CS541 – project 2

Wen-Ling Chi

Email: chi64@purdue.edu

Part 1 – Understanding the data

1-1 Is each column present in every data entry? List which columns do not appear in every data entry.

Filename	Columns	
	• city	
airport icon	airport_province	
airport.json	• island	
	• elevation	
	• population	
city.json	• elevation	
	• country_capital	
	• capital	
country.json	• province	
country-other	- otherweene	
localname.json	• othername	
countrypopulations.json	• capital	
	• gdp	
	agriculture	
economy.json	• service	
economy.json	• industry	
	• inflation	
	unemployment	
ethnicgroup.json	• country_capital	
etinicgroup.json	ethnic_group_percentage	
language.json	• country_capital	
language.json	percentage	
	• city	
organization.json	• country	
Organization.json	• province	
	established	
	• independence	
politics.json	wasdependent	
politics.jsori	dependent	
	• government	
	country_province	
population.json	population growth	
	infant mortality	
	• population	
province.json	• capital	
province.json	• capprov	

1-2 Given the file descriptions and the knowledge graph, are there any redundant columns in the files?

List which files have redundant columns. Name the columns and provide reasons why, in your opinion, those columns are redundant?

Filename	Redundant columns / Reason	
airnort ican	country_name: we have country_code as a foreign key to the	
airport.json	country table.	
	 country_name: we have country as a foreign key to the country 	
city.json	table.	
city.json	• country_capital: if we know the country code, then we should also	
	know the capital.	
	 country_area: since we know its city name, we can also know its 	
	country. Through country table, we can find the area of the country.	
cityothername.json	• country_capital: since we know its city name, we can also know its	
	country. Through country table, we can find the capital of the	
	country.	
	• province: we already have "capital" as the foreign key referencing	
country.json	the city table, and we can find out the province the capital belongs to	
	from "province" in the city	
	country_name: we have country_code as a foreign key to the	
countrypopulations.json	country table.	
3	• capital: if we know the country name, then we should also know the	
	capital.	
	• country: we have country_code as a foreign key to the country table.	
ethnicgroup.json	• country_capital: we have country_code as a foreign key to the	
	country table.	
	country_capital: we have country as a foreign key to the country	
language.json	table.	
language,json	• country_area: we have country as a foreign key to the country table.	
	country: this table doesn't connect to country table, so this column	
	is unnecessary. Also, if we know the city, then we should also know	
located-on.json	the country.	
	 province_area: this is unnecessary attribute for this data. 	
	country_name: we have country_code as a foreign key to the	
population.json	country table.	
	• country_province: if we know the country_code, we can also know	
	its province. Only "Gaza Strip" has no province, but in the country	
	table, it also has no province.	
	 country_name: we have country_code as a foreign key to the 	
religion.json	country table.	
	 country_population: if we know the country, then we should also 	
	know the population of the country.	

1-3 What logical constraints would you set for which columns in the files? For example, latitude should always be in the range of [-90, 90], and the longitude should be in [-180, 180]. Here you don't have to set constraints for every column of every file.

Filename	Constraints		
airmart is an	• latitude: [-90, 90]		
airport.json	• longitude: [-180, 180]		
attention a	• latitude: [-90, 90]		
city.json	• longitude: [-180, 180]		
country.json	• encompass_percentage: [0, 100]		
	• agriculture: [0, 100]		
	• service: [0, 100]		
economy.json	industry: [0, 100]unemployment: [0, 100]		
ethnicgroup.json	ethnic_group_percentage: [0, 100]		
language.json	• percentage: [0, 100]		
religion.json	• percentage: [0, 100]		

Part 2 - Relational database schema

2-1 Write an SQL script that contains the schema for the storage and manipulation of the data provided in the JSON files.

See ./sql_schema/project2_schema.sql for the Oracle SQL schema script.

(Use Oracle SQL Developer to run sql file)

2-2 Compile the entity-relationship diagram of your database.

