

Exam

Uncertain Data Management Université Paris-Saclay, M2 Data&Knowledge

February 6th, 2018

This is the final exam for the Uncertain Data Management class, which will determine your grade for this class in conjunction with the class project. The duration of the exam is 2 hours. The exam consists of 4 independent exercises: exercises 1, 2, and 4 are about uncertain data management (Antoine's part), and exercise 3 is about social data management (Silviu's part). **You must write your answer to exercise 3 on a separate sheet of paper.**

Write your name clearly on the top right of every sheet used for your exam answers, and number every page.

You are allowed up to two A4 sheets of personal notes (i.e., four page sides), printed or written by hand, with font size of 10 points at most. If you use such personal notes, you must hand them in along with your answers. You may not use any other written material.

The exam is strictly personal: any communication or influence between students, or use of outside help, is prohibited. No electronic devices such as calculators, computers, or mobile phones, are permitted. Any violation of the rules may result in a grade of 0 and/or disciplinary action.

No additional explanations will be given during the exam, and no questions will be answered. If you think you have found an error in the problem statement, you should report on your answer sheet what you believe to be the error, and how you chose to interpret the intent of the question to recover from the alleged error.

Exercise 1: From c-tables to Boolean c-tables (10 points)

Recall that a *c-table* is a relational table where tuples carry a *condition*. The *condition* is a Boolean formula on equalities and inequalities that involve constant values and variables (called *named nulls*). In this exercise, we consider c-tables where *nulls may only occur in conditions*, i.e., they do not occur elsewhere in the table.

Recall that a *valuation* of a c-table T is a function that maps each named null of T to a value. Given a valuation ν of T , we can apply it to T by replacing each named null with its value in ν , evaluating the equalities and inequalities, evaluating the Boolean formulae, and keeping the tuples of T whose condition evaluates to true. The result of this process is a relational table, called a *possible world* of T , and we say that it is *obtained* from ν . Remember that a given possible world may be obtained by several different valuations.

Consider the following c-table, which describes upcoming events in the distant future in some major world cities:

Events		
city	year	
London	2042	$\text{NULL}_1 = \text{“GB”} \vee \text{NULL}_2 = \text{“Europe”}$
Mexico	2042	$\text{NULL}_1 = \text{“MX”}$
Singapore	2044	$\text{NULL}_1 = \text{“SG”} \wedge \text{NULL}_2 = \text{“Asia”}$

One example of a valuation is the function mapping NULL_1 to “GB” and NULL_2 to “Africa”. The possible world obtained by this valuation is the following relational table:

city	year
London	2042

Question 1 (0.5 point). Consider the valuation mapping NULL_1 to “MX” and NULL_2 to “Europe”. Write down the possible world obtained by this valuation. (No justification is expected.)

Answer. The condition of the first and second tuples evaluate to true, and the condition of the third tuple evaluates to false, so we obtain:

city	year
London	2042
Mexico	2042

Question 2 (1 point). Is there a valuation that yields the following possible world? If there is one, give an example valuation; if there is none, explain why.

city	year
London	2042
Singapore	2044

Answer. There is no such valuation. Indeed, to obtain the “Singapore” tuple, we must have $\text{NULL}_1 = \text{“SG”}$ and $\text{NULL}_2 = \text{“Asia”}$, and then the condition of the “London” tuple is necessarily false.

Question 3 (1 point). How many different possible worlds are there? (No justification is expected.)

Answer. As in the previous question, it is clear that if we have the “Singapore” tuple in a possible world, then it is the only such tuple; and indeed there is a possible world containing this tuple. So there is one possible world with only the “Singapore” tuple, and the other possible worlds do not contain this tuple. Now, restricting our attention to the first two rows of the table, it is clearly possible to obtain no tuple, any one of the two tuples, or both of them, i.e., four possibilities. So there is a total of five possible worlds.

Question 4 (0.5 point). Consider the relational algebra query $Q := \Pi_{\text{year}}(\text{Events})$ that projects the table `Events` on the `year` attribute. Write down the result of evaluating this query on the possible world obtained at question 1.

Answer. We obtain:

year
2042

Question 5 (1 point). Is the following table a possible result of the query Q on table `Events`?

year
2042
2044

Answer. It is not, because the only way to have “2044” in the query result is to have the “Singapore” tuple, and we know from Question 2 that in this case we do not have any other tuple.

