

# COMP1314 Coursework

Philbert Poku (pkp1u24@soton.ac.uk)

Brief You have just been promoted to the position of “Data Manager” at the prestigious University of Northeasthampton. Unfortunately, their records are a mess, and it’s up to you to fix them.

## 1 Structured Data

The University of Northeasthampton currently stores all of their data in a pair of xml files. This would not be an ideal format at the best of times, but it is made worse by the way they have chosen to format their files.

## 1.1 Ex1

### Bash Script

```
1  #!/bin/bash
2
3  xml_file="$1"
4  csv_file="$2"
5
6  # Extract the start and end tags
7  start_tag=$(sed -n '3p' "$xml_file")
8  tagname=$(echo "$start_tag" | sed -E 's/^<([>]+)>$/\1/' | xargs)
9  # Extract and trim tag name
10 end_tag="</$tagname>"
11
12 # Extract headers and save to the CSV file
13 awk -v start_tag="$start_tag" -v end_tag="$end_tag" '
14     BEGIN { block_started = 0 }
15
16     $0 == start_tag { block_started = 1; next }
17     $0 == end_tag { exit }
18
19     block_started {
20         if (match($0, /<([^\>]+)(\/?)>/, matches)) {
21             tag = matches[1]
22             if (!(tag in seen)) {
23                 seen[tag] = 1
24                 headers = headers (headers ? "," : "") tag
25             }
26         }
27     }
28
29     END { print headers }
30 ' "$xml_file" > "$csv_file"
```

This code section begins with taking an xml file and csv file as positional arguments. We then extract the start tag and then dynamically create the end tag to know when each data block begins and ends. Firstly, we must extract the headers and we can do this using a for loop in awk where it matches what is in between the tags to something it has already seen before. If it has not been seen it is added to the string headers. It then prints the headers and adds it to the csv file.

## Bash Script

```
31 # Extract data from every line
32 grep -v '<?xml' "$xml_file" | \
33 awk '
34     # Process <faculty> blocks
35     /<faculty>/ && !/<\|/faculty>/ {
36         block_start=1
37         block=""
38         current_tag="faculty"
39         next
40     }
41     /<\|/faculty>/ && !/<faculty>/ && block_start && current_tag ==
42     ↪ "faculty" {
43         block_end=1
44     }
45     # Process <student> blocks
46     /<student>/ && !/<\|/student>/ {
47         block_start=1
48         block=""
49         current_tag="student"
50         gsub(/<student> */, "")
51     }
52     /<\|/student>/ && block_start && current_tag == "student" {
53         block_end=1
54     }
55     # Process completed blocks
56     block_start && block_end {
57         print block
58         block=""
59         block_start=block_end=0
60         next
61     }
62     block_start && current_tag == "faculty" {
63         line = gensub(/.*>([<]*)<./, "\\1", "g")
64         # Extract content between tags
65         block = block (block ? "," : "") line
66     }
```

## Bash Script

```
66 block_start && current_tag == "student" {
67     # Handle <address> tag - start
68     if ($0 ~ /<address>/) {
69         line = gensub(/<address>/, "\"", "g")
70         # Replace <address> with " for formatting purposes
71         getline next_line # Get the next line
72         gsub(/^[[:space:]]+|[[[:space:]]+$/, "", next_line)
73         next_line = gensub(/<\/address>/, "\"", "g", next_line)
74         # Replace </address> with " for formatting purposes
75         line = line " " next_line
76         # Concatenate the two lines
77         gsub(/[[[:space:]]+/, " ", line)
78         block = block (block ? "," : "") line
79         # Add the processed address to the block
80     }
81     # Treats null tags as empty fields
82     else if ($0 ~ /\>/) {
83         block = block (block ? "," : "") ""
84     }
85     # Extract content between tags for every other line
86     else {
87         line = gensub(/.*>([<]*)<./, "\\1", "g")
88         gsub(/[[[:space:]]+/, " ", line)
89         block = block (block ? "," : "") line
90         # Add the extracted data to the block
91     }
92 }
93 ' | \
94 # Remove leading whitespace and comma before the first character
95 sed 's/^ *, */,> >> "$csv_file" # Append the extracted data to the
96 ↪ CSV file
97
```

The next part of the code is the extraction process in which we divide it into processing the faculty csv and the students csv. When the opening or closing tag of faculty is encountered this means the start or end of a block. However, if faculty is encountered on the same line we extract between the faculty tags. If a line starts with the opening address tag, it spans two lines. This means that after the opening tag is saved to a variable and by getting the data on the next line they can be concatenated together. During that, the tags are replaced with speech marks for formatting purposes due to the address containing commas. For null tags, we continue onto the next line and for every other tag we just extract from between the tags. At the end remove the whitespace and the first comma for formatting purposes

and append to the csv file.

