

Chapter 4 - Update Anomalies and Normalization

1. [SQL Queries](#)
 - i. [A query that displays the city and country names, and populations from cities that have at least 200,000 people and less than 250,000 people.](#)
 - ii. [A query that shows the countries where 8 or more languages are spoken, and how many languages are spoken in that country.](#)
 - iii. [A query that displays the population, names and countries \(name and ID\) of the largest city in each country](#)
 - iv. [A query that displays the names of cities and their countries where the capital city is the largest of all cities listed for that country](#)
2. [Programming](#)
 - i. [Python Implementation](#)
3. [E4.5 from the book \(Page 121\)](#)
 - i. [Insertion anomaly](#)
 - ii. [Deletion anomaly](#)
 - iii. [Modification anomaly](#)
 - iv. [Full key functional dependencies, partial key functional dependencies, transitive functional dependencies](#)
 - v. [2NF Result](#)
 - vi. [3NF Result](#)
 - vii. [Using 3NF result, describe how anomalies are eliminated](#)
4. [Tools and Reference](#)
 - i. [Reference](#)
 - ii. [Tools](#)

This homework requires to have the "world" example Mysql database installed.

To get and install the sample database: [Setting Up the world Database](#)

SQL Queries

Note: All SQLs are developed with MYSQL, compatibilities with other RDBMS is not guaranteed.

Implement SQL queries and get the result sets of the following requirements.

A query that displays the city and country names, and populations from cities that have at least 200,000 people and less than 250,000 people.

```
1  -- Result set row count: 353
2  SELECT
3      city.name, country.name AS countryname, city.Population
4  FROM city
5  JOIN country ON city.countrycode = country.code;
```

```

6 | WHERE city.population >= 200000
7 | AND city.population < 250000;

```

A query that shows the countries where 8 or more languages are spoken, and how many languages are spoken in that country.

```

1 | -- Result set row count: 30
2 | SELECT
3 |     c.Code, c.Name, COUNT(c1.`Language`)
4 | FROM country c
5 | JOIN countrylanguage c1 ON c.Code = c1.CountryCode
6 | GROUP BY c.Code
7 | HAVING COUNT(c1.`Language`) >= 8;

```

A query that displays the population, names and countries (name and ID) of the largest city in each country

```

1 | -- Result set row count: 232
2 | SSELECT
3 |     cn.Code, cn.Name,
4 |     c.Name AS cityname,
5 |     c.ID AS cityid, MAX(c.Population) AS population
6 | FROM country cn
7 | JOIN city c ON cn.Code = c.CountryCode
8 | GROUP BY cn.code;

```

A query that displays the names of cities and their countries where the capital city is the largest of all cities listed for that country

```

1 | -- Result set row count: 184
2 | SELECT
3 |     c.ID, c.Name,
4 |     cn.Name AS countryname, c.Population
5 | FROM city c
6 | JOIN country cn ON (cn.Code = c.CountryCode AND cn.Capital = c.ID)
7 | JOIN (
8 |     SELECT c1.ID AS cityid, MAX(c1.Population)
9 |     FROM country cn1
10 |     JOIN city c1 ON cn1.Code = c1.CountryCode
11 |     GROUP BY cn1.code) largestcities ON largestcities.cityid = c.ID;

```

Programming

Imagine that you had two datafile cities.csv and countries.csv (text files, where each row is a row in the database, and each column is separated by commas) that represent the database world that we've been working with so far. An example piece of cities.csv looks like:

```

1 | 1,Kabul,AFG,Kabul,1780000

```

```
2 | 2,Qanadahar,AFG,Qandahar,237500
3 | ...
```

Without using any SQL, write a program (in any language) that answers #4 above. It will involve using multiple loops and variables (likely arrays/associative arrays) as well.

Python Implementation

```
1 | #!/usr/bin/env python
2 | # -*- coding: utf-8 -*-
3 |
4 | """[Homework 5]
5 | This homework will require you to use the World Database.
6 |
7 | This python module/program simulates an RDB SQL query which is an equivalent to
8 | the 4th query in this homework.
9 |
10 | [The 4th query]
11 | Write a query that displays the names of cities and their countries where
12 | the capital city in the largest of all cities listed for that country.
13 |
14 | SELECT c.ID, c.Name, cn.Name AS countryname, c.Population
15 | FROM city c
16 | JOIN country cn ON (cn.Code = c.CountryCode AND cn.Capital = c.ID)
17 | JOIN (
18 |     SELECT countrycode, ID AS cityid, MAX(Population)
19 |     FROM city
20 |     GROUP BY countrycode) largestcities ON largestcities.cityid = c.ID;
21 |
22 | [Python simulation]
23 | Imagine that you had two data file cities.csv and countries.csv that represent
24 | the database world that we've been working with so far.
25 | An example piece of cities.csv looks like:
26 | 1,Kabul,AFG,Kabul,1780000
27 | 2,Qanadahar,AFG,Qandahar,237500
28 | ...
29 |
30 | Without using any SQL, write a program (in any language) that answers #4 above.
31 | It will involve using multiple loops and variables (likely arrays/associative
32 | arrays) as well.
33 |
34 | [Note]
35 | Both cities.csv and countries.csv must be a full dump from their database
36 | tables.
37 |
38 | This program results the same number of cities as Q4 does: 184
39 | """
40 |
41 | import csv
42 | import os.path
43 | import io
```

```

44 import itertools
45 import traceback
46
47
48 _city_csv = "cities.csv"
49 _country_csv = "countries.csv"
50 _csv_dialect = 'excel'
51
52
53 def _convertToNumber(d):
54     """Convert dictionary string values into integer if the values are formed
55     with digits.
56     """
57     return dict((k, v if not v.isdigit() else int(v)) for (k, v) in d.items())
58
59
60 def _read_csv(csv_file):
61     """Read the given csv file and store the result sets into a list of
62     dictionary."""
63     if csv_file == None or csv_file.strip() == "":
64         raise ValueError("Missing file path")
65     if not os.path.exists(csv_file) and not os.path.isfile(csv_file):
66         raise IOError("File does not exist: %s" % csv_file)
67
68     result_set = None
69     with open(csv_file, "rb") as csvfile:
70         reader = csv.DictReader(csvfile, dialect=_csv_dialect)
71         result_set = [_convertToNumber(r) for r in reader]
72     # test
73     # print re
74     return result_set
75
76
77 def _groupby(result_set, groupby_key, sort_key=None, sort_desc=False):
78     """Slice the given list of dict into sublist groups by the given key.
79     If the sub groups should be sorted, specify sort_key.
80     """
81     grouped = []
82     if not result_set or len(result_set) == 0 or not groupby_key or len(groupby_key.strip()) == 0:
83         return grouped
84     for column, group in itertools.groupby(result_set, lambda d: d[groupby_key]):
85         g = list(group)
86         if sort_key and len(sort_key.strip()) != 0:
87             g = sorted(g, key=lambda k: k[sort_key], reverse=sort_desc)
88         grouped.append(g)
89     return grouped
90
91
92 def _group_max(result_set, groupby_key, sort_key):
93     """Get a list of dict records which have the largest value of sort_key in
94     groups.

