

VE472 LAB3

Wu Jiayao 517370910257

3. Verifying the data

Schema

```
1  create table name
2  (
3      nconst varchar(10) not null primary key,
4      primaryName text not null,
5      birthYear varchar(4) not null,
6      deathYear varchar(4),
7      primaryProfession text not null,
8      knownForTitles text not null
9  );
10
11 create table title
12 (
13     tconst varchar(10) not null primary key,
14     titleType varchar(64) not null,
15     primaryTitle text not null,
16     originalTitle text not null,
17     isAdult boolean not null,
18     startYear varchar(4) not null,
19     endYear varchar(4),
20     runtimeMinutes integer not null,
21     genres text not null
22 );
23
24 create table principal
25 (
26     tconst varchar(10) not null,
27     ordering integer not null,
28     nconst varchar(10) not null,
29     category text not null,
30     job text,
31     characters text
32 );
33 create table rating
34 (
35     tconst varchar(10) not null primary key,
36     averageRating double not null,
37     numVotes integer not null
38 );
39
40 .separator "\t"
41 .import name.basics.tsv name
42 .import title.basics.tsv title
43 .import title.principals.tsv principal
44 .import title.ratings.tsv rating
```

the oldest movie

```
1 select tconst, primaryTitle, startYear from title
2 where startYear <> "\N" and titleType = "movie"
3 order by startYear
4 limit 1;
```

1	tconst	primaryTitle	startYear
2	-----	-----	-----
3	tt2210499	Birmingham	1896

the longest in 2019

```
1 select tconst, primaryTitle, startYear, runtimeMinutes from title
2 where startYear = "2009" and runtimeMinutes <> "\N" and titleType = "movie"
3 order by runtimeMinutes desc
4 limit 1;
```

1	tconst	primaryTitle	startYear	runtimeMinutes
2	-----	-----	-----	-----
3	tt1458948	Native of Owhyhee	2009	390

The year with the most movies

```
1 select startYear, count(*) as count from title
2 where startYear <> "\N" and titleType = "movie"
3 group by startYear
4 order by count desc
5 limit 1;
```

1	startYear	count
2	-----	-----
3	2019	14205

The name of the person who contains in the most movies

```
1 select principal.nconst, name.primaryName, count(*) as contained from
principal
2 inner join title on principal.tconst = title.tconst
3 inner join name on principal.nconst = name.nconst
4 where title.titleType = "movie"
5 group by principal.nconst
6 order by contained desc
7 limit 1;
```

1	nconst	primaryName	contained
2	-----	-----	-----
3	nm0006137	Ilaiyaraaja	952

the principal crew

```
1 select tconst,nconst, name.primaryName, category from principal
2 left outer join name using (nconst)
3 where principal.tconst in
4 (
5     select rating.tconst from rating
6     where rating.numVotes > 500
7     order by rating.averageRating
8     desc
9     limit 1
10 );
```

	tconst	nconst	primaryName	category
1	tt11128054	nm1249167	Kiril Spaseski	production_designer
2	tt11128054	nm8262223		actor
3	tt11128054	nm9819225		actor
4	tt11128054	nm1426138		actor
5	tt11128054	nm0953471	Ivan Zaric	actor
6	tt11128054	nm0804080	Slobodan Skerl	director
7	tt11128054	nm5173970	Ivana Mikovic	producer
8	tt11128054	nm0559105	Mate Matic	composer
9	tt11128054	nm3409959	Maja Radosevic	cinematographer
10	tt11128054	nm5275099		production_designer

the count of each pair

```
1 select birthYear, deathYear, count(*) as pairCount from name
2 where birthYear <> "\N" and deathYear <> "\N"
3 group by birthYear, deathYear
4 order by pairCount desc;
```

5. Advanced

```
1 create table name_profession
2 (
3     nconst varchar(10) not null,
4     profession text not null,
5     foreign key(nconst) references name(nconst)
6 );
7
8 create table IF NOT EXISTS title_genre
9 (
10     tconst varchar(10) not null,
11     genre text not null,
12     foreign key(tconst) references title(tconst)
13 );
```

```
1 import sqlite3
2 conn = sqlite3.connect('var/imdb.sqlite3')
3
4 c = conn.cursor()
```

```

5 names = c.execute("select * from name")
6 insert_list = []
7 for name in names:
8     if name[4] != '\\N' and len(name[4])>0:
9         nconst = name[0]
10        professions = name[4].split(',')
11        for _profession in professions:
12            insert_list.append((nconst,_profession))
13 c.executemany("insert into name_profession values (?,?)",insert_list)
14
15 titles = c.execute("select * from title")
16 insert_list = []
17 for title in titles:
18     if title[8] != '\\N' and len(title[8])>0:
19         tconst = title[0]
20         genres = title[8].split(',')
21         for _genre in genres:
22             insert_list.append((tconst,_genre))
23 c.executemany("insert into title_genre values (?,?)",insert_list)
24 conn.commit()
25 conn.close()

```

top 3 common profession

```

1 select profession, count(*) as count, avg(deathYear - birthYear) as
   avgLifeSpan from name_profession
2 inner join name using (nconst)
3 where deathYear <> "\\N" and birthYear <> "\\N"
4 group by profession
5 order by count desc
6 limit 3;

```

1	profession	count	avgLifeSpan
2	-----	-----	-----
3	actor	60186	70.0475525869804
4	writer	30931	71.8120655652905
5	actress	26260	73.4029702970297

The top 3 most popular (received most votes) genres

```

1 select genre, sum(numVotes) as votes from title_genre
2 inner join rating using (tconst)
3 group by genre
4 order by votes desc
5 limit 3;

```

1	genre	votes
2	-----	-----
3	Drama	466974607
4	Action	302900340
5	Comedy	286806624

The average time span (endYear - startYear) of the titles for each person

```
1 select principal.nconst, name.primaryName, avg(title.endYear -  
   title.startYear) as avgTimeSpan from principal  
2 inner join title using (tconst)  
3 inner join name using (nconst)  
4 where title.startYear <> "\N" and title.endYear <> "\N"  
5 group by principal.nconst  
6 order by avgTimeSpan desc;
```