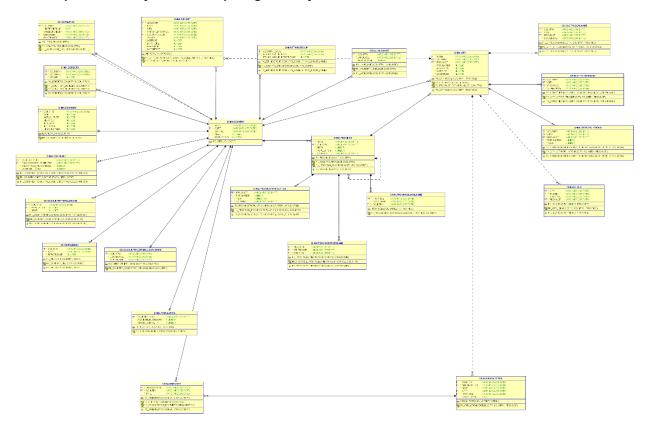


Fig. 1 - The entity-relationship diagram of my database

Part 3 – Inserting the data to Oracle DBMS

Environment

This part will use python3 and Oracle SQL Developer.

Execution Steps

 Execute python file to create insert values for Oracle database. The output will be in ./sql_schema/project2_data.sql
 See ./sql_schema/createInsertValue_sql.py

```
Executive Command:
python sql_schema/createInsertValue_sql.py
```

Execute sql file to insert values into Oracle SQL database.
 See ./sql_schema/project2_data.sql
 (Use Oracle SQL Developer to run sql file)
 (It is possible to take a long to insert values into database)

Part 4 - Using MongoDB

Environment

This part will use python3 and MongoDB(3.6). The reason for using MongoDB (3.6) is that we need to use a tool - mongoimport.

After downloading MongoDB(3.6), put ~/MongoDB/Server/3.6/bin into environment variables.

Execution Steps

Create MongoDB schema
 Execute python file to create MongoDB schema. It will generate 23 collections in MongoDB.
 See ./mongodb_schema/mongodb_schema.py

```
Executive Command:
python mongodb_schema/mongodb_schema.py
```

Create json files for the data for MongoDB
 Execute python file to create json files for inserting values in MongoDB. The output json files will be in ./mongodb_data directory.

 See ./mongodb_schema/insertValue_mongodb.py

```
Executive Command:
python mongodb_schema/insertValue_mongodb.py
```

Insert data into MongoDB
 Click bat file to insert data into MongoDB.
 See ./mongodb_schema/import_data_mongodb.bat

Executive command in bat file (Samples)

```
mongoimport --db project2 --collection ethnicgroup .\mongodb_data\ethnicgroup.json --jsonArray
mongoimport --db project2 --collection language .\mongodb_data\language.json --jsonArray
mongoimport --db project2 --collection politics .\mongodb_data\politics.json --jsonArray
mongoimport --db project2 --collection population .\mongodb_data\population.json --jsonArray
mongoimport --db project2 --collection religion .\mongodb_data\portion.json --jsonArray
mongoimport --db project2 --collection borders .\mongodb_data\portion.json --jsonArray
mongoimport --db project2 --collection province .\mongodb_data\province.json --jsonArray
mongoimport --db project2 --collection city .\mongodb_data\city.json --jsonArray
```

Fig. 2 – samples in import_data_mongodb.bat

Part 5 - Queries

Environment

This part will use python3 and Oracle SQL Developer.

5-1

Implement the following queries in both Oracle and MongoDB. For each of the following queries, capture the execution time for both DBMS environments (question 5-2).

а

Select the name of the country that shares the longest border with one or more other countries.

- Oracle: guery/oracle guery/ 5-1-a.sql
- MongoDB: <u>query/mongodb query/mongodb query.py</u>
 Executive Command: python query/mongodb_query/ mongodb_query.py 1

b

For each country, list the city that has the maximum average population.

- Oracle: query/oracle query/5-1-b.sql
- MongoDB: <u>query/mongodb query/mongodb query.py</u> Executive Command:

python query/mongodb_query/ mongodb_query.py 2

c

List the following information of the countries where agriculture contributes the most to their economy. Name of the country, its GDP, percentage contribution from agriculture, and inflation.

- Oracle: <u>query/oracle query/ 5-1-c.sql</u>
- MongoDB: <u>query/mongodb query/mongodb query.py</u>
 Executive Command: python query/mongodb_query/ mongodb_query.py 3

d

List the following information of the countries in the descending order of their ethnic diversity. Name of the country, number of ethnic groups, and the percentage of the major ethnicity. Ethnic diversity of a country increases as the number of ethnic groups that live in the country increases.

- Oracle: <u>query/oracle query/ 5-1-d.sql</u>
- MongoDB: <u>query/mongodb query/mongodb query.py</u> **Executive Command:**

python query/mongodb query/ mongodb query.py 4

e

Find the number of countries with the number of ethnic groups equal to the number of languages used.

