

CSC343H1 F 20209: INTRODUCTION TO DATABASES

YouTube Trending Video Database Analysis

Chuyi Hou: 1004197834
Yuchen Tong: 1003534669

October, 2020



Computer Science
UNIVERSITY OF TORONTO



Division of Engineering Science
UNIVERSITY OF TORONTO

Contents

1	Dataset and Relational Schema (Phase 1)	1
1.1	Domain	1
1.2	Dataset	1
1.3	Investigative Questions	1
1.4	Schema	1
2	Schema Implementation (Phase 2)	2
2.1	SQL Schema	2
2.2	Import Schema Demo on SQL	3
3	Data Cleaning and Import (Phase 3)	6
3.1	Decisions	6
3.2	Cleaning Steps	6
3.3	create_Category_from_json.py	7
3.4	data_cleaning.py	8
3.5	Data Import Demo on SQL	9
4	Queries and Results (Phase 4)	10
4.1	Investigative Questions Revision	10
4.2	What we have learned	11

1 Dataset and Relational Schema (Phase 1)

1.1 Domain

We want to study the trends of Youtube videos from different countries.

1.2 Dataset

The dataset can be downloaded manually via <https://www.kaggle.com/datasnaek/youtube-new> or using API: `kaggle datasets download -d datasnaek/youtube-new`

The raw dataset consists of ten json files for category ids and ten csv files of YouTube trending videos from ten different countries. Before implementing the schema, we should create country IDs and convert publish date so that it can be compared with trending date, split and count numbers of tags for each video, and find category names for each category ID.

We will be using the following columns of the csv files:

video_id, trending_date, title, channel_title, category_id, publish_time, tags, views, likes, dislikes

Last but not the least, we should get a sense of what those tags stand for and what categories of contents they represent.

1.3 Investigative Questions

- Find top n^* videos that are simultaneously trending in different countries given a specific time interval (Report: vID,vTitle,countryName,categName,tags,time interval).
- Given a specific time interval, find the trending categories that at least n^* countries have in common. (Report: categName, Time interval)
- Find the all time top n^* tags for each country. (Report: tag, countryName)

∗: Means variable.

1.4 Schema

Country(countryID, countryName)

Video(vID, vTitle, categID, countryID)

Category(categID, categName)

Channel(channelTitle, vID)

Publish(vID, pubTime, tags, categID, channelTitle)

Trending(trendingDate, vID)

Popularity(vID, views, likes, dislikes)

Publish[vID] \subseteq Video[vID]

Channel[vID] \subseteq Video[vID]

Popularity[vID] \subseteq Publish[vID]

Video[countryID] \subseteq Country[countryID]

Publish[categID] \subseteq Category[categID]

Publish[channelTitle, vID] \subseteq Channel[channelTitile, vID]

2 Schema Implementation (Phase 2)

2.1 SQL Schema

```
1 drop schema if exists TrendingYouTube cascade;
2 create schema TrendingYouTube;
3 set search_path to TrendingYouTube;
4
5 create domain ctryID as smallint
6     check (value >= 1 and value <= 10);
7
8 create domain hits as bigint
9     check (value >= 0);
10
11 create domain ctryName as text
12     check (value in ('Canada', 'Germany', 'France',
13         'United Kingdom', 'India', 'Japan', 'Korea',
14         'Mexico', 'Russia', 'United States'));
15
16 create domain catID as smallint
17     check (value >= 1 and value <= 30);
18
19 create table country(
20     countryID ctryID primary key not NULL,
21     countryName ctryName unique not NULL
22 );
23
24 create table category(
25     categID catID primary key not NULL,
26     categName text unique not NULL
27 );
28
29 create table video(
30     vID character(11) primary key not NULL,
31     vTitle text not NULL,
32     categID catID not NULL,
33     countryID ctryID not NULL,
34     foreign key (countryID) references country
35
36 );
37
38 create table channel(
39     channelTitle text not NULL,
40     vID character(11) not NULL,
41     primary key (channelTitle, vID),
42     foreign key (vID) references video
43 );
44
45 create table publish(
46     vID character(11) primary key not NULL,
47     pubTime date not NULL,
48     tags text not NULL,
49     categID catID not NULL,
50     channelTitle text not NULL,
```

```

51     foreign key (categID) references category,
52     foreign key (channelTitle, vID) references channel
53 );
54
55 create table trending(
56     trendingDate date not NULL,
57     vID character(11) references video not NULL,
58     primary key (trendingDate, vID)
59 );
60
61 create table popularity(
62     vID character(11) primary key not NULL,
63     view hits not NULL,
64     likes hits not NULL,
65     dislikes hits not NULL,
66     foreign key (vID) references publish
67 );

