculated.

My abort criterion is the change between the centroids. If they still change I will apply the algorithms again. Because of the computer accuracy I do not use 0, but  $10^{-15}$ . The algorithm finds after two rounds already good enough clusters. The distances in the first round are to see in 1

The new centroids are in 2

And the centroids in 4 The next round the distances are in 3

Cluster1 contains 5

Cluster 2 contains as in 6

Cluster 3 in 7

Tabelle 1: Distances first round		
distance_to_cluster_1	distance_to_cluster_2	distance_to_cluster_3
0	4	20
68	64	48
3	7	23
19	15	1
20	16	0
1	3	19
39	35	19
40	36	20
4	0	16
82	78	62
19	15	1
5	1	15
16	12	4
0	4	20
35	31	15
3	7	23
46	42	26
3	7	23
6	2	14
9	5	11

Tabelle 2: centroid		
centroid_1	centroid_2	centroid_3
7000	13000	15000
0,29309	0,515625	28,43964
6	10	74

Tabelle 3: Distances second round		
distance_to_cluster_1	distance_to_cluster_2	distance_to_cluster_3
0	4	68
68	64	0
3	7	71
19	15	49
20	16	48
1	3	67
39	35	29
40	36	28
4	0	64
82	78	14
19	15	49
5	1	63
16	12	52
0	4	68
35	31	33
3	7	71
46	42	22
3	7	71
6	2	62
9	5	59

Tabelle 4: centroids second round		
centroid_1	centroid_2	centroid_3
7000	13000	15000
0,29309	0,515625	28,43964
6	10	74

Tabelle 5: Cluster 1		
amount_req	case_duration	total_activities
7000	0,29309	6
23112	0,000532	3
2500	0,021134	7
6500	0,002049	6
10000	0,000799	3
8500	0,000486	3
7000	0,29309	6
23112	0,000532	3
2500	0,021134	7
6500	0,002049	6
10000	0,000799	3
8500	0,000486	3

Tabelle 6: Cluster 2

amount_req	case_duration	total_activities
13000	0,515625	10
7500	0,489861	11
17000	0,01375	12
3000	6,955382	15
6000	0,048206	25
10000	12,95023	26
13000	0,515625	10
5000	7,612419	25
7500	0,489861	11
6000	6,503808	22
25000	21,14506	41
17000	0,01375	12
3000	6,955382	15

Tabelle 7: cluster 3

amount_req	case_duration	total_activities
15000	28,43964	74
6000	0,048206	25
10000	12,95023	26
19000	19,78213	45
5000	29,51885	46
35000	19,74352	88
5000	7,612419	25
6000	6,503808	22
25000	21,14506	41
7800	19,1099	52
15000	28,43964	74
19000	19,78213	45
5000	29,51885	46
35000	19,74352	88
7800	19,1099	52



## Question 2

<new process\*>- RapidMiner Studio Trial 8.1.003 @ LAPTOP-SEGBH6SS

File Edit Process View Connections Cloud Settings Extensions Help

Show rules matching  
all of these conclusions: ▼  
A\_ACCEPTED  
Q\_ACCEPTED  
A\_APPROVED

Min. Criterion:  
confidence  
Min. Criterion Value:

Result History AssociationRules (Create Association Rules) ExampleSet (/Local Repository/data/BPI/DataHA1)

Data Graph Description Annotations

Premises	Conclusion	Support	Confidence	Lift
Q_ACCEPTED	A_ACCEPTED	0.400	0.889	1.778
Q_ACCEPTED	A_APPROVED	0.400	0.889	2.222
Q_ACCEPTED	A_ACCEPTED, A_APPROVED	0.400	0.889	2.222
A_APPROVED	A_ACCEPTED	0.400	1	2
A_APPROVED	Q_ACCEPTED	0.400	1	2.222
A_ACCEPTED, Q_ACCEPTED	A_APPROVED	0.400	1	2.500
A_APPROVED	A_ACCEPTED, Q_ACCEPTED	0.400	1	2.500
A_ACCEPTED, A_APPROVED	Q_ACCEPTED	0.400	1	2.222
Q_ACCEPTED, A_APPROVED	A_ACCEPTED	0.400	1	2

Find data, operators, etc. Search All Studio Auto Model Design Results

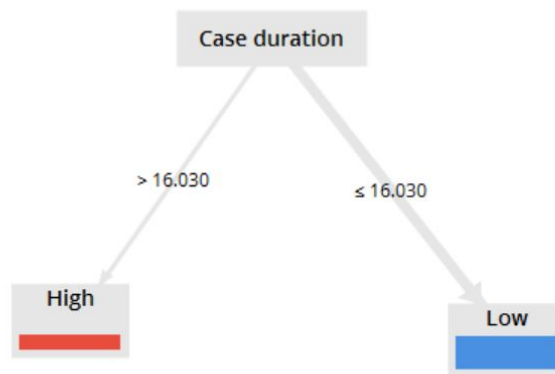
08:01 2-5-2018



I would pick the rule  $\{A\_ACCEPTED, O\_ACCEPTED\} \Rightarrow \{A\_APPROVED\}$  and  $\{A\_APPROVED\} \Rightarrow \{A\_ACCEPTED, O\_ACCEPTED\}$ , because they have the highest lift, confidence and support. When you have a closer look you will see that the sets are probably logical equivalent. Always pick the rule with the best lift, confidence and support.

### Question 3

The found decision tree is



If you check the Confusionmatrix

accuracy: 100.00%

	true Low	true High	class precision
pred. Low	14	0	100.00%
pred. High	0	6	100.00%
class recall	100.00%	100.00%	

you see, that this decision tree classifys the data perfectly. So you can predict by just knowing the case duration the total activities. If the case duration is higher, than also the total activities are high. This seems to be logical, if you have to do a lot this takes most of the times longer and otherwise around, if you do not need long you mostly did not do a lot of different things in the time.

### Question 4

1.

$$L1 = [\langle a, b, e, f \rangle, \langle a, b, e, c, d, b, f \rangle, \langle a, b, c, e, d, b, f \rangle, \langle a, b, c, d, e, b, f \rangle, \langle a, e, b, c, d, b, f \rangle]$$

The  $\alpha$ -Algorithm gives the following:

$$T_L = \{a, b, c, d, e, f\}$$

$$T_I = \{a\}$$

$$T_O = \{f\}$$

	a	b	c	d	e	f
a	#	→	#	#	→	#
b	#	#	→	#		→
c	#	#	#	→		#
d	#	→	#	#		#
e	#				#	→
f	#	←	#	#	#	←

$$X_L = \{(\{a\}, \{b\}), (\{a\}, \{e\}), (\{b\}, \{c\}), (\{b\}, \{f\}), (\{c\}, \{d\}), (\{b\}, \{c, f\}), (\{d\}, \{b\}), (\{e\}, \{f\}), (\{a, d\}, \{b\})\}$$

$$Y_L = \{(\{a\}, \{e\}), (\{b\}, \{c, f\}), (\{b\}, \{f\}), (\{c\}, \{d\}), (\{e\}, \{f\}), (\{a, d\}, \{b\})\}$$

