i Cover Page

Department of Computer Science

Examination paper for TDT4300 Data warehousing and data mining

Academic contact during examination: Tárik Saleh Salem, phone: 948 58 617

Examination date: 08-08-2019

Examination time (from-to): 09.00-13.00

Permitted examination support material: D: No tools allowed except an approved simple calculator

Other information:

Students will find the examination results in Studentweb. Please contact the department if you have questions about your results. The Examinations Office will not be able to answer this.

¹ Attribute Type (3 marks)

Which type of attribute is Fahrenheit temperature?

Select one alternative:

	4.5
יע	atio
Ra	7 U.V.

Interval

Nominal

Ordinal

Maximum marks: 3

² Dimensionality reduction (3 marks)

What are the purposes of dimensionality reduction?

Select one or more alternatives:

- Allow data to be more easily visualized
- May help to eliminate irrelevant features or reduce noise
- Remove outliers in data
- Reduce amount of time and memory required by data mining algorithms

3

Jaccard	coefficient	(3 marks)
---------	-------------	-----------

There are two bit vectors $ {f p}$ and $ {f q}$: ${f p} = [1110100011]$ ${f q} = [1101001011]$
What is the Jaccard coefficient for the bit vectors ${f p}$ and ${f q}$? Write your answer here
Note: the answer is a real-valued number.
Maximum marks: 3
Cosine similarity (3 marks)
There are two vectors $ {f p} $ and $ {f q} $:
$\mathbf{p} = [3,1,5]$
$\mathbf{q} = [6,7,2]$
What is the cosine similarity between ${f p}$ and ${f q}$? Write your answer here
Note the answer is a real-valued number (up to two decimals).

⁵ Modeling (10 marks)

Design the data warehouse for an electronics company. The data warehouse has to allow to analyze the company's situation at least with respect to the Product, Customers and Time. Moreover, the company needs to analyze:

- the product with respect to its name and type (office, home, accessories, wearable, photography, networking, etc.);
- the customers with respect to their spatial location, by considering at least cities, regions and states.

The company is interested in learning at least the quantity, income and discount of its sales.

One should be able to perform the following example analysis against the data warehouse:

- Total quantity for each year.
- Total quantity for each month.
- Average income for every day for each product type.

Create a star schema for the described case and define a concept hierarchy for each dimension.

Note: You have two options for drawing: 1) use a separate paper that will be scanned after the exam, 2) use the "Insert Drawing" tool in the toolbar.

Fill in your answer here

Format	- B	<u> Π</u> × ₂ × ² <u>T</u> _x [: Ω = 🖋 Σ 🔀	
				Words: 0

⁶ OLAP (10 marks)

Given a cube with dimensions:

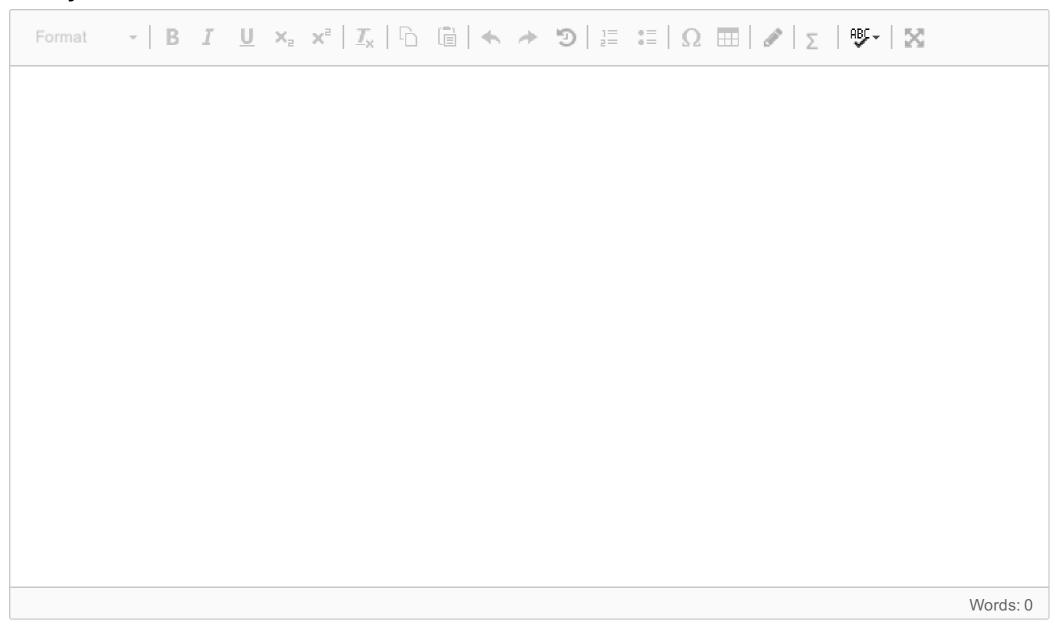
- Time(Day-Month-Year)
- Item(ItemName-Brand)
- Location(Street-City-ProvinceOrState-Country)