```

2.2 Import Schema Demo on SQL

```

1  csc343h-houchuyi=> \i schema.ddl
2  DROP SCHEMA
3  CREATE SCHEMA
4  SET
5  CREATE DOMAIN
6  CREATE DOMAIN
7  CREATE DOMAIN
8  CREATE DOMAIN
9  CREATE TABLE
10 CREATE TABLE
11 CREATE TABLE
12 CREATE TABLE
13 CREATE TABLE
14 CREATE TABLE
15 CREATE TABLE
16 csc343h-houchuyi=> \d
17                               List of relations
18 Schema | Name | Type | Owner
19 -----+-----+-----+-----
20 trendingyoutube | category | table | houchuyi
21 trendingyoutube | channel | table | houchuyi
22 trendingyoutube | country | table | houchuyi
23 trendingyoutube | popularity | table | houchuyi
24 trendingyoutube | publish | table | houchuyi
25 trendingyoutube | trending | table | houchuyi
26 trendingyoutube | video | table | houchuyi
27 (7 rows)
28
29 csc343h-houchuyi=> \d category
30                               Table "trendingyoutube.category"
31 Column | Type | Collation | Nullable | Default
32 -----+-----+-----+-----+-----
33 categid | catid | | not null |
34 catename | text | | not null |

```

```

35 Indexes:
36     "category_pkey" PRIMARY KEY, btree (catid)
37     "category_catename_key" UNIQUE CONSTRAINT, btree (catename)
38 Referenced by:
39     TABLE "publish" CONSTRAINT "publish_catetid_fkey" FOREIGN KEY
40     (catid) REFERENCES category(catetid)
41
42 csc343h-houchuyi=> \d channel
43         Table "trendingyoutube.channel"
44         Column      |      Type      | Collation | Nullable | Default
45 -----+-----+-----+-----+-----
46 channeltitle | text          |           | not null |
47 vid          | character(11) |           | not null |
48 Indexes:
49     "channel_pkey" PRIMARY KEY, btree (channeltitle, vid)
50 Foreign-key constraints:
51     "channel_vid_fkey" FOREIGN KEY (vid) REFERENCES video(vid)
52 Referenced by:
53     TABLE "publish" CONSTRAINT "publish_channeltitle_fkey" FOREIGN KEY
54     (channeltitle, vid) REFERENCES channel(channeltitle, vid)
55
56 csc343h-houchuyi=> \d country
57         Table "trendingyoutube.country"
58         Column      |      Type      | Collation | Nullable | Default
59 -----+-----+-----+-----+-----
60 countryid  | ctryid         |           | not null |
61 countryname | ctryname       |           | not null |
62 Indexes:
63     "country_pkey" PRIMARY KEY, btree (countryid)
64     "country_countryname_key" UNIQUE CONSTRAINT, btree (countryname)
65 Referenced by:
66     TABLE "video" CONSTRAINT "video_countryid_fkey" FOREIGN KEY (countryid)
67     REFERENCES country(countryid)
68
69 csc343h-houchuyi=> \d popularity
70         Table "trendingyoutube.popularity"
71         Column      |      Type      | Collation | Nullable | Default
72 -----+-----+-----+-----+-----
73 vid          | character(11) |           | not null |
74 view         | hits          |           | not null |
75 likes        | hits          |           | not null |
76 dislikes     | hits          |           | not null |
77 Indexes:
78     "popularity_pkey" PRIMARY KEY, btree (vid)
79 Foreign-key constraints:
80     "popularity_vid_fkey" FOREIGN KEY (vid) REFERENCES publish(vid)
81
82 csc343h-houchuyi=> \d publish
83         Table "trendingyoutube.publish"
84         Column      |      Type      | Collation | Nullable | Default
85 -----+-----+-----+-----+-----
86 vid          | character(11) |           | not null |
87 pubtime      | date          |           | not null |
88 tags         | text          |           | not null |

