



**BSc, BEng and MEng Degrees Examination 2020—21**

DEPARTMENT OF COMPUTER SCIENCE

**Data Analysis and Management**

**Time Allowed:** TWENTY-FOUR hours

**Time Recommended:** 2 hours 30 minutes

**Word limit:** NA

**Allocation of Marks:**

**This assessment is out of 100 marks:**

- Q1 (14 marks) | Q2 (14 marks) | Q3 (22 marks) | Q4 (9 marks)
- Q5 (11 marks) | Q6 (13 marks) | Q7 (17 marks)

**Instructions:**

- Answer **all** questions.
- For some of the questions you will need: Jupyter notebook for Python 3 and SQLite DB Browser (version 3.12); and some files, which are available on VLE (Assessment -- 24\_Hour\_Exam).
- Submit your answers to [the Department's Teaching Portal](#) as a zipped directory containing four files:
  - Data2.pdf (This file must include all answers (such as text, SQL commands, code, figures) for all questions)
  - Data2\_Q5.ipynb (This file must include the code as part of your answers for Q5)
  - Data2\_Q6.ipynb (This file must include the code as part of your answers for Q6)
  - Data2\_Q7.ipynb (This file must include the code as part of your answers for Q7)

**Note:** The 'ipynb' files will be used to test and run your solution

If a question is unclear, answer the question as best you can, and note the assumptions you have made to allow you to proceed.

### **A note on Academic Integrity**

We are treating this online examination as a time-limited open assessment, and you are therefore permitted to refer to written and online materials to aid you in your answers.

However, you must ensure that the work you submit is entirely your own, and for the whole time the assessment is live you must not:

- communicate with departmental staff on the topic of the assessment
- communicate with other students on the topic of this assessment.
- seek assistance with the assignment from the academic and/or disability support services, such as the Writing and Language Skills Centre, Maths Skills Centre and/or Disability Services. (The only exception to this will be for those students who have been recommended an exam support worker in a Student Support Plan. If this applies to you, you are advised to contact Disability Services as soon as possible to discuss the necessary arrangements.)
- seek advice or contribution from any third party, including proofreaders, friends, or family members.

We expect, and trust, that all our students will seek to maintain the integrity of the assessment, and of their award, through ensuring that these instructions are strictly followed. Failure to adhere to these requirements will be considered a breach of the Academic Misconduct regulations, where the offences of plagiarism, breach/cheating, collusion and commissioning are relevant - [see AM.1.2.1](#)” (Note this supersedes section 7.3 of the Guide to Assessment).

## Q1. [Total: 14 Marks]

Go through the following case study, and then answer the questions that follow:

ABC Ltd. is an online supplier who requires a database system for processing sales orders. Customer accounts are created when they place their first order, with many orders can be placed. According to the company's policy, each order can only be from one customer.

Each order placed by a customer can consist of one or many items, but a given item refers to exactly one order. Each order raises an invoice, and an invoice belongs to one order. A product could be part of one or many items, but each item refers to only one product.

Each order is then processed by a single employee on the sales accounts team, who prepares the shipment of the goods to the customer. An employee can be scheduled for processing many separate orders and preparing many shipments each day, but may not have any schedules for both tasks at all.

Items in each order can be packaged together into one shipment, or packaged separately and sent by several shipments, if required. A shipment can only belong to one item.

A customer may not have any complaints to make, but they can make as many they want. Each complaint comes from one customer, which is then dealt by one employee. Not all employees will handle the complaints, and sometimes an employee may have to deal with many complaints.

### Task:

- (i) [6 marks] You are required to draw the Entity Relationship Model (ERM) for the proposed system, clearly showing all entities, relationships among participating entities, multiplicity constraints and cardinalities.
- (ii) [6 marks] From (i):
  - (a) [2 marks] What are the entity pairs involved in the relationship for raising an order?
  - (b) [1 mark] From (a) what is the cardinality of this relationship?
  - (c) [1 mark] From (a) state with explanation which entity is the parent and which entity is the child? If there is not a parent and a child, you still need to provide a justification.
  - (d) [2 marks] From (a) when mapping the ERM to tables, how will the relationship between these two entities be represented?
- (iii) [2 marks] Are there any weak entities? If yes, list them.