Assume the following materialized cuboids:

- {Month, ItemName, City}
- {Month, Brand, Country}
- {Year, Brand, ProvinceOrState}
- {ItemName, City} where year = 2016

Given the following OLAP query: {Brand, City} with condition Month = June 2016, which cuboid(s) should be used? Explain your answer below.

Fill in your answer here



⁷ Apriori Algorithm (15 marks)

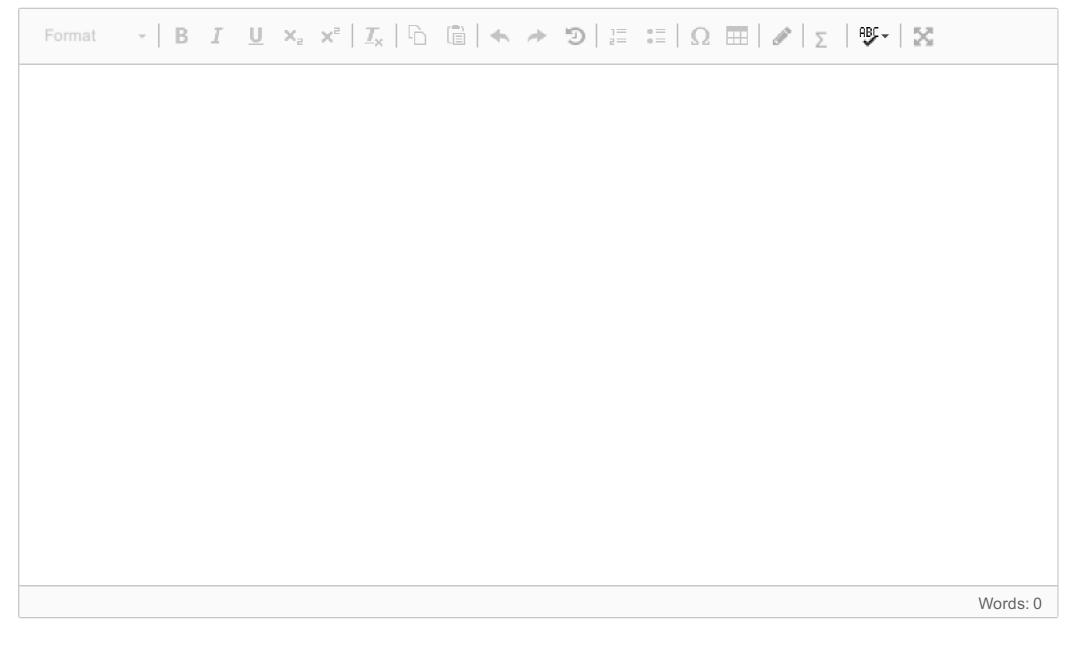
Assume the market basket data below. Use the Apriori algorithm to find all frequent itemsets with minimum support 33.33% (i.e. minimum support count is 2).