Question 6 (2 points). Construct a c-table T_6 with nulls appearing only in conditions that represents the result of evaluating the query Q on the c-table `Events`. In other words, for any valuation ν of `NULL1` and `NULL2`, the possible world of T_6 obtained for ν should be the result of evaluating Q on the possible world of `Events` obtained for ν . (No justification is expected.)

Answer. We can take, for instance:

T_6	
year	
2042	(<code>NULL₁ = “GB”</code> \vee <code>NULL₂ = “Europe”</code>) \vee (<code>NULL₁ = “MX”</code>)
2044	<code>NULL₁ = “SG”</code> \wedge <code>NULL₂ = “Asia”</code>

Question 7 (3 points). Recall that a *Boolean c-table* is a c-table whose nulls can only be replaced by two values, `True` and `False`; in other words, we only consider valuations whose domain is $\{\text{True}, \text{False}\}$, unlike the c-tables considered so far where nulls can be replaced by arbitrary values.

Using the construction shown in class, construct a Boolean c-table T_7 whose possible worlds are exactly the possible worlds of the `Events` c-table. Explain briefly each step of your construction, and try to simplify the conditions that occur in T_7 .

Answer. We write down the list of possible worlds, we number them in binary, and we write down the corresponding Boolean condition on three variables `NULL1`, `NULL2`, and `NULL3`:

000	001	010	011	100
city year	city year	city year	city year	city year
	London 2042	Mexico 2042	London 2042 Mexico 2042	Singapore 2044
$\text{NULL}_1 = \text{False}$	$\text{NULL}_1 = \text{False}$	$\text{NULL}_1 = \text{False}$	$\text{NULL}_1 = \text{False}$	$\text{NULL}_1 = \text{True}$
$\wedge \text{NULL}_2 = \text{False}$	$\wedge \text{NULL}_2 = \text{False}$	$\wedge \text{NULL}_2 = \text{True}$	$\wedge \text{NULL}_2 = \text{True}$	$\wedge \text{NULL}_2 = \text{False}$
$\wedge \text{NULL}_3 = \text{False}$	$\wedge \text{NULL}_3 = \text{True}$	$\wedge \text{NULL}_3 = \text{False}$	$\wedge \text{NULL}_3 = \text{True}$	$\wedge \text{NULL}_3 = \text{False}$

We now construct the table, annotating each row with the condition corresponding to the possible worlds where it occurs, and simplifying the conditions. (As it turns out, they simplify quite well.)

city	year	
London	2042	$\text{NULL}_1 = \text{False} \wedge \text{NULL}_3 = \text{True}$
Mexico	2042	$\text{NULL}_1 = \text{False} \wedge \text{NULL}_2 = \text{True}$
Singapore	2044	$\text{NULL}_1 = \text{True}$

Question 8 (1 point). Construct a Boolean c-table T_8 that represents the result of evaluating the query Q on the c-table Events. (You do not need to use a minimal number of variables.)

Answer. One possibility is to apply the construction of Question 7 again to T_6 , in which case the result will use two variables. Another possibility is to evaluate the query Q directly on T_7 . With the latter method, we obtain the following Boolean c-table, using three variables:

city	year	
2042	$\text{NULL}_1 = \text{False} \wedge (\text{NULL}_3 = \text{True} \vee \text{NULL}_2 = \text{True})$	
2044	$\text{NULL}_1 = \text{True}$	

Exercise 2: Probabilistic modeling (5 points)

You are part of a secret service which monitors which countries buy and sell which types of radioactive materials. As you do not expect to be sure of which country is doing what, you would like to represent the information as two TID tables, Buy and Sell, that would look like the following, with some probabilities p_1, p_2, p_3, p_4 :

Buy			Sell		
country	material		country	material	
Syldavia	uranium	p_1	Borduria	uranium	p_4
Syldavia	plutonium	p_2			
Borduria	uranium	p_3			

You do not know what the probabilities p_1, p_2, p_3, p_4 should be. However, your field agents have been able to figure out the probability of some query results. Your job is to choose the right probabilities for p_1, p_2, p_3, p_4 so that the query results have the right probability.