## 1.2 Ex2

### Bash Script

```
1  #!/bin/bash
2
3  csv_file=$1
4  txt_file=$2
5
6  # delete the header and sort the rest based on the first field
7  # write the sorted data to the text file
8  tail -n +2 "$csv_file" | awk -F ',' '{print $1}' | sort | uniq >
   ↪ "$txt_file"
```

This Bash script takes in a csv file and a text file as positional arguments and deletes the first line of the csv file. It does this by using the tail command to display the contents of the file and only starts from the second line effectively deleting the header. It then uses the cut command to extract the first field of the csv file which is the student name and pipes this to the sort command to sort it alphabetically. It then outputs the result to a new text file.

## 2 Relational Model

Now that the data can be more easily examined, the next task is to inspect it. students.csv and faculty.csv can be considered as two relations.

### 2.1 Ex3

We can express each students and faculty relation in the form of  $R(A, B, C, D, \dots)$  by expressing the elements of each csv header as attributes. The faculty relation can be expressed as: faculty(faculty, building, room, capacity). The students relation can be expressed as:

```
students(student_name,student_id,student_email,programme,year,address,
contact, module_id,module_name,module_leader,lecturer1,lecturer2,faculty,
building,room,exam_mark,coursework1,coursework2,coursework3)
```

### 2.2 Ex4

For each relation, in order to list the minimal set of Functional Dependencies in the form  $A \rightarrow B, C$ , we have to find which attributes uniquely determine other attributes. For the faculty relation we have 4 attributes and we can identify that a faculty and a building determines a room and a room determines the capacity. Therefore the minimal set of Functional Dependencies for faculty are:

```
building → faculty
building, room → faculty
building, room → capacity
```

Similarly, for the student relation we have 19 attributes and we can identify that the `student_id` determines all the attributes related to information about the student and the `module_id` likewise for the module. However both ids can give you information about the various coursework and the exam mark. Therefore the minimal set of functional dependencies for students are: bullet point

```
student_id → student_name
student_id → student_email
student_id → year
student_id → address
student_id → contact
module_id → module_name
module_id → module_leader
module_id → lecturer1
module_id → lecturer2
module_id → faculty
module_id → building
module_id → room
student_id, module_id → exam_mark
student_id, module_id → coursework1
student_id, module_id → coursework2
student_id, module_id → coursework3
```

## 2.3 Ex5

A key is a set of attributes for which no two tuples in a relation instance have the same values for all attributes of the key. Therefore, to find a candidate key of our relation we must discover the minimal superkeys of our relation as all possible candidate keys are subsets of the superkey. To do this we can use a closure algorithm, which makes it easier since we already know all of the minimal functional dependencies. From the faculty relation we can start with all the possible subsets that contain more than 2 attributes from (faculty, building, room, capacity) starting with {building,room}. The closure of {building,room} is {faculty, building, room, capacity} because from building and room you can derive both the capacity and the faculty, so {building,room} is a superkey. From {building} you can uniquely determine the faculty but you cannot uniquely determine the room therefore it is not a superkey. Overall, there is 1 superkey for the faculty relation {room, building} meaning there is only one candidate key.

From the students relation we realise that there are a lot of possible subsets  $2^{19}$  to be exact therefore having the minimal set of functional dependencies is useful here for the closure algorithm. From our functional dependencies we know {student\_id} can derive the student\_name, student\_email, programme, year, address and contact however this does not include module information or grade information. From {module\_id} we can derive the

module\_name, module\_leader, lecturer1, lecturer2, faculty, building and room which does not include student information or grade information. For the last functional dependency {student\_id, module\_id} we can derive all attributes related to student, module and grade information making it a candidate key. As it is a superkey, it must be the only candidate key.

