

25.01.2019

1

First name	Name	Matr.-Nr.

RWTH Aachen
Lehrstuhl für Informatik 9
Prof. Dr. van der Aalst

Trial Exam
Introduction to Data Science
January 25th, 2019

Study course:

- Diplom Informatik • Master Informatik • Master Data Science
- Master SSE • Master Media Informatics • Other: _____

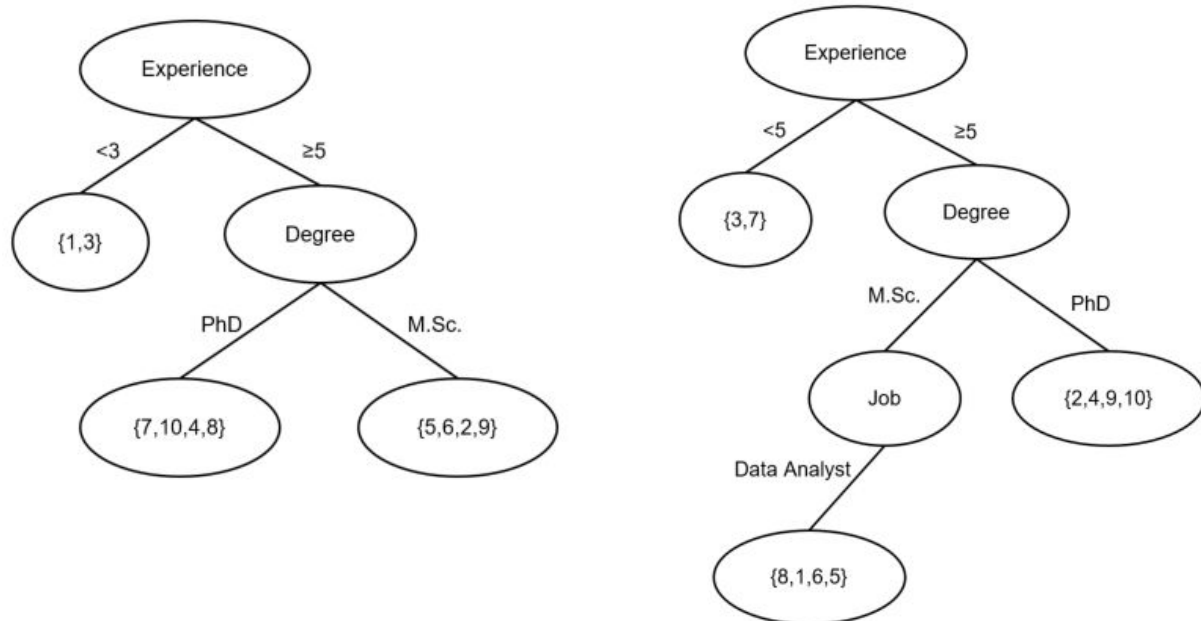
- Duration of the exam:
- Write your first name, name and Matrikelnummer on each sheet.
- Give your solutions in readable and traceable manner. **Solutions will be graded based on completeness and correctness of description/application of the algorithm/method.**
- Give your solutions on the exam sheets only. If you need extra paper, use only paper provided by the supervisors. Make sure to specify your name and Matrikelnummer on all papers.
- Please cross out those things you do not wish to be graded.
- In case of attempted deception, your exam will be graded as **failed**.
- At the end of the exam hand in your complete copy. Do not separate any sheets by removing the staples.
- **You may only use a black or blue pen, and a basic calculator; no additional material (e.g. books, cell phones, laptop, etc.) is allowed.**
- **Only answers that are given in English will be graded.**

Signature _____

First name	Name	Matr.-Nr.

1. Decision Tree

Suppose that we have two different decision trees to classify people with respect to their salaries (leaves of the trees). Which classifier is better when the impurity is considered as the variance within a leaf?



Solution:

They need to calculate variance of each leaf and the mean of the variances for each classifier. If the mean of the variances is lower then the corresponding classifier is better.

(This question can be made quite tricky, if in the question it is mentioned that whether the number of instances in each leaf matters or not. In that case, if the number of instances in each leaf matters, weighted mean should be used. Otherwise, simple mean should be used.)

$$\text{var}(a) = \frac{\sum_{i=1}^n (a_i - \bar{a})^2}{n - 1}$$

For the first classifier:

25.01.2019

3

First name	Name	Matr.-Nr.