For instance, consider the Boolean query $Q_0 := \Pi_{\emptyset}(\sigma_{\mathbf{country}=\text{“Syldavia”} \wedge \mathbf{material}=\text{“uranium”}}(\mathbf{Buy}))$. This query asks whether Syldavia is buying uranium. Your field agents have determined that this query has 50% probability of being true. To represent accurately this information in your TID tables, you will simply set $p_1 := 0.5$. This ensures that Q_0 has 50% probability of being true. In the rest of the exercise, we will determine p_2, p_3, p_4 in a similar way.

Question 1 (1 point). A new report from your field agents indicates that the Boolean query $Q_1 := \Pi_{\emptyset}(\sigma_{\mathbf{country}=\text{“Syldavia”}}(\mathbf{Buy}))$ has 80% probability of being true. What is the value of p_2 ? (Remember that $p_1 = 0.5$.)

Answer. We know that the probability of this query is $1 - (1 - p_1) \times (1 - p_2)$. So we have the following equation:

$$1 - (1 - p_1) \times (1 - p_2) = 0.8$$

By elementary arithmetic we then have:

$$p_2 = 1 - \frac{1 - 0.8}{1 - p_1}$$

Substituting the value of p_1 and evaluating, we obtain:

$$p_2 = 0.6$$

Question 2 (1 point). Your field agents now report that there is probability $\frac{5}{8}$ that Borduria is transacting uranium, i.e., it is either buying or selling it (or both). Write down the relational algebra query Q_2 corresponding to this field report.

Answer. The Boolean query is:

$$Q_2 := \Pi_{\emptyset}(\sigma_{\mathbf{country}=\text{“Borduria”} \wedge \mathbf{material}=\text{“uranium”}}(\mathbf{Buy} \cup \mathbf{Sell}))$$

Question 3 (1 point). Is the field report of Question 2 sufficient to determine the probabilities p_3 and p_4 ? Explain why, or why not.

Answer. The report tells us the following:

$$1 - (1 - p_3) \times (1 - p_4) = \frac{5}{8}$$

This is not sufficient to deduce the values of p_3 and p_4 , because the system is under-specified: there are two variables but only one equation.

For instance, here are two possible choices (there are infinitely many):

- $p_3 = 0$ and $p_4 = \frac{5}{8}$;
- $p_3 = \frac{5}{8}$ and $p_4 = 0$;

Question 4 (2 points). A new field report arrives: there is, in addition, a $\frac{1}{8}$ probability that Borduria is both selling and buying uranium. Give all possible solutions for p_3 and p_4 .

Hint: remember that quadratic equations of the form $ax^2 + bx + c = 0$ can be solved with the quadratic formula: compute $\Delta := b^2 - 4ac$, and then there are three cases:

- If $\Delta < 0$ then there are no solutions;
- If $\Delta > 0$ there are two solutions: $x_1 = \frac{-b+\sqrt{\Delta}}{2a}$ and $x_2 = \frac{-b-\sqrt{\Delta}}{2a}$;
- If $\Delta = 0$ there is one solution $x_1 = x_2$ given by the formulas above.

Answer. The new report tells us that:

$$p_3 \times p_4 = \frac{1}{8}$$

This clearly implies that $p_3 \neq 0$ and $p_4 \neq 0$. So let us rewrite this as follows:

$$p_4 = \frac{1}{8p_3}$$

And let us substitute this in the equation of question 2:

$$1 - (1 - p_3) \times \left(1 - \frac{1}{8p_3}\right) = \frac{5}{8}$$

We rewrite this and obtain:

$$(1 - p_3) \times \left(1 - \frac{1}{8p_3}\right) = \frac{3}{8}$$

Multiplying each side by 8 and expanding, we obtain:

$$8 - \frac{1}{p_3} - 8p_3 + 1 = 3$$

We reorder and multiply by p_3 (which is non-zero) and obtain the quadratic equation:

$$8p_3^2 - 6p_3 + 1 = 0$$

We compute $\Delta = 4$, so $\sqrt{\Delta} = 2$, and the solutions are: $p_3 = \frac{1}{4}$ and $p_3 = \frac{1}{2}$. The corresponding values for p_4 are respectively $\frac{1}{2}$ and $\frac{1}{4}$.

Exercise 3: Uncertain Graphs (3 points)

Consider an uncertain graph $\mathcal{G} = (V, E)$ as defined in class, where each edge $e \in E$ has attached a probability $p_e \in (0, 1)$, encoding the probability of e existing in the graph, and where all p_e 's are considered independent.