## 2.4 Ex6

Now for each relation we need to identify a suitable primary key. A primary key is always one of the candidate keys therefore for the faculty relation we have a choice from only one candidate key which is {room, building}. {room, building} has all the attributes needed to uniquely identify a room and faculty. Thus {building, room} should be chosen as it has the least information required to identify a single row of attributes uniquely.

For the students relation we only have one choice and that is {student\_id, module\_id} as it has all the attributes needed to uniquely identify student information, module information and grade related information.

## 3 Normalisation

A clearer picture of the current state of the data has now been established, making it the right time to begin improving it.

### 3.1 Ex7

For the faculty relation the data is already in First Normal Form as the relation contains only atomic values which cannot be broken down any further single values. In addition there are no composite, multi-valued or nested values nor are there objections and arrays. There are also no repeating groups as every attribute is unique.

However, some attributes can be composed into new relations for the students relation. For example, lecturer1 and lecturer2 can be made into one attribute lecturers as it is a repeating group. We can then decompose this into a new relation cohort\_lecturers(student\_id, module\_id, lecturers) using our primary key student\_id, module\_id}. The same can be done for coursework1, coursework2, and coursework3 where we decompose it into a new attribute coursework. Using our primary key, we can create another new relation courseworks(student\_id, module\_id, exam\_mark, coursework) where we can add the exam\_mark to the relation as it depends on the coursework. Therefore our original students relation is now

```
students(student\_name, student\_id, student\_email, programme, year, address,
contact, module\_id, module\_name, module\_leader, faculty, building, room)
```

### 3.2 Ex8

To resolve any conflicts or irregularities in the data decomposition was used for example in the students relation for there not to be a redundancy lecturers attribute it was moved into a

new relation `cohort_lecturers(student_id, module_id, lecturers)` and to make sure that it does not have any non-atomic values. The same was done for the `courseworks` relation where the redundancy was handled by creating a new relation.

### 3.3 Ex9

A partial key dependency is where every non-key attribute (that is not part of any candidate key) is dependent on all attributes of all candidate keys. What this means for the `students` relation is that all the attributes in the functional dependencies: `student_id → student_name, student_email, programme, year, address, contact` and `module_id → module_name, module_leader, faculty, building, room` are partial key dependencies as these attributes depend only on `student_id` and `module_id` and these create the only candidate key.

For the `faculty` relation, the functional dependency `building → faculty` is a partial key dependency as `faculty` depends on `building` which is a proper subset of the primary key `{building, room}` and `faculty` does not appear on any candidate key.

### 3.4 Ex10

To achieve 2nd Normal Form for the `faculty` relation we can introduce more relations through decomposition. For example, the functional dependency `building → faculty` can be decomposed into new relations by copying the determinant which is the part of the key they are dependent on. A new relation `location(building, faculty)` can be created leaving the `faculty` relation as `faculty(building, room, capacity)`.

To achieve 2nd Normal Form for the `students` relation we can introduce more relations. For example, the functional dependencies for `student_id` and `module_id` can be decomposed into new relations by copying the determinant which is the part of the key they are dependent on. A new relation `student_info` can be created

```
student\_info(student\_name,student\_id,student\_email,programme,year,address,contact)
```

leaving the `students` relation as

```
students(student\_id, module\_id, module\_name, module\_leader, faculty, building, room)
```

The same can be done for `module_id` where a new relation `module_info` can be created

```
module\_info(module\_id,module\_name,module\_leader,faculty,building, room)
```

leaving the `students` relation as `students(student_id, module_id)` which is the primary key.

### 3.5 Ex11

To achieve 2nd Normal Form, the `faculty` relation can be decomposed based on its functional dependencies while preserving the appropriate primary key attributes `building` and `room`. Therefore this is why `faculty` was moved into the `location` relation. The `students` relation was also decomposed based on its functional dependencies while preserving the appropriate primary key attributes `student_info` and `module_info`. New relations were created where the attributes were partial key dependencies for example `student_info` and `module_info`.



