

## CITP – Project Portfolio Source data

The data provided for the CITP project portfolio includes a subset of the Internet Movie Database IMDb's publicly available dataset, a supplementary dataset collected from the Open Movie Database OMDb, as well as a preliminary inverted word index covering these data. Find description and instructions on how to access and load the data into your own database below.

### The IMDb dataset

The IMDb dataset is available in the PostgreSQL database backup file (dump) **imdb.backup**. It is a subset of around 130000 of the movies in IMDb's publicly available dataset (<https://www.imdb.com/interfaces/>). To create a database on your local PostgreSQL database server and load the data into that new database, you can simply issue the following two psql commands in your Windows CMD or Mac Terminal command prompt (assuming that your current directory is where the file is located)<sup>1</sup>:

```
psql -U postgres -c "create database imdb"
psql -U postgres -d imdb -f imdb.backup
```

The result of loading the **imdb.backup** data in this way will be a new database, called **imdb** with the 7 tables shown in figure 1. The schema / table structure corresponds closely to the 7 tsv-files (tab-separated-values files) publicly available as IMDb's dataset (<https://www.imdb.com/interfaces/>). Observe that a **movie** (and an episode in a series) **is** in IMDb's terminology **called a title**. To prepare we created a database with the 7 tables, loaded data from the tsv-files into these, reduced the content to cover around 130.000 movies and dumped the reduced database to the file **imdb.backup**. The content of the reduced dataset is fine for our purposes and is easier to work with than the full dataset. The 7-table database is not in good shape. A first problem is: it calls for a thorough redesign to become a well-designed relational database. During the reduction of the data, we started with the title\_basics and aimed at reducing the other tables correspondingly, so that for instance all titles included in title\_ratings are also included in title\_basics. However, during this process we did not change any column values. So, a second problem is: some columns may contain inconsistent data due to dangling references.

The description of the table columns below is a slightly modified version of the description of the tsv-files at the IMDb download page.

---

<sup>1</sup> If you alternatively (or in combination) prefer to use the remote server on cit.ruc.dk, you can use the following command (replace citXX with your group cit01, cit02, cit03 ...):

```
psql -h cit.ruc.dk -p 5432 -U citXX -W -f imdb.backup
```

On cit.ruc.dk you don't have privilege to create databases, so you must use your default database, citXX.