```

```

89  categid      | catid        |          | not null |
90  channeltitle | text         |          | not null |
91 Indexes:
92     "publish_pkey" PRIMARY KEY, btree (vid)
93 Foreign-key constraints:
94     "publish_categid_fkey" FOREIGN KEY (categid) REFERENCES category(categid)
95     "publish_channeltitle_fkey" FOREIGN KEY (channeltitle, vid) REFERENCES
96     channel(channeltitle, vid)
97 Referenced by:
98     TABLE "popularity" CONSTRAINT "popularity_vid_fkey" FOREIGN KEY (vid)
99     REFERENCES publish(vid)
100
101 csc343h-houchuyi=> \d trending
102           Table "trendingyoutube.trending"
103   Column      |      Type      | Collation | Nullable | Default
104 -----+-----+-----+-----+-----
105 trendingdate | date           |           | not null |
106 vid           | character(11)  |           | not null |
107 Indexes:
108     "trending_pkey" PRIMARY KEY, btree (trendingdate, vid)
109 Foreign-key constraints:
110     "trending_vid_fkey" FOREIGN KEY (vid) REFERENCES video(vid)
111
112 csc343h-houchuyi=> \d video
113           Table "trendingyoutube.video"
114   Column      |      Type      | Collation | Nullable | Default
115 -----+-----+-----+-----+-----
116 vid           | character(11)  |           | not null |
117 vtitle        | text           |           | not null |
118 categid       | catid          |           | not null |
119 countryid     | ctryid         |           | not null |
120 Indexes:
121     "video_pkey" PRIMARY KEY, btree (vid)
122 Foreign-key constraints:
123     "video_countryid_fkey" FOREIGN KEY (countryid) REFERENCES country(countryid)
124 Referenced by:
125     TABLE "channel" CONSTRAINT "channel_vid_fkey" FOREIGN KEY (vid)
126     REFERENCES video(vid)
127     TABLE "trending" CONSTRAINT "trending_vid_fkey" FOREIGN KEY (vid)
128     REFERENCES video(vid)

```

3 Data Cleaning and Import (Phase 3)

3.1 Decisions

Before we execute the data cleaning and import process, we decided to make some modifications to our schema and choose a narrowed selection of the YouTube trending dataset. For the changes made in schema, since countryID for videos is representing that video being trending in that country, we changed:

$$Video(\underline{vID}, vTitle, cateID, countryID) \rightarrow Video(\underline{vID}, vTitle, cateID)$$

and

$$Trending(trendingDate, \underline{vID}) \rightarrow Trending(trendingDate, \underline{vID}, countryID)$$

Moreover, we needed to change the foreign key accordingly:

$$Video[countryID] \subseteq Country[countryID] \rightarrow Trending[countryID] \subseteq Country[countryID]$$

For the new dataset, it only included three countries instead of 10 as initially proposed. They are Canada, the United States, and the United Kingdom. Notice that they are all English-speaking countries, hence we can access the dataset in a more understandable manner in terms of the output result that we are going to write queries on.

Next, in the process of data cleaning and import, we created two pieces of python code. *create_Category_from_json.py* is responsible for extracting category ids and their corresponding names from json files, and *data_cleaning.py* is to extract relevant data from each country's csv file. During data cleaning, we encountered an issue with load some countries' csv files. This might due to the csv decoder being unable to process some foreign languages' characters. This issue was avoided since we now only considered three English-speaking countries.

3.2 Cleaning Steps

- For *create_Category_from_json.py*
 - Read all json files which contains category ID and category name mapping for different countries.
 - Create a overall key(Category ID) and value(Category name) mapping by adding each item from previous mapping to a python dictionary.
 - If there are exist same category name mapped by different category ID in different country, suffix a number(how many times the name is occuring) to avoid duplicates.
- For *data_cleaning.py*
 - We have 3 csv files containing YouTube trending video data for 3 countries respectively.
 - First, we create the Country.csv by assigning 1 to 3 to these 3 countries respectively, and then add a countryID column to the country's csv file (i.e. for Canada, its csv file will be added a countryID column with all ones).
 - Then, combine all countries's csv by stacking one on up of another and select (project) relevant columns. (notice that there could be cases where one video can be trending in multiple countries, hence duplicates might occur in the combined csv).
 - Finally, load the combined csv (3 countries data are in side with relevant columns) and extract columns for SQL tables accordingly. Drop duplicates based on keys that were defined in the schema before write to files.
 - Output cleaned csv files: {Video.csv, Channel.csv, Publish.csv, Trending.csv, Popularity.csv}
- Country.csv, with columns being countryID (1 to 3) and countryName,
- All cleaned files are: {Country.csv, Category.csv, Video.csv, Channel.csv, Publish.csv, Trending.csv, Popularity.csv}