### 3.6 Ex12

Transitive dependencies occur when a non-prime attribute (an attribute not part of any candidate key) depends on another non-prime attribute through an intermediary attribute. In the faculty relations there are none because all relations faculty(building, room, capacity) and location(building, faculty) have attributes that are candidate keys therefore, there are no non-key attributes that are transitively dependent on any candidate key. For the students relations students(student\_id, module\_id) student\_info(student\_name, student\_id, student\_email, programme, year, address, contact) module\_info(module\_id, module\_name, module\_leader, faculty, building, room,) cohort\_lecturers(student\_id, module\_id, lecturers) courseworks(student\_id, module\_id, exam\_mark, coursework) there are 2 transitive dependencies: exam\_mark  $\rightarrow$  coursework and building, room  $\rightarrow$  faculty.

### 3.7 Ex13

As we have already established in the faculty relation that building and room determine faculty we can say the same thing here in the module\_info relation. We can decompose it into new relations module\_info(module\_id, module\_name, module\_leader, building, room) and site(building, room, faculty). In the site relation building and room are the keys and all the fields have type TEXT. In the module\_info module\_id remains as the key and all the fields have type TEXT. Finally, as exam\_mark determines coursework in the courseworks relation we can create a new relation grade(exam\_mark, coursework) leaving the courseworks relation as courseworks(student\_id, module\_id, exam\_mark). In the courseworks relation student\_id and module\_id are of type TEXT but exam\_mark is of type INTEGER. In the grade relation all of the fields are of type INTEGER.

## 4 Modelling

The existing data must now be transformed into an SQLite database based on the specified relations.

### 4.1 Ex14

#### students

Attribute	Datatype
student_id	TEXT
module_id	TEXT

#### student\_info

Attribute	Datatype
student_name	TEXT

student_id	TEXT
student_email	TEXT
programme	TEXT
year	INTEGER
address	TEXT
contact	TEXT

### **cohort\_lecturers**

<b>Attribute</b>	<b>Datatype</b>
student_id	TEXT
module_id	TEXT
lecturers	TEXT

### **module\_info**

<b>Attribute</b>	<b>Datatype</b>
module_id	TEXT
module_name	TEXT
module_leader	TEXT
building	TEXT
room	TEXT

### **site**

<b>Attribute</b>	<b>Datatype</b>
building	TEXT
room	TEXT
faculty	TEXT

### **grade**

<b>Attribute</b>	<b>Datatype</b>
exam_mark	INTEGER
coursework	TEXT

### **courseworks**

<b>Attribute</b>	<b>Datatype</b>
------------------	-----------------

student_id	TEXT
module_id	TEXT
exam_mark	INTEGER

### faculty

Attribute	Datatype
building	TEXT
room	TEXT
capacity	INTEGER

### location

Attribute	Datatype
building	TEXT
room	TEXT
faculty	TEXT

## 4.2 Ex15

### Linux shell

```
1  sqlite3 university.db
2  --this creates a new database
3  --you must now create the tables: studentscsv and facultycsv
4  CREATE TABLE IF NOT EXISTS studentscsv (
5      student_name TEXT,
6      student_id INTEGER,
7      student_email TEXT,
8      programme TEXT,
9      year INTEGER,
10     address TEXT,
11     contact TEXT,
12     module_id TEXT,
13     module_name TEXT,
14     module_leader TEXT,
15     lecturer1 TEXT,
16     lecturer2 TEXT,
17     faculty TEXT,
18     building TEXT,
19     room TEXT,
20     exam_mark INTEGER,
21     coursework1 INTEGER,
22     coursework2 INTEGER,
23     coursework3 INTEGER
24 );
25
26 CREATE TABLE IF NOT EXISTS facultycsv (
27     faculty TEXT,
28     building TEXT,
29     room TEXT,
30     capacity INTEGER
31 );
32 --The data in the csv files are then imported to populate the tables
33 .mode csv
34 .import /home/pkp1u24/comp1314-cs/students.csv studentscsv
35 .import /home/pkp1u24/comp1314-cs/faculty.csv facultycsv
36 --then we dump the database into as ex15.sql
37 .output /home/pkp1u24/comp1314-cs/ex15.sql
38 .dump
```