$\text{var}(\text{Leaf1}) = 2, \text{var}(\text{Leaf2}) = 6.25, \text{var}(\text{Leaf3}) = 8.3$

Mean = 5.51

Weighted mean = 6.2

For the second classifier:

$\text{var}(\text{Leaf1}) = 8, \text{var}(\text{Leaf2}) = 8.6, \text{var}(\text{Leaf3}) = 14.91$

Mean = 10.5

Weighted mean = 11

Since the mean of the variances for the first classifier is lower, the first classifier is better.

First name	Name	Matr.-Nr.

2. Regression

Television time / day in hours	Length of deep sleep/day in h
0.3	5.8
2.2	4.4
0.5	6.5
1.8	5.0
0.2	6

- a. Based on the data minimize the squared error function to calculate the linear regression function that predicts the length of deep sleep based on the television time.

Solution:

We use the following two formulas to calculate the linear regression function:

$$b = \bar{y} - m\bar{x}$$

$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\bar{x} = 1 \text{ and } \bar{y} = 5.54$$

$$m = -0.818$$

$$b = 6.36$$

$$\text{Linear regression function : } y = 6.36 - 0.818x$$

- b. Interpret the regression function obtained in question a). What is the relation between television time and the length of deep sleep described by the function?

Solution:

With a television time of 0 hours (no television watching during the day) the amount of deep sleep is predicted to be 6.63 hours. As the amount of television time increases the amount of deep sleep decreases. As the value of x increases, from a certain value on the hours of deep

First name	Name	Matr.-Nr.

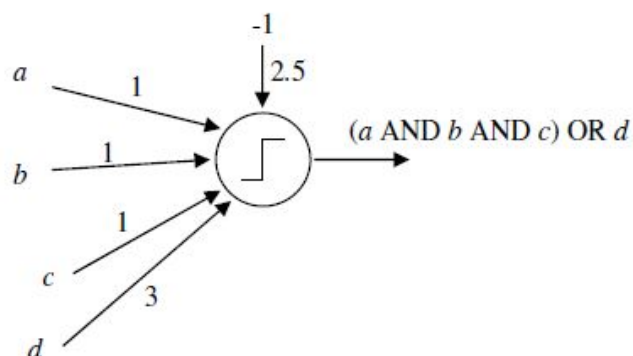
sleep will be predicted negative. Given the context, we know that these values are not realistic and the regression model should only be considered for positive values.

3. Neural Networks

- a. Draw a neuron or a neural network for the Boolean function ((a AND b AND c) OR d). Is it possible to have the function with one neuron? Specify the weights of the network. (Activation function is step function with threshold of zero(0))

Solution:

yes, it is possible to have an neuron with 5 inputs and one output to perform the above function.



First name	Name	Matr.-Nr.

4. Evaluation of Supervised Learning Algorithms

- a. The following table summarizes the results of a statistical regression. Compute the Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) metrics.

Forecast	Observation
7.1	6.3
10.9	10.2
12.3	14.0
10.5	10.9
10.3	7.7
8.8	5.2

Solution:

RMSE = 1.99, MAE = 1.63

First name	Name	Matr.-Nr.

5. Clustering

Let C_1 and C_2 be two clusters generated in the n^{th} iteration by k-means algorithm, where C_1 contains the elements (11, 68), (37, 52), (28, 70), (37, 55), (44, 56), (22, 57) and (45, 52), C_2 contains the elements (42, 66), (59, 62), (43, 63), (48, 64), (42, 61), (53, 61), (47, 61), (34, 67) and (48, 53). Each element is characterized by two values. When the $(n+1)^{\text{th}}$ iteration of k-means algorithm finishes:

- Please write down the elements for cluster C_1 and C_2 respectively (use Euclidean Distance).
- Please write down the centroids for cluster C_1 and C_2 (use Euclidean Distance).
The results should be accurate up to 2 decimal places

Solution:

Cluster 1: (11, 68), (37, 52), (28, 70), (37, 55), (34, 67), (22, 57)

Cluster 2: (42, 66), (59, 62), (43, 63), (48, 64), (42, 61), (53, 61), (44, 56), (47, 61),
(48, 53), (45, 52)

Center of cluster 1: (28.17, 61.5)

Center of cluster 2: (47.1, 59.9)

First name	Name	Matr.-Nr.

