

# CSCI620 Phase 2 : IMDB Dataset

Ketan Kokane  
Rochester Institute of  
Technology  
1 Lomb Memorial Dr,  
Rochester, NY 14623  
kk7471@rit.edu

Siddarth Sargunaraj  
Rochester Institute of  
Technology  
1 Lomb Memorial Dr,  
Rochester, NY 14623  
sxs2469@rit.edu

Ameya Nagnur  
Rochester Institute of  
Technology  
1 Lomb Memorial Dr,  
Rochester, NY 14623  
an4920@rit.edu

Kavya Kotian  
Rochester Institute of  
Technology  
1 Lomb Memorial Dr,  
Rochester, NY 14623  
kk2014@rit.edu

## ABSTRACT

This *project* analyses the dataset used by IMDb. There is an ER diagram representing the entities constituting the dataset and relationships between them. Based on this ER diagram, we build a Relational Schema that helps us write queries to create tables for this database. We improve the database with indexing and create queries to provide efficient access of data to the users. Made use of python to translate tsv files to sql

## Keywords

Data Management, ER Diagram, Relational Scheme, Functional Dependencies, Indexing

## 1. WHAT HAS BEEN DONE?

- Analysis of the entities in the IMDB dataset
- Convert the dataset entities to an ER diagram
- Improve upon the ER Diagram
- Create a Relational Schema
- Found out the functional dependencies in the dataset and to normalize the data
- SQL scripts to create and populate tables
- Test the schema performance with sql queries

## 2. ER DIAGRAM

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Group 14 CSCI620

Copyright 2019 ACM X-XXXXX-XX-X/XX/XX .

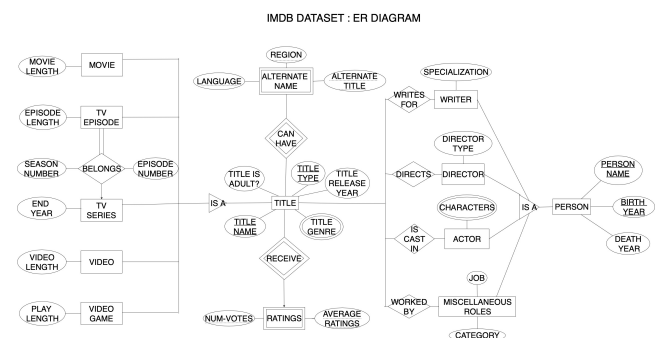


Figure 1: ER Diagram

## 3. ANALYSIS OF THE ENTITIES IN THE IMDB DATASET

The entities defined in the dataset are Title, Title Genres, Movie, Video, Video games, TV Series, TV Episode, Alternate Name, Ratings, Person, Writer, Director, Actor, Character, Miscellaneous Roles.

Every title has a Title name and Title type which are unique for every Movie, Video, Video games, TV Series and TV Episode. A title can belong to multiple Genres. A title can belong to multiple Genres and notifies if the content is adult. **IMDb provides data about Movie, Video, Video games, TV Series, TV Episode which are instances of the Title entity.**

Movie is a Title and gives the Movie Length. Video is a Title with an additional attribute Video Length. Video Games is a Title with an additional attribute Video Length. Every TV Episode belongs to a TV Series and a TV Series has multiple TV Episodes. Alternate Names and Ratings depend on Title and their attributes serve as Discriminators. **IMDb provides data about every cast members of every title**

Every Person has a Person name and Birth Year which are unique for every Person, Writer, Director, Actor and Miscellaneous Roles. A Person entity also tells us about the persons Death Year.