- Oracle: query/oracle query/ 5-1-e.sql
- MongoDB: query/mongodb query.py **Executive Command:**

python query/mongodb_query/ mongodb_query.py 5

f

The GDP and IMR (infant mortality rate) together define the development of a country. Find the top 10 countries with the highest GDP (in decreasing order) and their corresponding IMR value.

- Oracle: <u>query/oracle query/ 5-1-f.sql</u>
- MongoDB: <u>query/mongodb query/mongodb query.py</u> **Executive Command:**

python query/mongodb_query/ mongodb_query.py 6

g

Find the following information of the country with the highest religious freedom. Name of the country, the religions that are practiced in the country. Note that the larger the number of religions practiced in a country, the higher its religious freedom is.

- Oracle: <u>query/oracle query/ 5-1-q.sql</u>
- MongoDB: <u>query/mongodb query/mongodb query.py</u> **Executive Command:**

python query/mongodb_query/ mongodb_query.py 7

h

What is the proportion of commonwealth countries in the top 100 countries in terms of GDP.

- Oracle: query/oracle query/ 5-1-h.sql
- MongoDB: guery/mongodb guery.py **Executive Command:**

python query/mongodb_query/ mongodb_query.py 8

i

Find the country with the largest population density in each continent. List the name of the country and its population density against the name of the continent. The population density of a country is the ratio between its total population and area.

```
    Oracle: <a href="query/oracle_query/5-1-i.sql">query/oracle_query/5-1-i.sql</a>
    MongoDB: <a href="query/mongodb_query/mongodb_query.py">query/mongodb_query/mongodb_query.py</a>
    j
```

Find the names of countries and their capitals where the capital has more than one airport.

- Oracle: query/oracle query/5-1-j.sql
- MongoDB: <u>query/mongodb query.py</u> Executive Command:

python query/mongodb_query/ mongodb_query.py 10

5-2

Create a table to summarize the execution time in the two DBMS for each of the queries.

Query Number	Execution time		
	Oracle (ms)	MongoDB (ms)	
5-1-a	259	97	
5-1-b	650	416	
5-1-c	99	45	
5-1-d	316	81	
5-1-e	109	359	
5-1-f	95	42	
5-1-g	54	35	
5-1-h	71	509	
5-1-i	114	106	
5-1-j	98	188	

5-3

Comment on the execution times you have obtained. Give educated reasons why one DBMS has performed better than the other?

In the above form, most of MongoDB's execution times are better than Oracle because MongoDB has less constraints. Thus, compare to Oracle SQL, we need less checks when we execute the query in MongoDB.

For the few cases, MongoDB's execution times are worse than Oracle. From my observation, MongoDB performs worse than Oracle SQL if we need to use multiple collections/tables. It is possible that the Oracle SQL Database will cache and create automatically indexing internally.

5-4

Suggest one performance improvement technique that you can use on each DBMS to improve the performance. Implement your technique on each DBMS and report your results in a table similar to questions 5-2.

Oracle

Use the "CREATE INMEMORY JOIN GROUP" statement to create a join group, which is an object that specifies frequently joined columns from the same table or different tables. When you create a join group, Oracle Database stores special metadata for the columns in the global dictionary, which enables the database to optimize join queries for the columns.

Because code in country table and country in city table are frequently used by the queries, so I use the above method to create a join group.

See query/oracle_query/optimized.sql

MongoDB

<u>Create Indexes to Support Queries</u>. If a query searches multiple fields, create a compound index. Scanning an index is much faster than scanning a collection.

See the function - optimizedDB(db) in query/mongodb_query/mongodb_query.py

• Execution time

Query Number	Execution time		
	Oracle (ms)	MongoDB (ms)	
5-1-a	181	58	
5-1-b	101	61	
5-1-c	91	36	
5-1-d	83	79	
5-1-e	72	69	
5-1-f	61	35	
5-1-g	38	28	
5-1-h	54	32	
5-1-i	91	70	
5-1-j	39	55	

5-5

Comment on the execution times you have obtained after implementing your performance improvement technique. How much improvement do you observe compared to your results in 5-2?

Oracle

All queries reduce their execution time. Because almost every query uses country table - code or city table - country. The higher the frequency they use, the more the execution time is reduced.

MongoDB

All queries reduce their execution time. If the query needs to scan different collections, the execution time will be reduced even more. Because I store the fields in the index, scanning the index can reduce the execution time.