

Task 2 – Movie Database

Design Document

Federico Fregosi, Mirko Laruina,
Riccardo Mancini, Gianmarco Petrelli

May 12, 2020

Contents

1	Specifications	1
1.1	Application Overview	1
1.2	Actors	1
1.3	Requirement Analysis	1
2	Design	2
2.1	Use-case diagram	2
2.2	Class diagram	2
2.3	Data model	5
2.4	Software Architecture	7
3	Database	8
3.1	Creation of the database	8
3.2	Scraping	8
3.3	Updating the database	8
4	Backend Implementation	8
4.1	Java classes overview	8
4.2	Authentication	8
4.3	Movies	9
4.4	Ratings	13
4.5	Users	15

1 Specifications

1.1 Application Overview

The application is an aggregator of movies and movie ratings with the purpose of providing logged users statistics and informations about a large set of movies. Logged-in user can also rate movies they have watched while not logged-in users may still use the service to browse movie rankings and statistics but they are not able to give their rate. Only movies released in Italy are considered.

All users can search a movie and view its details (e.g., title, original title, duration, cast, ...) along with its average rating from users and from external sources.

In addition, all users can browse the list of movies sorting and filtering it by many parameters (e.g. year, genre, country, actors, ...).

System administrators can view all user profile pages and ban users. In order to do that, he can check the full history of ratings. Once a user is banned, he can no longer log in and his username and email cannot be used by new users.

The movie database will be built upon the publicly available IMDb dataset.

The ratings will be gathered also by periodically scraping external websites (e.g., Rotten Tomatoes, Coming Soon, MyMovies).

1.2 Actors

Anonymous user, registered user, administrator and updater “bot”.

1.3 Requirement Analysis

1.3.1 Functional Requirements

An **anonymous user** must be able to register in order to become a *registered user*. Login is carried out using username and password selected by the user when registering. Username must be unique. A valid email address is also required in order to register. An email cannot be used more than once.

Both **anonymous user** and **registered user** must be able to:

- view details and average rating of a specific movie
- view a list of movies and filter it by many parameters. Combined filters are also allowed
- view aggregated statistics about movies: the user can choose on which field to aggregate movies (year, country, actor, director, genre) and additional filters (like the movie browsing feature). E.g. the user might want to see the ranking of the countries with the best movies in the last 10 years.

A **registered user** must be able to rate a movie, in addition to what anonymous user can do. A registered user must also be able to manage his profile. In the profile a registered user can:

- check, add and modify his personal data
- browse the history of his rates

- view aggregated statistics about his profile (i.e. most viewed genre, most recurrent actor, etc...) based on his rated movies
- delete the account

Finally, a registered user can logout in any moment.

An **administrator** is a special registered user who must be able to ban users. In order to do that, an administrator can check a global rating history to retrieve information about all the application's activity, and to check every user's profile. Banned user's rating are automatically removed from the database. