

Explore Databases Application Final Report
INF551
Mingjun Liu(1321657359)
Minghong Xu(3989201972)

Overview

Our final web application has the website: <https://exploredatabase.herokuapp.com/>

Project demo Youtube link: https://www.youtube.com/watch?v=iTNC5oMu_8E

More related materials included in the google drive. Link of google drive:
<https://drive.google.com/open?id=1bkkWJfZOVSNMaxs8EIMGyJMRVJ5MYBS>

Project idea

In this project, we are going to implement a web application that allows our user to explore the databases we choose. The user is able to search one or multiple keywords within any one of our databases and navigate between the tables by hyperlinks. The application will involve techniques about both front end and back end. We need to design an user interface that allows the user to choose the database to explore and submit the keywords to the back end, which employs forms from html. The hyperlink of the foreign key in the tables should allow the user to jump to the tables with the foreign key as an primary key, which means each table should have at least one foreign key relationship to the other tables. Back end responsibility includes importing databases from MySQL and uploading databases to Firebase, realizing the import function, search function and the navigation function, which are described with detailed explanation in the following parts. Since both of the members have lacked experience on designing user interfaces, we used the web application helper herokuapp.com to complete the UI. At this time, the three databases we chose are World, Movies, and FIFA World Cup.

The first database is the world database provided by the professor. The database contains three tables: city, country and countrylanguage. Code is the primary key in the country table, which is named as CountryCode and acts as foreign key in city table and countrylanguage table.

The second database was adapted from MovieLens, which is a movie recommendation service. The database contains three tables, which are movies, tags and ratings. The movie table includes movieId, movie titles and genres the movies belong to. The tags table includes userId, movieId and tags given to the movie from the users. The rating table includes userId, movieId and rating to the movies. User ids are consistent between ratings table and tags table, and movie ids are consistent among all three tables.

movies				tags				ratings			
movieId	title	genres		userId	movieId	tag	timestamp	userId	movieId	rating	timestamp
1	Toy Story (1995)	Adventure Animation Chil		2	60756	funny	1445714994	1	1	4	964982703
2	Jumanji (1995)	Adventure Children Fanta		2	60756	Highly quotable	1445714996	1	3	4	964981247
3	Grumpier Old Men	Comedy Romance		2	60756	will ferrell	1445714992	1	6	4	964982224
4	Waiting to Exhale	Comedy Drama Romance		2	89774	Boxing story	1445715207	1	47	5	964983815
5	Father of the Bride	Comedy		2	89774	MMA	1445715200	1	50	5	964982931
6	Heat (1995)	Action Crime Thriller		2	89774	Tom Hardy	1445715205	1	70	3	964982400
7	Sabrina (1995)	Comedy Romance		2	106782	drugs	1445715054	1	101	5	964980868
8	Tom and Huck (19	Adventure Children		2	106782	Leonardo DiCapi	1445715051	1	110	4	964982176
9	Sudden Death (19	Action		2	106782	Martin Scorsese	1445715056	1	151	5	964984041
10	GoldenEye (1995)	Action Adventure Thriller		7	48516	way too long	1169687325	1	157	5	964984100
11	American Preside	Comedy Drama Romance		18	431	Al Pacino	1462138765	1	163	5	964983650
12	Dracula: Dead and	Comedy Horror		18	431	gangster	1462138749	1	216	5	964981208
13	Balto (1995)	Adventure Animation Chil		18	431	mafia	1462138755	1	223	3	964980985
14	Nixon (1995)	Drama		18	1221	Al Pacino	1461699306	1	231	5	964981179
15	Cutthroat Island (1	Action Adventure Roman									

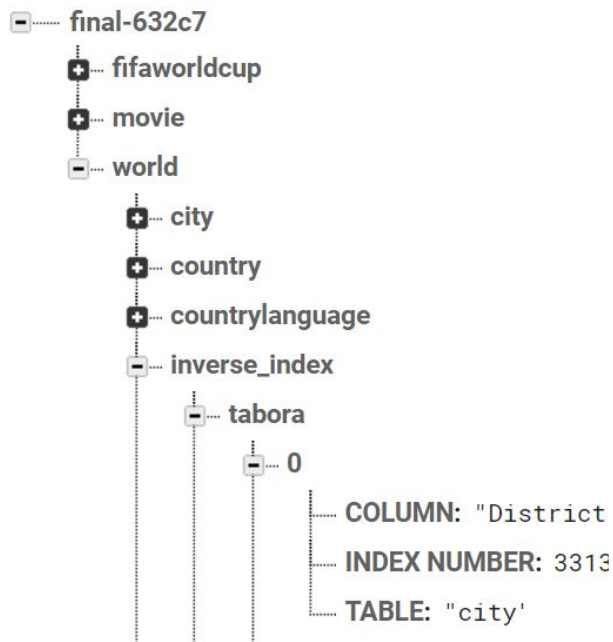
The third database is the FIFA World Cup database from Kaggle Data. The database contains three tables: WorldCupMatches, WorldCupPlayers and WorldCups. This database has a more complicated foreign key relationship between the tables. For example, Year is the primary key in the WorldCup table, and the foreign key in WorldCupMatches table; matchID is the primary key in the WorldCupMatches table, and the foreign key in WorldCupPlayers.