## Q2. [Total: 14 Marks]

The following sample table (refer to Figure 1) details students who are being enrolled on one or more modules, their grades, and the tutors who are teaching these modules. The composite key will be (sID, moduleNo).

sID	moduleNo	moduleName	grade	sName	tutorID	tutorName	tutorOffice
S101	M1001	AI	A	TC	T1	ZD	R121
S102	M1002	HCI	B	SR	T2	AS	R220
S103	M1003	DI	C	SG	T3	BN	R010
S104	M1003	DI	B	KP	T3	BN	R010
S101	M1004	SE	A	TC	T4	CF	R111
S103	M1002	HCI	A	SG	T2	AS	R220

Figure 1: Students, grades, and tutors

### Task:

- [2 marks] Which columns in the table contain redundant data?
- [4 marks] List all functional dependencies which exist in the table.
- [3 marks] Of the functional dependencies that **actually exist** which, if any, violate the rules of the 2NF?
- [2 marks] How could you resolve the issues in (iii), whilst maintaining the required relationship?
- [3 marks] Having answered (iv), show all the tables that are required, clearly showing the primary keys, foreign keys, and any composite keys.

For Question 3 you need SQLite 3.12 DB Browser

## Q3. [Total: 22 Marks]

DVD XYZ is a simple database designed to support the business activities of a DVD rental library. The multiplicity constraints are as follows:

- A distribution center has one or many members of staff. Each staff member works at one distribution center.
- A staff member manages zero or one distribution center. Each distribution center is managed by staff member.
- A DVD has one or many copies. Each DVD copy belongs to one DVD.
- A supplier supplies one or many DVDs. Each DVD is supplied by one supplier.
- A distribution center has one or many DVD copies. Each DVD copy is at one distribution center.

With regards to the above information, a set of tables has been designed, consisting of the following five tables.

### A\_DistributionCenter Table

dCenterNo	dStreet	dCity	dState	dZipCode	staffNo
B001	8 Jefferson Way	Portland	OR	97201	S1500
B002	City Center Plazza	Seatle	WA	98122	S0010
B003	14 -8th Avenue	New York	NY	10012	S0415
B004	16 -14th Avenue	Seatle	WA	98128	S2250

### B\_Staff Table

staffNo	name	position	salary	dCenterNo
S0003	Sally Adams	Snr Assistant	30000	B001
S0010	Mary Martinez	Manager	50000	B002
S0415	Art Peters	Manager	41000	B003
S1500	Tom Daniels	Manager	46000	B001
S2250	Sally Stern	Manager	48000	B004
S2350	Robert Chin	Supervisor	32000	B002

### C\_Supplier Table

supplierNo	name	address	telNo	status
S01	Universal Home Videos	100 Universal City Plaza	8188666000	OK
S02	MGM Home Videos	2500 Broadway St, Santa Monica, CA, 90404	8189002000	OK
S03	Buena Vista Pictures	1100 Santa Monica Bivd, CA, 90041	3208406500	OK
S04	Paramount Pictures	5555 Melrose Avenue, Hollywood, CA, 90038	3238621130	OK
S05	20th Century Fox Home Video	900 Center Plaza, Beverly Hills, CA, 90213	6007772300	OK

### D\_DVD Table

catalogNo	title	genre	rating	supplierNo
207132	Casino Royale	Action	PG-13	S02
330553	Lord of the Rings III	Action	PG-13	S04
445624	Mission Impossible III	Action	PG-13	S03
634817	War of the Worlds	Sci-Fi	PG-13	S05
781132	Shrek 2	Children	PG	S03
902355	Harry Potter	Children	PG	S01

### E\_DVDCOPY Table

videoNo	available	catalogNo	dCenterNo
178643	False	634817	B001
199004	True	207132	B001
200900	True	330553	B002
210087	True	902355	B002
243431	True	634817	B002
245456	True	207132	B002
245457	True	207132	B002
317411	True	781132	B003

The tables are in the design phase, and they need to be implemented.

### Task:

Using **SQLite 3.12 DB Browser**, you are required to:

- (i) [3 marks] Give the SQL commands to create each table, clearly showing the primary and foreign keys.
- (ii) [3 marks] Give the SQL commands to populate each table with all the records.
- (iii) [16 marks] Having answered ((i)-(ii)), you are now required to write and show your commands for each of the following queries:
  - (a) [2 marks] Write a query so that the result lists the titles of all DVDs along with their length in descending order of length.
  - (b) [2 marks] Write a query so that the result lists the full details of all DVD copies (discs) held in the Portland and New York distribution centers.
  - (c) [3 marks] Write a query so that the result lists all distribution centers and their total staff salary costs but only where such costs are greater than 50,000.
  - (d) [3 marks] Write a query so that the result lists the total number of DVD copies (identified by video number) available for rental in each of the distribution centers, with rows displayed in descending order of copies available.
  - (e) [3 marks] Write a query so that the result lists all of the DVD copies supplied by 20th Century Fox Home Videos along with their titles and current availability for rental.
  - (f) [1 mark] From the results of your last query, what is the video number of the copy that is available?
  - (g) [2 marks] Sally Adams (staff number = S0003) has been given a raise of 1000 per year. She has also just got married and her surname has changed to Daniels. Write a query so that these changes are reflected in the database.

**Q4. [Total: 9 Marks]**

Below is a dataset (see Figure 2) on students and grades. Write an XML document for this data. You should clearly show the use of internal DTD to define the structure of the XML document.

		Maths	English	IT	
ID	Student	Test	Essay	Activity 1	Test 2
1643022	Pace, Camden G.	54%	C+	Pass	32%
1647021	Weber, Lucy X.	32%	C-	Fail	73%
1606023	Branch, Caesar A.	73%	C-	Pass	63%

Figure 2: Students and grades

For Question 5 you need Jupyter notebook for Python 3

**Q5. [Total: 11 Marks]**

Consider the given dataset (refer to **kmeans.csv**) on VLE (Assessment -- 24\_Hour\_Exam). Using **Jupyter notebook for Python 3** you are required to use the k-means algorithm to cluster this dataset into 3 groups. Specifically, you need to show:

- (i) [7 marks] The coordinates of the point that belong to each cluster.
- (ii) [4 marks] Graphically, the points that belong to each cluster, including the means. The means should be of a different shape and colour to those provided for the points.

**Include your code in your answer.**

For Question 6 you need Jupyter Notebook for Python 3

**Q6. [Total: 13 Marks]**

You will be using data generated from a survey conducted to explore the prevalence and impact of sleep problems on aspects of people's lives. Refer to the file (**survey\_responses.csv**) on VLE (Assessment -- 24\_Hour\_Exam). Below (refer to Figure 3) is a description of the data source and how the collected data has been coded.

Description of variable	Variable name	Coding instructions
Gender	gender	male; female
Rate quality of sleep	qualslp	1=very poor, 2=poor, 3=fair, 4=good, 5=very good, 6=excellent
Hours sleep/week nights	hourwnit	Hours sleep on average each weeknight
Hours sleep/week ends	hourwend	Hours sleep on average each weekend night

Figure 3: Codebook

With regards to the data, suppose we want to examine the following: “Is there a significant difference in quality of sleep reported by men and women?”

### Task

To answer this question:

- (i) [6 marks] Show the code (**written in Jupyter notebook for Python 3**) to justify whether a parametric or non-parametric test would be needed.
- (ii) [3 marks] Having answered (i), which statistical test will you apply? Justify your answer.
- (iii) [4 marks] From (ii), apply the appropriate statistical test and clearly explain your result. You need to show the code (**written in Jupyter notebook for Python 3**) that you have used to reach your answer.

**For Question 7 you need Jupyter Notebook for Python 3**

### Q7. [Total: 17 Marks]

Consider the training dataset (refer to **training\_data.csv** on VLE) (Assessment -- 24\_Hour\_Exam) which describes a set of objects using eight attributes, A1-A8. The dataset also lists whether each object is a ‘0’ (as negative) or ‘1’ (as positive) example of a certain, unnamed concept (see column ‘Output\_Class\_Label’ in the csv file). We want to build a logistic regression model in Python to try to predict this class.

### Task

**Using Jupiter Notebook for Python 3:**

- (i) [1 mark] Show the code to read the file ‘training\_data’.
- (ii) [1 mark] We do not want the ID column in our analysis. Show the code to drop this column.
- (iii) [5 marks] Show the code to clean the dataset and run the logistic regression. You should use 80% of the data for training and 20% for testing.



- (iv) [2 marks] Show the magnitudes of the coefficients. You must show your code. Which attribute corresponds to which weight must be clearly shown when running the code?
- (v) [1 mark] Which attribute is impacting the result the most?
- (vi) [2 marks] Show the code to display the confusion matrix. The confusion matrix should clearly show the following labels: 0 and 1, Actual (or True), and Predicted.
- (vii) [2 marks] Interpret the confusion matrix.
- (viii) [3 marks] From (vi) calculate the Sensitivity, Positive Predicted Value, and F1 score. Show your workings including the formulas.

**END OF PAPER**