## 4. RELATIONAL SCHEMA

**Title** (Title\_Name, Title\_Type, Title\_Release\_Year, Title\_Is\_Adult?)

**PK:** (Title\_Name, Title\_Release\_Year, Title\_Type)

**Title\_genres** (Title\_Name, Title\_Type, Genre)

**PK:** (Title\_Name, Title\_Type, Genre)

**FK:** (Title\_Name, Title\_Type)

**Movie** (Title\_Name, Title\_Type, Movie\_Length)

**PK:** (Title\_Name, Title\_Type)

**FK:** Title\_Name, Title\_Type)

**Video** (Title\_Name, Title\_Type, Video\_Length)

**PK:** (Title\_Name, Title\_Type)

**FK:** (Title\_Name, Title\_Type)

**Video\_Games** (Title\_Name, Title\_Type, Play\_Length)

**PK:** (Title\_Name, Title\_Type, Play\_Length)

**FK:** (Title\_Name, Title\_Type)

**TV\_Series** (Title\_Name, Title\_Type, End\_Year)

**PK:** (Title\_Name, Title\_Type)

**FK:** Title\_Name, Title\_Type)

**TV\_Episode** (Title\_Name, Title\_Type, Episode\_Length)

**PK:** (Title\_Name, Title\_Type)

**FK:** Title\_Name, Title\_Type)

**Belongs** (TV\_Series\_Name, TV\_Series\_Type, TV\_Episode\_Name, TV\_Episode\_Type, Season\_Number, Episode\_Number)

**FK\_TV\_Episode** = (TV\_Series\_Name, TV\_Series\_Type)

**FK\_TV\_Series** = (TV\_Episode\_Name, TV\_Episode\_Type)

**Alternate\_Name** (Title\_Name, Title\_Type, Region, Language, Alternate\_Title)

**PK:** (Title\_Name, Title\_Type, Region, Language, Alternate\_Title)

**FK:** (Title\_Name, Title\_Type)

**Ratings** (Title\_Name, Title\_Type, Num\_Votes, Average\_Ratings)

**PK:** (Title\_Name, Title\_Type, Num\_Votes, Average\_Ratings)

**FK:** (Title\_Name, Title\_Type)

**Person** (Person\_Name, Birth\_Year, Death\_Year)

**PK:** (Person\_Name, Birth\_Year)

**Writer** (Person\_Name, Birth\_Year, Specialization)

**FK:** (Person\_Name, Birth\_Year)

**Writes\_For** (Title\_Name, Title\_Type, Person\_Name, Birth\_Year)

**FK\_Person:** (Person\_Name, Birth\_Year)

**FK\_title:** Title\_Name, Title\_Type)

**Director** (Person\_Name, Birth\_Year, Director\_type)

**FK:** (Person\_Name, Birth\_Year)

**Directs** (Title\_Name, Title\_Type, Person\_Name, Birth\_Year)

**FK\_Person:** (Person\_Name, Birth\_Year)

**FK\_title:** Title\_Name, Title\_Type)

**Actor** (Person\_Name, Birth\_Year)

**FK:** (Person\_Name, Birth\_Year)

**Character** (Person\_Name, Birth\_Year, Character)

**FK:** (Person\_Name, Birth\_Year)

**Is\_Cast\_In** (Title\_Name, Title\_Type, Person\_Name, Birth\_Year)

**FK\_Person:** (Person\_Name, Birth\_Year)

**FK\_title:** Title\_Name, Title\_Type)

**Miscellaneous\_Roles** (Person\_Name, Birth\_Year, Category, Job)

**FK:** (Person\_Name, Birth\_Year)

**Worked\_By** (Title\_Name, Title\_Type, Person\_Name, Birth\_Year)

**FK\_Person:** (Person\_Name, Birth\_Year)

**FK\_title:** (Title\_Name, Title\_Type)

## 5. FUNCTIONAL DEPENDANCIES

### 5.1 Assumptions

- (Title\_Name, Title\_Type ) be represented as  $A$
- (Title\_Name, Title\_Type ) for TV\_Series be represented as  $A_{series}$
- (Title\_Name, Title\_Type ) for TV\_Episodes be represented as  $A_{episodes}$
- (Person\_Name, Birth\_Year) be represented as  $B$

### 5.2 Dependencies

$A \rightarrow$  Title\_Release\_Year, Title\_Is\_Adult?

$A \rightarrow$  Movie\_Length

$A \rightarrow$  Video\_Length

$A \rightarrow$  Play\_Length

$A \rightarrow$  End\_Year

$A \rightarrow$  Episode\_Length

$A_{series} \rightarrow A_{episodes}$

$A_{episodes} \rightarrow$  Season\_Number, Episode\_Number

$B \rightarrow$  Death\_Year

$B \rightarrow$  Specialization

$B \rightarrow$  Director\_type

$B \rightarrow$  Character

B → Category, Job

## 6. DATABASE POPULATION

Attached *create\_all.sql* and *populate\_all.sql* file which are generated using the python script which parses the entire data-set of IMDB and samples 10000 titles from and loads the dataset.

## 7. PYTHON SCRIPT TO READ AND LOAD THE DATASET

Attached *code.py* file which uses pandas framework to analyze the dataset and based on the functional dependencies generates the sql files.

## 8. QUERIES

### 8.1

#### 8.1.1 Description

The query creates a join of the tables movie and Titlegenres to get the movie titles in the genre 'Comedy'.

#### 8.1.2 Query

```
SELECT M.titlename,
       G.genre
FROM   movie M,
       titlegenres G
WHERE  M.titlename = G.titlename
       AND G.genre LIKE 'Comedy';
```

#### 8.1.3 Output

TitleName	Genre
41 el hombre perfecto	Comedy
A Bizarre Love Triangle	Comedy
Adventures of Death	Comedy
All Wrong	Comedy
America's Most Retarded	Comedy

(few output values)

#### 8.1.4 Query justification

Joins 2 tables based on titlename i.e. an attribute with most unique entries in the database and fetches information based on genres.