## 4.3 Ex16

### Linux shell

```
1  sqlite3 university.db
2  --this opens the database
3  .output ex16.sql
4  --we then want to output our results to ex16.sql
5  --we have to create the tables to satisfy our normalised relations
6  --This is based on the existing data in university.db
7  CREATE TABLE IF NOT EXISTS students (
8  student_id INTEGER,
9  module_id TEXT,
10 PRIMARY KEY (student_id, module_id)
11 );
12
13 CREATE TABLE IF NOT EXISTS student_info (
14 student_id INTEGER PRIMARY KEY,
15 student_name TEXT,
16 student_email TEXT,
17 programme TEXT,
18 year INTEGER,
19 address TEXT,
20 contact TEXT
21 );
22
23 CREATE TABLE IF NOT EXISTS cohort_lecturers (
24 student_id INTEGER,
25 module_id TEXT,
26 lecturer TEXT,
27 PRIMARY KEY (student_id, module_id, lecturer)
28 );
29
30 CREATE TABLE IF NOT EXISTS module_info (
31 module_id TEXT PRIMARY KEY,
32 module_name TEXT,
33 module_leader TEXT,
34 building TEXT,
35 room TEXT
36 );
37
38 CREATE TABLE IF NOT EXISTS site (
39 building TEXT,
40 room TEXT,
41 faculty TEXT
42 );
43
44 CREATE TABLE IF NOT EXISTS grade (
45 exam_mark INTEGER,
46 coursework INTEGER
47 );
```

## Linux shell

```
48 CREATE TABLE IF NOT EXISTS courseworks (  
49     student_id TEXT,  
50     module_id TEXT,  
51     exam_mark INTEGER,  
52     PRIMARY KEY (student_id, module_id)  
53 );  
54  
55 CREATE TABLE IF NOT EXISTS faculty (  
56     building TEXT,  
57     room TEXT,  
58     capacity INTEGER,  
59     PRIMARY KEY (building, room)  
60 );  
61  
62 CREATE TABLE IF NOT EXISTS location (  
63     building TEXT,  
64     room TEXT,  
65     faculty TEXT,  
66     PRIMARY KEY (building, room)  
67 );  
68  
69  
70 INSERT INTO students (student_id, module_id)  
71 SELECT DISTINCT student_id, module_id  
72 FROM studentcsv;  
73  
74 INSERT INTO student_info (student_id, student_name, student_email,  
75     ↪ programme, year, address, contact)  
76 SELECT DISTINCT student_id, student_name, student_email, programme,  
77     ↪ year, address, contact  
78 FROM studentcsv;  
79  
80 INSERT INTO cohort_lecturers (student_id, module_id, lecturer)  
81 SELECT DISTINCT student_id, module_id, lecturer1  
82 FROM studentcsv  
83 UNION  
84 SELECT DISTINCT student_id, module_id, lecturer2  
85 FROM studentcsv;  
86  
87 INSERT INTO module_info (module_id, module_name, module_leader,  
88     ↪ building, room)  
89 SELECT DISTINCT module_id, module_name, module_leader, building, room  
90 FROM studentcsv;  
91  
92 INSERT INTO site (building, room, faculty)  
93 SELECT DISTINCT building, room, faculty  
94 FROM studentcsv;
```

### Linux shell

```
92 INSERT INTO grade (exam_mark, coursework)
93 SELECT DISTINCT exam_mark, coursework1
94 FROM studentcsv
95 UNION ALL
96 SELECT DISTINCT exam_mark, coursework2
97 FROM studentcsv
98 UNION ALL
99 SELECT DISTINCT exam_mark, coursework3
100 FROM studentcsv;
101
102 INSERT INTO courseworks (student_id, module_id, exam_mark)
103 SELECT DISTINCT student_id, module_id, exam_mark
104 FROM studentcsv;
105
106 INSERT INTO faculty (building, room, capacity)
107 SELECT DISTINCT building, room, capacity
108 FROM facultycsv;
109
110 INSERT INTO location (building, room, faculty)
111 SELECT DISTINCT building, room, faculty
112 FROM facultycsv;
```

This creates the necessary tables using the Create Table statement and populates them by using insert statements. It also uses Union to combine elements from a table and union all to combine all the elements into one element.