Question 1 (1 point). How many *possible worlds* does the graph have (give the formula)? State which graph variables it depends on (number of nodes, number of edges, etc.). If we were to store this graph in a probabilistic relational database, which of the block-independent or tuple-independent models would we need?

Answer. The graph has $2^{|E|}$ possible worlds, i.e., all combinations of edges existing or not. Since the probabilities are only on edges and are independent, the formula depends only on the number of edges in the uncertain graph.

Each edge can be considered a tuple in a binary relation (the graph). Since the probabilities on edges are all independent, a tuple-independent model is sufficient.

Question 2 (2 points). Consider now an uncertain graph where each edge e is annotated with a Boolean formula on named nulls (the formulas are constructed like in Exercise 1). We assume that the graph has at least 1 and at most $n = |E|$ named nulls. Give the minimal number and the maximal number of possible worlds, and state on which variable(s) it depends on.

Answer. The minimal number of possible worlds is 1, corresponding to the case where we have only one named null, but we annotate each edge using either **True** or **False**; the maximal number is 2^n , corresponding to the case where we have n named nulls, and each named null occurs on one edge only.

The formula for computing the number of possible worlds depends on the number of named nulls we have in our graph.

Exercise 4: From TID to BID (2 points or more)

Note: this exercise is substantially harder than the previous ones. It is recommended not to attempt it until you have solved the other exercises.

Recall that a relational algebra query is called *monotone* if it does not use the difference operator.

Question 1. Prove that, for every monotone relational algebra query Q and TID instance I , if the probabilistic instance $Q(I)$ cannot be represented by a TID instance, then it cannot be represented by a BID instance either.

Answer. It can be shown by an immediate induction on monotone relational algebra operators that any monotone relational algebra query Q is also monotone in the following sense: for any relational instances $I_1 \subseteq I_2$, we have $Q(I_1) \subseteq Q(I_2)$. Now, for any TID instance I , letting I_+ be the possible world of I containing all tuples of I , it is obvious that we have $I' \subseteq I_+$ for any possible world I' of I . In other words, I_+ is the maximal possible world. Thanks to the monotonicity of Q , this implies that, for any possible world J of the probabilistic instance $Q(I)$, letting $J_+ := Q(I_+)$, we have $J \subseteq J_+$. In other words, $Q(I)$ also has a maximal possible world J_+ .

Assume now that $Q(I)$ can be represented by a BID instance B . We show that there are no two tuples $t_1 \neq t_2$ of B in the same block. Indeed, assuming to the contrary that such tuples t_1 and t_2 exist, consider any possible world B_1 of B that includes t_1 ,

and the possible world B_2 obtained from B by replacing t_1 by t_2 ; it is still a possible world of B . Now, the maximal possible world J_+ of B should be such that $B_1 \subseteq J_+$ and $B_2 \subseteq J_+$, hence, J_+ should contain both t_1 and t_2 , which is impossible because they were assumed to be in the same block. Hence, B consists only of singleton blocks. This means that B can be represented by a TID instance. So indeed we have shown that if $Q(I)$ can be represented by a BID instance then it can be represented by a TID instance, which is the desired result.

Question 2. Give an example of a relational algebra query Q and TID instance I such that the probabilistic instance $Q(I)$ cannot be represented by a TID instance but can be represented by a BID instance.

Answer. We consider the TID instance:

R		
a	b	
x	y	0.5
x	z	1
y	z	1

We let the query Q be:

$$Q = \sigma_{\mathbf{a}=x}(\mathbf{R}) - \Pi_{\mathbf{a}, \mathbf{b}}((\rho_{\mathbf{b} \rightarrow \mathbf{c}}(\mathbf{R}) \bowtie \rho_{\mathbf{a} \rightarrow \mathbf{c}}(\mathbf{R})))$$

There are two possible worlds of \mathbf{R} :

R		R	
a	b	a	b
x	y		
x	z	x	z
y	z	y	z
0.5		0.5	

The respective query results are:

R		R	
a	b	a	b
x	y		
		x	z
0.5		0.5	

It is clear that this probabilistic instance cannot be represented by a TID instance, because it has no maximal possible world. However, it can be represented by the following BID instance:

R		
<u>a</u>	b	
<i>x</i>	<i>y</i>	0.5
<i>x</i>	<i>z</i>	0.5