6. Frequent Items Sets and Association Rules

The table below shows a set of frequent itemsets L_2 of length 2. Let C_3 be the set of candidate itemsets of length 3 which is generated from L_2 by using Apriori algorithm.

Itemset	Support
{Apple, Beer}	3
{Apple, Coffee}	6
{Coffee, Orange}	10
{Coffee, Muffin}	2
{Muffin, Orange}	5

- a. Please show the two versions of C_3 generated after the itemsets merging step and the pruning step.

Solution:

	Itemset
C_3 generated after itemsets merging step:	{Apple, Beer, Coffee}
	{Coffee, Orange, Muffin}

	Itemset
C_3 generated after the pruning step:	{Coffee, Orange, Muffin}

First name	Name	Matr.-Nr.

Presume that there are three association rules from the transaction set shown in the table below, which are $\{\text{Beef, Coffee}\} \Rightarrow \{\text{Orange}\}$, $\{\text{Coffee}\} \Rightarrow \{\text{Apple, Orange}\}$ and $\{\text{Apple}\} \Rightarrow \{\text{Orange}\}$.

TID	Set of items
1	{Beef, Coffee, Muffin, Orange}
2	{Coffee, Orange, Apple}
3	{Beef, Apple, Juice, Orange}
4	{Coffee, Apple, Orange}
5	{Apple, Beef, Coffee, Orange}
6	{Apple, Beef, Juice, Orange}
7	{Coffee, Juice, Orange}
8	{Apple, Juice}

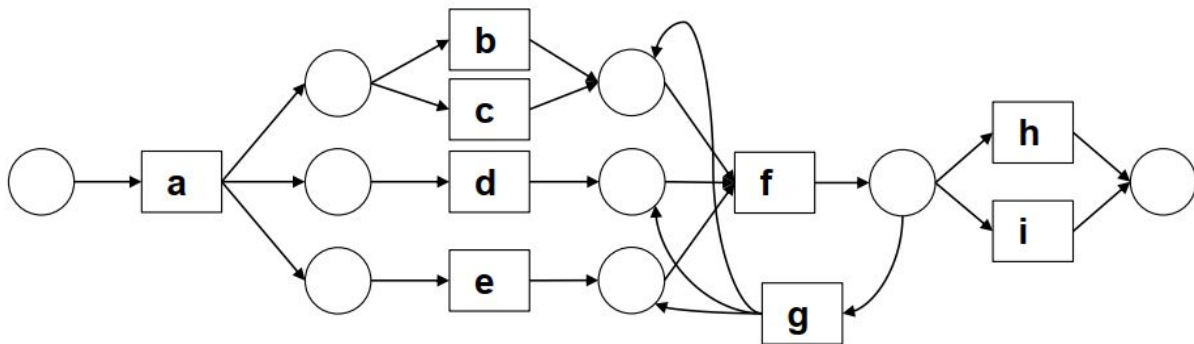
- b. Please calculate and write down the confidence and the lift of each of the three association rules. The results should be accurate up to 3 decimal places.

solution:

Rule	Confidence	Lift
$\{\text{Beef, Coffee}\} \Rightarrow \{\text{Orange}\}$	1	1.143
$\{\text{Coffee}\} \Rightarrow \{\text{Apple, Orange}\}$	0.6	0.96
$\{\text{Apple}\} \Rightarrow \{\text{Orange}\}$	0.833	0.952

First name	Name	Matr.-Nr.

8. Process Mining



- a. Compute the log fitness of the following event log on the given Petri net using token-based replay. For each trace, indicate the token counts.

[<a, d, c, e, f, g, f, h>⁷⁰,

<a, d, c, f, h>²⁰,

<a, d, c, e, f, g, f, h, f>¹⁰]

Solution:

First variant:

produced = 13
consumed = 13
missing = 0
remaining = 0

Second variant:

produced = 8
consumed = 8
missing = 1
remaining = 1

Third variant:

produced = 14
consumed = 16
missing = 3
remaining = 1

Log fitness:

$$\frac{1}{2} \left(1 - \frac{20 \times 1 + 10 \times 3}{70 \times 13 + 20 \times 8 + 10 \times 16} \right) + \frac{1}{2} \left(1 - \frac{20 \times 1 + 10 \times 1}{70 \times 13 + 20 \times 8 + 10 \times 14} \right)$$

First name	Name	Matr.-Nr.