Transaction ID	Items	
T1	a, b	
T2	a, b, c	
Т3	a, d, e	
T4	a, d, e, f	
T5	c, e	
Т6	d, e	

- 1. Show how the **frequent itemsets** are generated.
- 2. $\{a,d,e\}$ is one of the frequent itemsets. Find all association rules based on this set, given confidence threshold 60% (it is not necessary to use Apriori to find the association rules, but show how confidence is calculated for each of the candidate rules that are based on $\{a,d,e\}$).

Note: you have two options for drawing: 1) use a separate paper that will be scanned after the exam, 2) use the "Insert Drawing" tool in the toolbar.

Fill in your answer here



8 FP-growth Algorithm (15 marks)

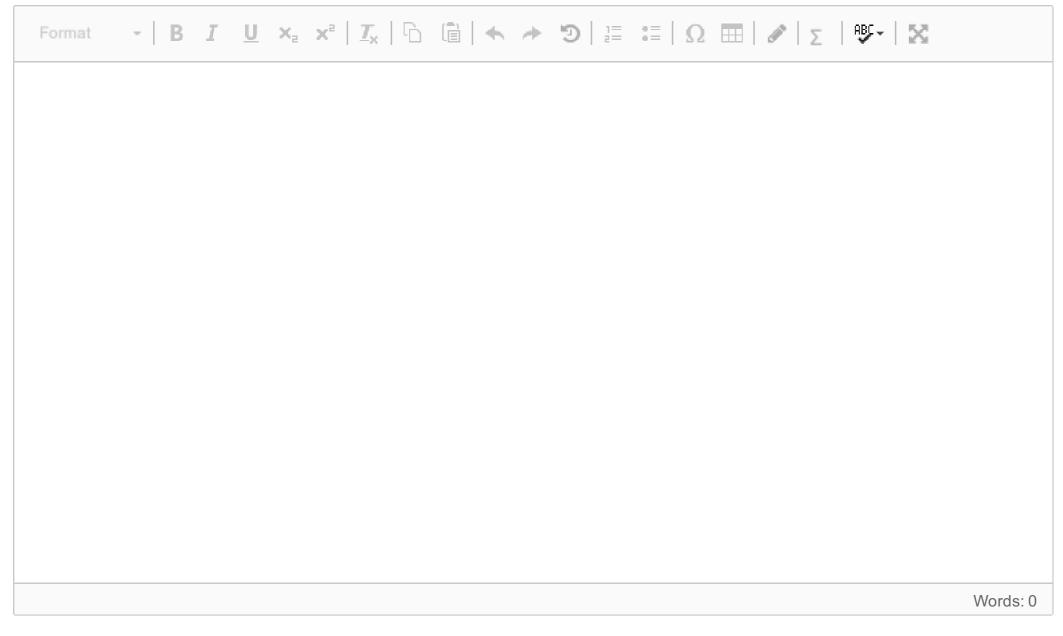
Assume the market basket data below. You are now going to use the FP-growth algorithm in order to find all frequent itemsets with minimum support of 22% (i.e., minimum support count is 2).

Transaction ID	Items	
T1	D, E, H	
T2	D, H	
Т3	D, F, G, H	
T4	B, G, H	
T5	C, G, H	
Т6	D, G, I	
Т7	D, F, G, H	
Т8	G, H, J	
Т9	A, D, F, G, H	

- 1) Construct a FP tree based on the dataset.
- 2) Find frequent itemsets using the FP-growth algorithm. Use table notation with the following columns in order to show the result:
 - Item
 - Conditional pattern base
 - Conditional FP-tree
 - Frequent itemsets

Note: you have two options for drawing: 1) use a separate paper that will be scanned after the exam, 2) use the "Insert Drawing" tool in the toolbar.