## 5 Querying

Now that the data is finally in a sensible format, you can now use it to answer questions! For each exercise, write an SQL query against your newly created normalised tables in your database. You must not create further tables specifically for these queries. Each SQL statement should be written in the report, as well as saved as `exjnumberi.sql` (e.g. `ex16.sql`). Write an SQLite statement to find:

## 5.1 Ex17

### Linux shell

```
1 sqlite3 university.db
2 --this opens the database
3 .output /home/pkp1u24/comp1314-cs/ex17.sql
4 --redirects the output to ex17.sql
5 SELECT building, SUM(capacity) AS total_capacity
6 FROM faculty
7 GROUP BY building;
8 --This query then calculates the total capacity of each building in
   ↪ the university
```

## 5.2 Ex18

### Linux shell

```
1 sqlite3 university.db
2 --this opens the database
3 .output /home/pkp1u24/comp1314-cs/ex18.sql
4 --redirects the output to ex18.sql
5 SELECT student.student_id, student.student_name, AVG(cw.exam_mark) AS
   ↪ average_exam_mark
6 FROM student_info student
7 JOIN courseworks cw ON student.student_id = cw.student_id
8 WHERE student.programme = 'Computer Science' AND student.year = 1
9 GROUP BY student.student_id, student.student_name
10 ORDER BY average_exam_mark DESC;
11 --This query then gets the average exam mark for each Year 1 Computer
   ↪ Science student by grouping them with their id and name and sorts
   ↪ it in descending order
```



## 5.3 Ex19

### Linux shell

```
1 sqlite3 university.db
2 .output /home/pkp1u24/comp1314-cs/ex19.sql
3 -- Calculate average marks for each module and faculty
4 SELECT mi.module_id, mi.module_leader, s.faculty, AVG(g.exam_mark +
   ↪ g.coursework) AS average_mark
5 FROM module_info mi
6 JOIN site s ON mi.building = s.building AND mi.room = s.room
7 JOIN courseworks cw ON mi.module_id = cw.module_id
8 JOIN grade g ON cw.exam_mark = g.exam_mark
9 GROUP BY mi.module_id, mi.module_leader, s.faculty
10 HAVING AVG(g.exam_mark + g.coursework) = (
11 SELECT MAX(module_avg)
12 FROM (
13 SELECT AVG(g_inner.exam_mark + g_inner.coursework) AS module_avg,
   ↪ s_inner.faculty AS faculty
14 FROM module_info mi_inner
15 JOIN site s_inner ON mi_inner.building = s_inner.building AND
   ↪ mi_inner.room = s_inner.room
16 JOIN courseworks cw_inner ON mi_inner.module_id = cw_inner.module_id
17 JOIN grade g_inner ON cw_inner.exam_mark = g_inner.exam_mark
18 WHERE s_inner.faculty = s.faculty
19 GROUP BY mi_inner.module_id
20 )
21 )
22 ORDER BY average_mark DESC;
23
```

This code calculates the average marks for each module within each faculty. Finds the maximum average mark within each faculty. Joins the two queries to find the modules with the highest average marks for their respective faculty. Orders the results by the highest average mark.

The output from ex19.sql was :

```
ANIM3403|Michelle Mueller|Faculty of Art|50.0423796259451
ELEC4404|Margaret Sanders|Faculty of Technology|49.4737901381392
ENG1402|Bruce Schneider|Faculty of Languages|47.4435638152501
```

## 5.4 Ex20

### Linux shell

```
1 sqlite3 university.db
2 --this opens the database
3 .output /home/pkp1u24/comp1314-cs/ex20.sql
4 SELECT mi.module_id, s.room, s.building, f.capacity,
5 COUNT(cw.student_id) AS student_count
6 FROM module_info mi
7 JOIN courseworks cw ON mi.module_id = cw.module_id
8 JOIN site s ON mi.building = s.building AND mi.room = s.room
9 JOIN faculty f ON s.building = f.building AND s.room = f.room
10 GROUP BY mi.module_id, s.room, s.building, f.capacity
11 HAVING COUNT(cw.student_id) > f.capacity
12 ORDER BY student_count DESC;
```

This code finds the tables Where student count exceeds room capacity by joining module\_info, courseworks, site, and faculty. It uses COUNT(cw.student\_id) to calculate the number of students for each module. It then compares the student count to room capacity and groups the data by module, room, building, and capacity. The HAVING statement ensures only rows where the student count exceeds the room's capacity are included. The ORDER BY statement sorts the results by student\_count in descending order