*Duration/FetchTime : 0.359sec/0.000054sec*

### 8.2

#### 8.2.1 Description

The query creates a join of the tables belongs (Maps TvSeries name as S\_TitleName with it's episode names as E\_TitleName), Titleratings and tvepisode to get the total number of votes for tv series based on the ratings for it's episodes.

#### 8.2.2 Query

```
SELECT B.s_titlename,
       Sum(R.numvotes)
FROM   belongs B,
       titleratings R,
```

```
       tvepisode E
WHERE  B.e_titlename = E.titlename
       AND E.titlename = R.titlename
GROUP BY B.s_titlename;
```

#### 8.2.3 Output

S_TitleName	sum(R.NumVotes)
Ferris Bueller	8
Neighbours	8

#### 8.2.4 Query justification

Joins 3 tables again based on titlename but only of episodes and tv series. We use the group by clause to get episodes for every tv series and sum their votes making it more complex than simply fetching votes or ratings for a title.

*Duration/FetchTime : 0.041sec/0.000011sec*

### 8.3

#### 8.3.1 Description

The query creates a join of the tables titleAlternateName and Title to get the alternate titles for any Movie, TvSeries or other type of titles for a particular region. ("India" here)

#### 8.3.2 Query

```
SELECT T.TitleName,A.alternate_title
FROM   titlealternatename A,
       title T
WHERE  T.titlename = A.titlename
       AND A.region = 'India';
```

#### 8.3.3 Output

TitleName	alternate_title
There & Back	Kalia

#### 8.3.4 Query justification

Joins 2 tables and uses a table that provides region specific information. The join is done again on titlename for any type of title.

*Duration/FetchTime : 0.0034sec/0.0000088sec*

### 8.4

#### 8.4.1 Description

- The query is a nested query with the inner query joining the tables Titleratings, director and title\_director (Maps director names with title names they have worked for) to fetch the average ratings for all titles for which director named Aaron Kodz has worked.
- The outer query fetches all the directors and the titles they have worked for if the title has an average rating > than all ratings fetched by inner query.

#### 8.4.2 Query

```
SELECT R.titlename,
       D.personname,
       R.averageratings
FROM   titleratings R,
       director D,
       title_director TD
```

```

WHERE TD.titlename = R.titlename
AND D.personname = TD.personname
AND R.averageratings
    >ALL (SELECT TR.averageratings
        FROM titleratings TR,
            director DR,
            title_director TDR
        WHERE TR.titlename = TDR.titlename
            AND DR.personname = TDR.personname
            AND DR.personname = 'Aaron Kodz');

```

The IMdb dataset contains data about 12 different categories of titles, which can be relased in multiple regions in different languages. Maps persons who works in the movie.

### 8.4.3 Query justification

Has a nested query. Inner query joins 3 tables based on 2 different attributes that fetches average rating based on titles worked on by a specific director. Outer query joins the same three tables and checks for a value larger than all fetched by inner query. Equivalent to a nested loop in any programming language.

*Duration/FetchTime : 0.214sec/0.286sec*

### 8.4.4 Output

TitleName	PersonName	AverageRatings
111 Girls	Bijan Zmanpira	74
111 Girls	Nahid Ghobadi	74
41 el hombre perfecto	Pepe Romy	12
A Aa E Ee	Sabapathy	26
A Bizarre Love Triangle	Mu-yeong Lee	131

(few output values)

## 8.5 Query justification summary

We have tested queries with complexities of joining 2 or more tables, with a nested structure and with group by. We have covered all key tables hence covering a major chunk of the data ranging from a simple fetch from a 2 table join to a comparison with all entries of a table returned by an inner query. The remaining tables like Writer can easily be switched with Director like in query 4 but the result would be similar. Hence, all types of data is covered and less likely but complex queries are tested along with simple ones.

## 9. INDEXES

Used Cluster Index generated by the MySQL as the only indexes for the database, as design of the database allwoed to query the majority of the data just by using cluter index.

## 10. WORK DISTRIBUTION

Task	Handled By
Analysis of the entities in the IMDB dataset	Ameya, Kavya
Convert the dataset entities to an ER diagram	Ketan, Siddarth
Improve upon the ER Diagram	Ketan, Siddarth
Create a Relational Schema	Ameya, Kavya
Found out the functional dependencies in the dataset and to normalize the data	Ameya, Kavya
SQL scripts to create and populate tables	All
Test the schema performance with sql queries	All

## 11. CONCLUSION

Analyzed the imdb dataset and learned how to model the entities in the dataset, relate two entities using relation, mapping the cardinalities of the relations.