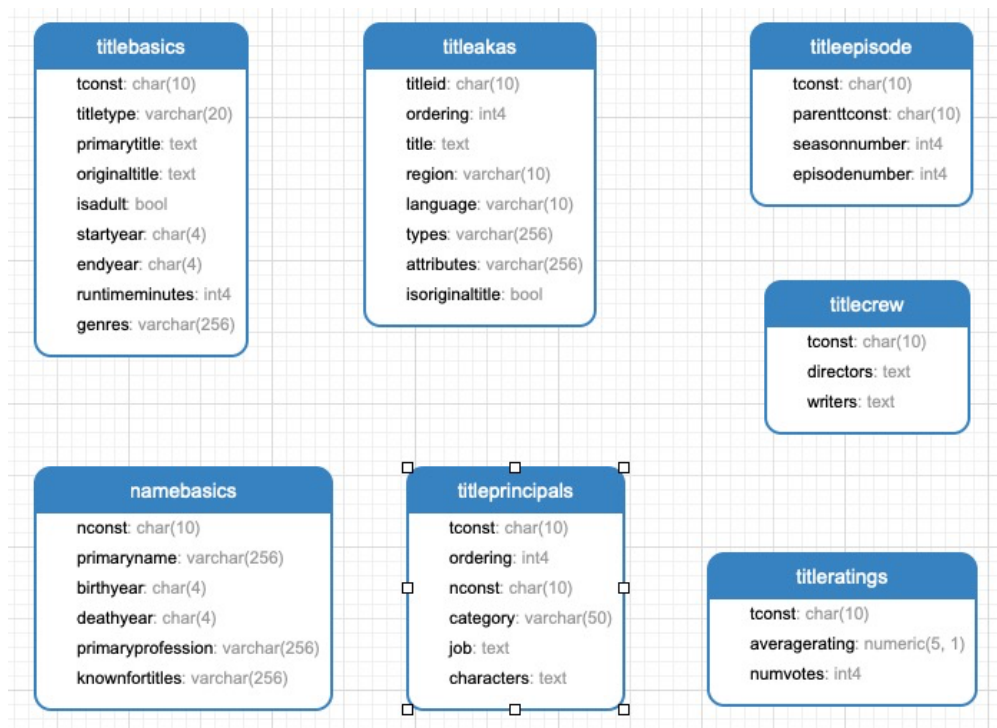


Figure 1: Schemas for the 7 IMDb dataset tables included in *imdb.backup*

**title\_akas** - Contains the following information for titles:

- ❑ titleId (string) - a tconst, an alphanumeric unique identifier of the title
- ❑ ordering (integer) - a number to uniquely identify rows for a given titleId
- ❑ title (string) - the localized title
- ❑ region (string) - the region for this version of the title
- ❑ language (string) - the language of the title
- ❑ types (strings) - Enumerated set of attributes for this alternative title. One or more of the following: "alternative", "dvd", "festival", "tv", "video", "working", "original", "imdbDisplay".
- ❑ attributes (strings) - Additional terms to describe this alternative title, not enumerated
- ❑ isOriginalTitle (boolean)

**title\_basics** - Contains the following information for titles:

- ❑ tconst (string) - alphanumeric unique identifier of the title
- ❑ titleType (string) - the type/format of the title (e.g. movie, short, tvseries, tvepisode, video, etc)
- ❑ primaryTitle (string) - the more popular title / the title used by the filmmakers on promotional materials at the point of release
- ❑ originalTitle (string) - original title, in the original language
- ❑ isAdult (boolean) - 0: non-adult title; 1: adult title
- ❑ startYear (YYYY) - represents the release year of a title. In the case of TV Series, it is the series start year
- ❑ endYear (YYYY) - TV Series end year. Only specified for TV Series title types
- ❑ runtimeMinutes (integer) - primary runtime of the title, in minutes
- ❑ genres (strings) - includes up to three genres associated with the title

**title\_crew** Contains director and writer information for all the titles in IMDb. Fields include:

- ❑ tconst (string) - alphanumeric unique identifier of the title

- ❑ directors (strings / nconsts) - director(s) of the given title
- ❑ writers (strings / nconsts) – writer(s) of the given title

**title\_episode** – Contains the tv episode information. Fields include:

- ❑ tconst (string) - alphanumeric identifier of episode
- ❑ parentTconst (string) - alphanumeric identifier of the parent TV Series
- ❑ seasonNumber (integer) – season number the episode belongs to
- ❑ episodeNumber (integer) – episode number of the tconst in the TV series

**title\_principals** – Contains the principal cast/crew for titles

- ❑ tconst (string) - alphanumeric unique identifier of the title
- ❑ ordering (integer) – a number to uniquely identify rows for a given titleId
- ❑ nconst (string) - alphanumeric unique identifier of the name/person
- ❑ category (string) - the category of the job that the person was in
- ❑ job (string) - the specific job title if applicable
- ❑ characters (string) - the name of the character played if applicable

**title\_ratings** – Contains the IMDb rating and votes information for titles

- ❑ tconst (string) - alphanumeric unique identifier of the title
- ❑ averageRating – weighted average of all the individual user ratings
- ❑ numVotes - number of votes the title has received

**name\_basics** – Contains the following information for names:

- ❑ nconst (string) - alphanumeric unique identifier of the name/person
- ❑ primaryName (string)– name by which the person is most often credited
- ❑ birthYear – in YYYY format
- ❑ deathYear – in YYYY format if applicable
- ❑ primaryProfession (strings)– the top-3 professions of the person
- ❑ knownForTitles (strings/tconsts) – titles the person is known for

### The OMDb dataset

IMDb does not include all their data in the public dataset, but some of what they don't disclose can be found elsewhere. We will use supplementary data downloaded from the OMDb API (<http://www.omdbapi.com>). This data is partly overlapping with what can be found in the IMDb dataset, but some useful unique additions are included.