Fill in your answer here



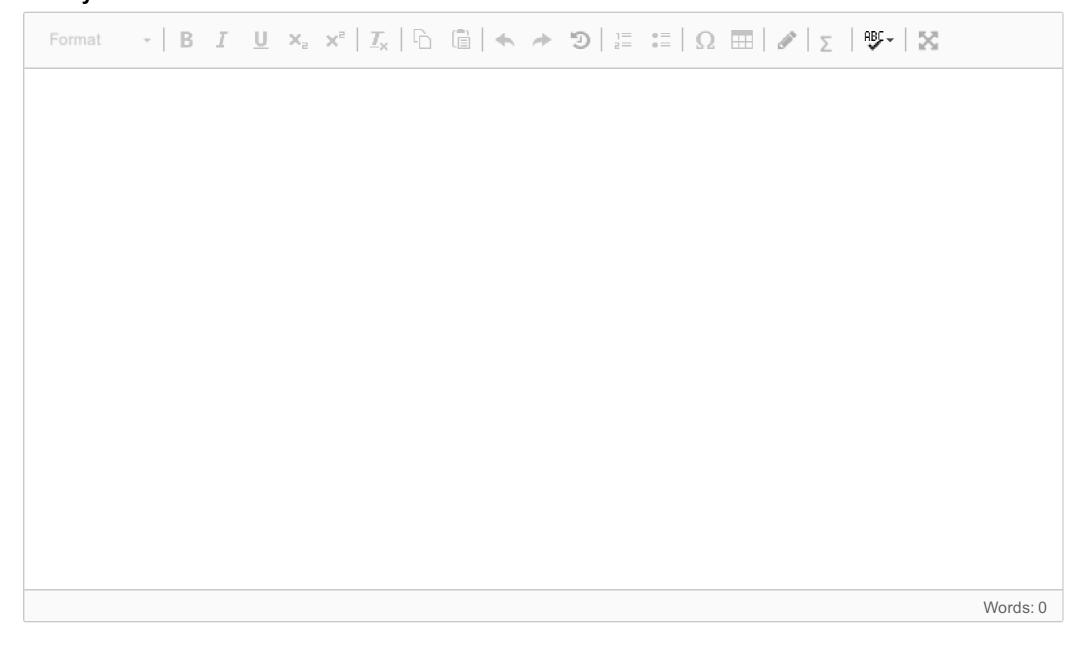
9 K-means Clustering (15 marks)

Do three iterations of the Lloyd's algorithm for K-means clustering on the 2-dimensional data below, using Euclidean distances. Use K=2 clusters and the initial prototype vectors (i.e. mean vectors) $\mathbf{m}_1=(3.0,5,0)$, $\mathbf{m}_2=(2.0,1.0)$. Write down calculation procedure and the cluster memberships as well as mean vectors after each iteration. Draw the data points, cluster means and cluster boundary after each iteration.

t	$\mathbf{x}^{(t)}$
1	(5.0, 2.0)
2	(4.0,1.0)
3	(2.0,4.0)
4	(1.0,3.0)
5	(0.0, 4.0)

Note: you have two options for drawing: 1) use a separate paper, that will be scanned after the exam, 2) use the "Insert Drawing" tool in the toolbar.

Fill in your answer here



¹⁰ Hierarchical clustering pros and cons (3 marks)

When does single linkage hierarchical clustering probably not perform well? **Select one or more alternatives**:

- Data contains outliers.
- Clusters have non-elliptical shapes
- Data is high-dimensional.
- There are large amounts of data points