= 0.97

9. Text Mining

- a. Given the following document corpus, compute the bigram probability estimates for the following bigrams. Make sure to add padding to the text with reserved symbols for the start and end of a document.

D1: "Process Mining is a discipline that perform analysis on process data."

D2: "Process Mining is a subfield of Data Science."

D3: "Data Science provides methods and techniques to obtain analysis from data."

D4: "Process Mining relates to other Process Science disciplines like Business Process Management."

D5: "Process Mining and Data Mining are related disciplines."

B1: process mining

B2: data science

B3: <s> process

Solution:

Probability of the bigram "process mining": #occurrences of "process mining" / #occurrences of "process"

$$P(B1) = 4/7 = 0.57$$

$$P(B2) = 2/5 = 0.4$$

$$P(B3) = 4/5 = 0.8$$

First name	Name	Matr.-Nr.

10. Responsible Data Science

For the following data:

- Given {"Age", "Gender", "State of domicile", "Religion"} as the Quasi-Identifier, does this data have 2-anonymity? If so, Identify the equivalence classes.
- What is the maximum value for l to have distinct l -diversity?
- What is the maximum value for l to have entropy l -diversity?

$\log(2) = 1$, $\log(1.9)=0.92$, $\log(1.8)=0.84$, $\log(1.7)=0.76$, $\log(1.6)=0.67$, $\log(1.5)=0.58$

Name	Age	Gender	State of domicile	Religion	Product
*	$20 < \text{Age} \leq 25$	Female	*	Hindu	Pea
*	$20 < \text{Age} \leq 25$	Female	*	Hindu	Bean
*	$20 < \text{Age} \leq 25$	Female	*	Muslim	Peanut
*	$20 < \text{Age} \leq 25$	Male	*	Buddhist	Pea
*	$20 < \text{Age} \leq 25$	Female	*	Muslim	Bean
*	$20 < \text{Age} \leq 25$	Male	*	Buddhist	Lentil
*	$\text{Age} \leq 20$	Male	*	Christian	Peanut
*	$20 < \text{Age} \leq 25$	Male	*	Buddhist	Lentil
*	$\text{Age} \leq 20$	Male	*	Christian	Peanut
*	$\text{Age} \leq 20$	Male	*	Christian	Pea

Solution:

(a)

First name	Name	Matr.-Nr.

Name	Age	Gender	State of domicile	Religion	Product
*	20 < Age ≤ 25	Female	*	Hindu	Pea
*	20 < Age ≤ 25	Female	*	Hindu	Bean
*	20 < Age ≤ 25	Female	*	Muslim	Peanut
*	20 < Age ≤ 25	Male	*	Buddhist	Pea
*	20 < Age ≤ 25	Female	*	Muslim	Bean
*	20 < Age ≤ 25	Male	*	Buddhist	Lentil
*	Age ≤ 20	Male	*	Christian	Peanut
*	20 < Age ≤ 25	Male	*	Buddhist	Lentil
*	Age ≤ 20	Male	*	Christian	Peanut
*	Age ≤ 20	Male	*	Christian	Pea

(b) It is clear that the maximum value for distinct I-diversity is 2.

(c)

Name	Age	Gender	State of domicile	Religion	Product	
*	20 < Age ≤ 25	Female	*	Hindu	Pea	→ Entropy = 1
*	20 < Age ≤ 25	Female	*	Hindu	Bean	
*	20 < Age ≤ 25	Female	*	Muslim	Peanut	→ Entropy = 1
*	20 < Age ≤ 25	Male	*	Buddhist	Pea	
*	20 < Age ≤ 25	Female	*	Muslim	Bean	→ Entropy = 0.92
*	20 < Age ≤ 25	Male	*	Buddhist	Lentil	
*	Age ≤ 20	Male	*	Christian	Peanut	→ Entropy = 0.92
*	20 < Age ≤ 25	Male	*	Buddhist	Lentil	
*	Age ≤ 20	Male	*	Christian	Peanut	
*	Age ≤ 20	Male	*	Christian	Pea	

$$Entropy(E) \geq \log(l) \quad \log(l) = 0.92 \quad l = 1.9$$