3.3 create_Category_from_json.py

```
1 import json
2 import pandas as pd
3
4 filenames = ['CA_category_id','DE_category_id','FR_category_id','
    ↪ GB_category_id', 'IN_category_id','JP_category_id','KR_category_id',
    ↪ 'US_category_id']
5
6 for f in filenames:
7
8     with open('./dataset/'+f+'.json') as file:
9
10         data_dict = json.load(file)
11
12         category_id = {}
13
14         for i in range(0, len(data_dict['items'])):
15             id = data_dict['items'][i]['id']
16             category = data_dict['items'][i]['snippet']['title']
17
18             category_id[id] = category
19
20 df = pd.DataFrame(category_id.items(),columns=['categID','categName'])
21
22 df.to_csv('Category.csv',index=False)
```

3.4 data_cleaning.py

```
1 import pandas as pd
2
3 # Country.csv and Category.csv are made manually
4
5 all_filenames = ['CAvideos', 'DEvideos', 'FRvideos', 'GBvideos', 'INvideos', '
    ↳ JPvideos', 'KRvideos', 'USvideos']
6 # 'MXvideos' 'RUvideos'
7 # first add one additional column 'countryID' to each csv files
8 countryID = 1
9 print('Start adding countryID column to each countries csv files')
10 for file in all_filenames:
11     df = pd.read_csv('./dataset/'+file+'.csv',
12                     usecols = ['video_id', 'trending_date', 'title', '
    ↳ channel_title', 'category_id', 'publish_time', '
    ↳ tags', 'views', 'likes', 'dislikes'])
13     df["countryID"] = countryID
14     df.to_csv('./dataset/' + file + 'ID.csv', index=False)
15     countryID += 1
16 print('Adding succesfully')
17
18 print('Start Combining all countries csv files')
19 # next, we combine all csv files
20 combined_csv = pd.concat([pd.read_csv('./dataset/' + f + 'ID.csv') for f
    ↳ in all_filenames])
21 #export to csv
22 combined_csv.to_csv('./dataset/combined.csv', index=False)
23 print('Successfully Combined')
24
25
26 # then project desired columns and write
27 # define columns for each csv files
28 data = {'Video': ['video_id', 'title', 'category_id', 'countryID'],
29         'Channel': ['channel_title', 'video_id'],
30         'Publish': ['video_id', 'publish_time', 'tags', 'category_id', '
    ↳ channel_title'],
31         'Trending': ['trending_date', 'video_id'],
32         'Popularity': ['video_id', 'views', 'likes', 'dislikes']
33     }
34
35 cols = [['vID', 'vTitle', 'categID', 'countryID'],
36         ['channelTitle', 'vID'],
37         ['vID', 'pubTime', 'tags', 'categID', 'channelTitle'],
38         ['trendingDate', 'vID'],
39         ['vID', 'views', 'likes', 'dislikes']]
40
41 keys = ['video_id', ['channel_title', 'video_id'], 'video_id', ['
    ↳ trending_date', 'video_id'], 'video_id']
42 i = 0
43 print('Start Extracting Columns From the Combined csv to Each Tables
    ↳ Needed for SQL Import')
44 for table in data:
45     df = pd.read_csv('./dataset/combined.csv')
```

```

46     df = df[data[table]]
47
48     # sorting by key
49     df.sort_values(keys[i], inplace = True)
50
51     # dropping ALL duplicate rows
52     df.drop_duplicates(subset = keys[i], keep = False, inplace = True)
53
54     # modify date datastyle for Trending table
55     if table == 'Trending':
56         dates = df.trending_date
57
58         new_col = []
59         for date in dates:
60             ymd = date.split('.')
61             new_col.append('20'+ymd[0]+'-'+ymd[1]+'-'+ymd[2])
62
63         new = pd.DataFrame(new_col, columns=['trendingDate'])
64
65         df.drop(columns=['trending_date'])
66
67     # changing columns using .columns()
68     df.columns = cols[i]
69
70     df.to_csv(table+'.csv', index=False)
71     i+=1
72 print('All Tables are Successfully Created')

```

3.5 Data Import Demo on SQL

4 Queries and Results (Phase 4)

4.1 Investigative Questions Revision

- Find top n^* videos that are simultaneously trending in different countries given a specific time interval (Report: vID,vTitle,categName,tags,time_interval,number_of_trendings).
- Find the monthly (from 2017-11 to 2018-06) top trending category that at least n^* countries have in common. (report: year, month, categName, number_of_trendings)
- Find channels that are trending in every month within a time interval (e.g. 2017-11 to 2018-03), and whose monthly total view is non-decreasing. (report: channelname)

*: Means variable.

4.2 What we have learned