Working components

Part 1: import.py

After running the execution syntax, the program will show “inverse index uploaded to Firebase” if everything works well. The link to Firebase will create a node “world” that contains all the tables in the database and an inverse_index that contains the index location of a word.

```
C:\Users\marga\Dropbox\USC\2020SPRING\INF551\project>python import.py world world
use information_schema
select TABLE_NAME from tables where TABLE_SCHEMA = 'world'
select COLUMN_NAME from COLUMNS where TABLE_NAME = 'city' and TABLE_SCHEMA = 'world'
select COLUMN_NAME from COLUMNS where TABLE_NAME = 'country' and TABLE_SCHEMA = 'world'
select COLUMN_NAME from COLUMNS where TABLE_NAME = 'countrylanguage' and TABLE_SCHEMA = 'world'
use world
select * from city
select * from country
select * from countrylanguage
****Return database data****
***upload table city ***
***upload table country ***
***upload table countrylanguage ***
***Inverse index uploaded to Firebase***
```



Part 2: Web Application

Our main page:

Database Exploration

← → ↻ exploredatabase.herokuapp.com

★ M 🔔

Home Page

Database Exploration

This web application is a simple search application. Users can type in the words they want to search. And after choosing the database they want to search words at, they can click 'confirm' to search. In the searching results, users can click the foreign key to navigate to the related tables.

[Back to the Homepage](#)

Search result for 'China Chinese hangzhou' of World database:

https://exploredatabase.herokuapp.com

← → ↻ exploredatabase.herokuapp.com/search

☆ M 🔔

country

Code	Name	Continent	Region	SurfaceArea	IndepYear	Population	LifeExpectancy	GNP	GNPOld	LocalName	GovernmentForm	HeadOfState	Capital	Code2
CHN	China	Asia	Eastern Asia	9572900.00	1523	1277558000	71.4	982268.00	917719.00	Zhongquo	People's Republic	Jiang Zemin	1891	CN
HKG	Hong Kong	Asia	Eastern Asia	1075.00		6782000	79.5	166448.00	173610.00	Xianggang/Hong Kong	Special Administrative Region of China	Jiang Zemin	937	HK
MAC	Macao	Asia	Eastern Asia	18.00		473000	81.6	5749.00	5940.00	Macau/Aomen	Special Administrative Region of China	Jiang Zemin	2454	MO

city

ID	Name	CountryCode	District	Population
1905	Hangzhou	CHN	Zhejiang	2190500

countrylanguage

CountryCode	Language	IsOfficial	Percentage
AUS	Canton Chinese	F	1.1
BRN	Chinese	F	9.3
CAN	Chinese	F	2.5
CHN	Chinese	T	92.0
CRI	Chinese	F	0.2
CXR	Chinese	F	0.0
HKG	Canton Chinese	F	88.7
JPN	Chinese	F	0.2
KHM	Chinese	F	3.1
KOR	Chinese	F	0.1
MAC	Canton Chinese	F	85.6

Search result for 'iconic fiction' of Movie database:

https://exploredatabase.herokuapp.com/search

movies

movieId	title	genres
296	Pulp Fiction (1994)	Comedy Crime Drama Thriller
46976	Stranger than Fiction (2006)	Comedy Drama Fantasy Romance

ratings

Index	userId	movieId	rating	timestamp
-------	--------	---------	--------	-----------

tags

Index	userId	movieId	tag
3468	599	296	iconic
453	125	7254	science fiction

[Back to the Homepage](#)

Search result for 'China persie' of FIFA World Cup database:

