# ECE4721J - Lab 3 Report

Methods and Tools for Big Data

Kexuan Huang

June 6, 2022

## 3. Verifying the Data

### The oldest movie

Query:

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

Output:

```
1  primaryTitle  startYear
2  ------------  ---------
3  Birmingham    1896
```

### The longest movie in 2009

Query:

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

Output:

```
1  primaryTitle      runtimeMinutes
2  ----------------  --------------
3  Native of Owhyhee  390
```

### The year with the most movies

Query:

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

Output:

```
1  startYear   count
2  ---------   -----
3  2021        15898
```

### The name of the person who contains in the most movies

Query:

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

Output:

```
1  primaryName   contained
2  -----------   ---------
3  Ilaiyaraaja   949
```

**The principal crew of the movie with highest average ratings and more than 500 votes**

Query:

```
1  select name.primaryName,
2      principal.category
3  from name,
4      principal
5  where name.nconst = principal.nconst
6      and principal.tconst in (
7          select rating.tconst
8          from rating
9          where rating.numVotes > 500
10         order by rating.averageRating desc
11         limit 1
12     );
```

Output:

```
1  primaryName       category
2  ---------------   --------
3  Melanie Zanetti   actress
4  David McCormack   actor
5  Joe Brumm         writer
6  David Barber      composer
```

**The count of each `Pair<BirthYear, DeathYear>` of the people**

Query:

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

Output:

> Too long, see `query.out`

## 4. Interaction with SQLite in Java / Python

> Please refer to `insert.py`

## 5. Advanced Analysis with the new Tables

**The top 3 most common professions among these people and also the average life span of these three professions**

Query:

```sql
 1  select profession,
 2      count(*) as count,
 3      avg(deathYear - birthYear) as avgLifeSpan
 4  from name,
 5      name_profession
 6  where name.nconst = name_profession.nconst
 7      and deathYear <> "\N"
 8      and birthYear <> "\N"
 9  group by profession
10  order by count desc
11  limit 3;
```

Output:

```
 1  profession   count    avgLifeSpan
 2  ----------   ------   ----------------
 3  actor        126066   70.0966160582552
 4  writer       64452    71.9769440824179
 5  actress      55228    73.5529803722749
```

**The top 3 most popular (received most votes) genres**

Query:

```
1  select genre,
2      sum(numVotes) as votes
3  from rating, title_genre
4  where rating.tconst = title_genre.tconst
5  group by genre
6  order by votes desc
7  limit 3;
```

Output:

```
1  genre    votes
2  ------   ---------
3  Drama    532586552
4  Action   355071204
5  Comedy   326066948
```

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

Query:

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

Output:

Too long, see `query.out`