



UNIVERSITÀ DI PISA

Large-Scale and Multi-structured Databases -
Project: Rotten Movies

Fabio Piras, Giacomo Volpi, Guillaume Quint

Contents

| | | |
|----------|---------------------|----------|
| 1 | Introduction | 2 |
|----------|---------------------|----------|

Chapter 1

Introduction

Rotten Movies is a web service where user can keep track of the general liking for movies, they can also share their thoughts with the world after signing up. In addition user can follow renowned top critic to be constantly updated with their latest review.

Guest user, as well as all other, can browse or search movies based on filters like: title, year of release and actor who worked in it. They can also view a Hall of Fame for the most positive review genres, production houses and years in which the best movie were released.

For this project, the application has been developed in Java Spring using IntelliJ to implement a web GUI. The application use a document database to store the main information about the movies, user, review and actor; a graph database is instead used to keep track of the relationship between normal user and top critic in terms of who follows who and between the movie and the user who reviewed it.

Chapter 2

Feasibility Study

The very first step to perform is to look up for what type of data we need to create the application and how to handle it by performing a feasibility study.

2.1 Dataset analysis and creation

Initially we search online for available dataset, we ended up with five different files coming from Kaggle and imdb with a combined storage of 4.18 Gb:

- title_principal.tsv
- name_basic.tsv
- title_basic.tsv
- title_crew.tsv
- rotten_movies.csv
- rotten_reviews.csv

The first three were found on the imdb site containing general data about the title of the entries, type (not all were movies), cast, crew and other useful information. The last two came from Kaggle and they contain data scraped from the rotten tomatoes site regarding the movies rating and their reviews. All of these dataset were organized in a relational manner.

The next step was to delete all useless information and perform a join operation to generate one single final dataset for the document database. These operations were performed through a python script with the help of the pandas library and Google Colab, the end result was a single json file of 270 Mb. The structure of the document is the following:

```
1 {  
2   "_id" : ObjectId(<<id_field>>),
```

```

3      "primaryTitle": "The_first_ever_movie",
4      "year": 1989,
5      "runtimeMinutes": 70,
6      "genres": ["Crime", "Drama", "Romantic"],
7      "productionCompany": "Dingo_Picture_Production" ,
8      "personnel": [
9          {
10             "name": "John_Doe",
11             "category": "producer",
12             "job": "writer"
13         },
14         {
15             "name": "Christopher_Lee",
16             "category": "actor",
17             "character": ["The_one"]
18         },
19         ...
20     ],
21     "top_critic_fresh_count": "4",
22     "top_critic_rotten_count": 0,
23     "user_fresh_count": 14,
24     "user_rotten_count": 1,
25     "top_critic_rating": 100,
26     "top_critic_status": "Fresh",
27     "user_rating": 93,
28     "user_status": "Fresh",
29     "review":
30     [
31         {
32             "critic_name": "AntonyE",
33             "review_date": "2018-12-07",
34             "review_type": "Rotten",
35             "top_critic": false,
36             "review_content": "I_really_didn't_liked_it!",
37             "review_score": "1/10"
38         },
39         {
40             "critic_name": "AntonyE",
41             "review_date": "2015-12-07",
42             "review_type": "Fresh",
43             "top_critic": true,
44             "review_content": "I_really_liked_it!",
45             "review_score": "A-"
46         },
47         ...
48     ]
49 }

```