Year	Datetime	Stage	Stadium	City	Home Team Name	Home Team Goals	Away Team Goals	Away Team Name	Attendance	Half-time Home Goals	Half-time Away Goals	Referee	Assistant 1	Assistant 2	RoundID	MatchID	Home Team Initials	Away Team Initials
2002.0	04 June 2002 - 15:30	Group C	Gwangju World Cup Stadium	Gwangju	China PR	0.0	2.0	Costa Rica	27217.0	0.0	0.0	VASSARAS Kyros (GRE)	MATOS Carlos (POR)	POOL Jaap (NED)	43950100.0	43950012.0	CHN	CRC
2002.0	08 Jun 2002 - 20:30	Group C	Jeju World Cup Stadium	Jeju	Brazil	4.0	0.0	China PR	36750.0	3.0	0.0	FRISK Anders (SWE)	LINDBERG Leif (SWE)	FIERRO Bomer (ECU)	43950100.0	43950026.0	BRA	CHN
2002.0	13 Jun 2002 - 15:30	Group C	Seoul World Cup Stadium	Seoul	Turkey	3.0	0.0	China PR	43605.0	2.0	0.0	RUIZ Oscar (COL)	TOMUSANGE Ali (UGA)	CHARLES Curtis (ATG)	43950100.0	43950042.0	TUR	CHN

Implementation details

Most of our work is built based on python. For the data, we write an import file, which uses python to script the data from MySQL, clean the data and upload the data to Firebase, so that the web application can retrieve data from Firebase directly without downloading data to the local disk. While reading the databases into python, since we don't know the columns of each table, we need to use the `information_schema` database to help us read the columns. We created a list that contains the name of the tables and corresponding names of the columns. After that, the cells of the database are saved by lists in memory. To make locating the words easier, when creating the index dictionary for each database, we used the index number as the primary key, which was generated by Firebase when we uploaded the database.

Furthermore, we write backend functions with python including searching words and retrieving data from tables with specific keys. In order to make the search function work, we also create an index table at Firebase for each database, which creates an index for each word appearing in a database. Some tuples of the table may contain more than one search word the users type in. So we also rank the tuple returned by the table with the times of search key appeared.

The last step is to realize the application on the website. We use the Flask Framework in python to write the website. Flask is a very useful package in python. It transfers many complex HTML functions and connection stuff to python code. Although Flask helps us a lot, we still write many HTML files to implement the templates of web pages. We also create hyperlink for every word functioned as foreign key that users can click to navigate to other tables with the key. And users can navigate continuously when there are foreign keys in the returned table. For the words that are foreign keys, they show with blue color. Finally, we use heroku to publish our website.

Performance analysis

We use the time function in python to record the time each search used. And use the time to analyze the query processing performance. As for the query mentioned above, we repeated several times to increase the reliability.

Search for 'China Chinese hangzhou' of World database takes 6.49s ~ 6.71s

Search for 'iconic fiction' of Movie database takes 0.93s ~ 1.27s

Search for 'China persie' of FIFA World Cup database takes 0.95s ~ 1.03s

It is obvious that searching in the world database takes much more time. The reason is that we have 3 searching words for that database. We tried to search only two words, and the time decreased to around 1s. Therefore, in general, all three databases behave similarly, which has nothing to do with the sizes of databases. The reason behind that is we use python to get data from Firebase, and transfer to the data type of dictionary. Dictionary is a hash table. Therefore, the time to get data is linear.

To improve the performance, we might find enhancement at the number of going through the database. At our existing searching code, for each search word the user types in, the function would go through the whole database one time. So if the user searches 10 words, the function would go through the database 10 times, which is a huge consumption. One way to improve is to cache all the search words in memory, and search them together in one iteration, that would save a lot of time.

Responsibility and work

There are two members in our team: Mingjun Liu and Minghong Xu. We do not have any abilities related to web development or mobile apps before. So we learnt web development by ourselves. We design and develop the web application together.

Conclusion

The database exploration web application consists of the knowledge we learned during the semester, and let us apply different kinds of data management techniques together. During the coursework, we learned how to download the database from MySQL into a csv file and upload a csv file and create an index for the keywords into Firebase. This project combines these skills together to realize the data management in this application. It also allows us to learn how HTML works. Although both of us have no experience working on user interfaces, we design our search and navigation program with the help of herokuapp. We are now able to navigate the tables by the foreign key relationships, and also improve the performance of the search program by creating a reference table to keywords. The experience of developing the web application let us understand how each data management technique works and how to apply them on the databases based on our needs.

Source:

<https://www.kaggle.com/abecklas/fifa-world-cup>

<http://files.grouplens.org/datasets/movielens/ml-latest-small-README.html>

<https://dev.mysql.com/doc/index-other.html>