```

```

95     This method does the SQL part:
96
97     SELECT *, max(sort_key)
98     FROM table
99     GROUP BY groupby_key
100     """
101     if not result_set or len(result_set) == 0:
102         return []
103     return [x[0] for x in _groupby(result_set, groupby_key, sort_key, True)]
104
105
106 def _join(src, to_join, on_src_key, on_to_join_key):
107     """Get a sublist of src which is in an intersection of 2 lists of dictionary
108     on one key.
109     This method simulates a SQL JOIN, such as:
110
111     SELECT t.*
112     FROM table t
113     JOIN table1 t1 ON t.key = t1.key
114     """
115     on_values = set(d[on_to_join_key] for d in to_join)
116     return [d for d in src if d[on_src_key] in on_values]
117
118
119 def _select_capitals_with_max_population(cities, countries):
120     """Get a list of cities which both are capital and have the most population
121     in the country"""
122     groupby_key = "CountryCode"
123     sort_key = "Population"
124     capital_key = "Capital"
125     city_id = "ID"
126     # cities with max(population)
127     grouped_cities_max_population = _group_max(cities, groupby_key, sort_key)
128     # capital cities
129     capital_cities = _join(cities, countries, city_id, capital_key)
130     return _join(capital_cities, grouped_cities_max_population, city_id, city_id)
131
132
133 # Main entrance
134 if __name__ == '__main__':
135     try:
136         cities = _read_csv("cities.csv")
137         countries = _read_csv("countries.csv")
138         if not cities or len(cities) == 0:
139             raise RuntimeError("No data read from city file")
140         if not countries or len(countries) == 0:
141             raise RuntimeError("No data read from country file")
142         capitals_max_population = \
143             _select_capitals_with_max_population(cities, countries)
144         for row in capitals_max_population:
145             print row

```

```
146         print "Cities: %d\n" % len(capitals_max_population)
147     except Exception as e:
148         print "Runtime error happened: %s" % e
149         traceback.print_exc()
150         exit(1)
151     exit(0)
```

E4.5 from the book (Page 121)

AIRPORT KLX TABLE

Date	AirlineID	Airline Name	TerminalID	NumberOfGates	NumberOfDepartingFlights
------	-----------	--------------	------------	---------------	--------------------------

Insertion anomaly

Cannot insert a new terminal, say (D:15), into the table without an airline information.

Deletion anomaly

Cannot delete airline SWA from the table without deleting terminal C's information.

Modification anomaly

Changing terminal B's NumberOfGates will make change to 3 records.

Full key functional dependancies, partial key functional dependancies, transitive functional dependancies

- 1 Full key functional dependancies: Date, AirlineID --> NumberOfDepartingFlights
- 2 Partial key functional dependancies: AirlineID --> AirlineName, TerminalID, NumberOfGates
- 3 Transitive functional dependancies: TerminalID --> NumberOfGates

2NF Result

TRAFFIC TABLE

Date	AirlineID	NumberOfDepartingFlights
------	-----------	--------------------------

AIRLINE TABLE

AirlineID	Airline Name	TerminalID	NumberOfGates
-----------	--------------	------------	---------------

3NF Result

TRAFFIC TABLE

Date	AirlineID	NumberOfDepartingFlights
------	-----------	--------------------------

AIRLINE TABLE

TerminalID	AirlineID	Airline Name
------------	-----------	--------------

TERMINAL TABLE

TerminalID	NumberOfGates
TerminalID	NumberOfGates

Using 3NF result, describe how anomalies are eliminated

- 1 Insertion: No anomaly because insertion a new terminal won't require full information of traffic and airline information.
- 2 Modification: No anomaly because update upon terminal B's NumberOfGates affect only 1 record in TERMINAL table.
- 3 Deletion: No anomaly because deleting terminal C in TERMINAL table will not affect AIRLINE table.

Tools and Reference

Reference

- SQL JOINS: [LEFT JOIN vs. LEFT OUTER JOIN](#)

Tools

- HeidiSQL, A lightweight and handy tool for SQL development: [HeidiSQL - MySQL, MSSQL and PostgreSQL made easy](#)

---- Jason, 10/12/2016