To simplify the inclusion of the OMDb data in your relational database, the data has been transformed and made available as a single table **omdb\_data** that can be imported (see figure 2). When building your own database, you can select and include columns of your own preference from this table. Only two are required: **poster** and **plot** providing respectively a link to the main poster for the movie and a detailed description of the plot for the title.

**omdb\_data** – supplementary data from OMDb

- ❑ poster - link to the main poster
- ❑ plot – a detailed description of the plot for the title
- ❑ ... and many more.

The data is provided as a PostgreSQL database backup file **omdb\_data.backup** and the **omdb\_data** table will be imported into your **imdb** database when you issue the command<sup>2</sup>:

```
psql -U postgres -d imdb -f omdb_data.backup
```

Notice that **tconst** is a column in **title\_basics** as well as in **omdb\_data**, so you can use this to combine the data. More than 26.000 titles in **title\_basics** don't have a corresponding row in **omdb\_data** (data for these titles is not available in OMDb).

omdb_data	wi
tconst: char(10)	tconst: char(10)
episode: varchar(80)	word: text
awards: varchar(80)	field: char(1)
plot: text	lexeme: text
seriesid: varchar(80)	
rated: varchar(80)	
imdbrating: varchar(80)	
runtime: varchar(80)	
language: varchar(200)	
released: varchar(80)	
response: varchar(80)	
writer: text	
genre: varchar(80)	
title: varchar(256)	
country: varchar(256)	
dvd: varchar(80)	
production: varchar(80)	
season: varchar(80)	
type: varchar(80)	
poster: varchar(180)	
ratings: varchar(180)	
imdbvotes: varchar(100)	
boxoffice: varchar(100)	
actors: varchar(256)	
director: text	
year: varchar(100)	
website: varchar(100)	
metascore: varchar(100)	
totalseasons: varchar(100)	

Figure 2: The source tables **omdb\_data** and **wi** available in **omdb\_data.backup** and **wi.backup**

<sup>2</sup> If you are using the remote server on cit.ruc.dk, you can use the following command (replace citXX with your group cit01, cit02, cit03 ...):

```
psql -h cit.ruc.dk -p 5432 -U citXX -W -f omdb_data.backup
```

### A supplementary word index for textual data

To support selected queries on textual data a so-called inverted index is provided. The inverted index is a word index that can be used to lookup a word and retrieve the documents (movies) that this word appears in. The word index data is provided in a single table **wi** shown in figure 2 and described below. The indexed columns are primarytitle, plot, characters and primaryname (indicated with the letters 't', 'p', 'c' and 'n' respectively).

- wi** – inverted word index on columns primarytitle, plot, characters and primaryname
- ❑ tconst (string) - alphanumeric unique identifier of the title
- ❑ word – the word
- ❑ field – the column where the word occurs, values are 't', 'p', 'c' and 'n' that indicates primarytitle, plot, characters and primaryname respectively
- ❑ lexeme – a lexeme derived for the word. Will be null for fields 'c' and 'n'.

Most important in this table are the columns tconst and word and it's fine just to use these. The field column can be used if you need to be selective regarding what columns to match when using the inverted index. Maybe you want to remove all the 'c' and 'n' words because you find them disturbing. You could also consider replacing them with "words" that are full primarynames / full character names. The lexeme column can be used to provide a search where different forms of words (inflexions) are harmonized by the use of a separate lexeme index. A lexeme is a string, just like a word, but it provides a kind of normalized form, where different inflexion forms, such as aim, aime, aimed, aiming and aims, normalize to a single lexeme, like aim.

To import the wi table into your own imdb database, download the file **wi.backup** and run the command<sup>3</sup>:

```
psql -U postgres -d imdb -f wi.backup
```

---

<sup>3</sup> On the remote server on cit.ruc.dk, the command is:  
`psql -h cit.ruc.dk -p 5432 -U citXX -W -f wi.backup`