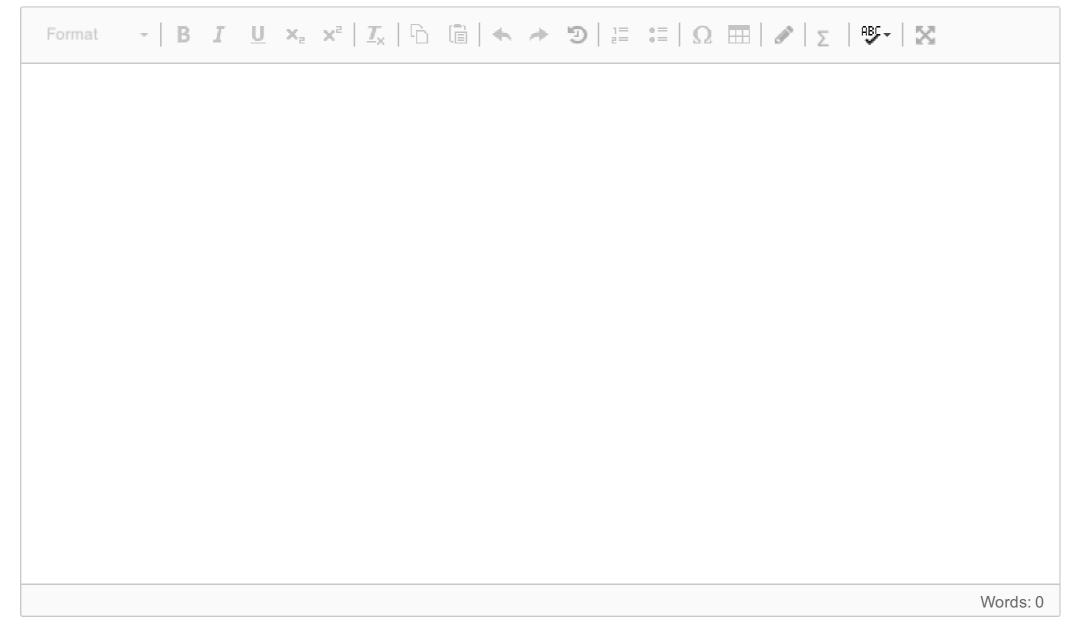
Maximum marks: 3

¹¹ Cross validation (5 marks)

Explain cross validation and what this technique is used for.

Note: if needed, you have two options for drawing: 1) use a separate paper, that will be scanned after the exam, 2) use the "Insert Drawing" tool in the toolbar.

Fill in your answer here



Maximum marks: 5

12 Decision tree (15 marks)

You are going to predict whether mushrooms are edible. You have the following data:

Example	Smooth	Spotted	Smelly	NotHeavy	Edible
А	1	1	0	1	0
В	1	0	0	0	0
С	1	0	1	1	0
				 	

TDT4300 Spring 2019 (August exam)

,					
D	0	0	1	0	0
E	0	1	1	1	0
F	1	0	1	0	1
G	0	0	0	1	1
Н	0	1	0	1	1
U	0	0	1	1	?
V	1	1	1	0	?
W	1	0	1	1	?

For mushrooms A through H, you know whether it is edible (1) or not edible (0), but you do not know about U through W.

You should use decision tree as a classification method. You will use the examples A through H as the training data. To decide the best split, you need to use ${f Entropy}$ for a node ${m t}$, given by

 $ext{Entropy}(t) = -\sum_j p(j|t) \log_2 p(j|t)$, where p(j|t) is the probability for class j given node t (i.e. the portion of class j in the node t). For each split, the "information gain" is defined by

$$ext{GAIN} = ext{Entropy}(p) - \left(\sum_{i=1}^k rac{n_i}{n} ext{Entropy}(i)
ight)$$
 , where n_i is the number of element in node i and n

is the total number of elements in the parent node p.

For Tasks 1 and 2, consider only mushrooms A through H. Tasks:

- 1. Which attribute should you choose as the root of a decision tree? Justify your choice by calculating the information gains of the attributes.
- 2. Build an ID3 decision tree to classify mushrooms as edible or not.
- 3. Classify mushroom U, V, and W using the decision tree to be edible or not edible.

Note: if needed, you have two options for drawing: 1) use a separate paper that will be scanned after the exam, 2) use the "Insert Drawing" tool in the toolbar.

Fill in your answer here

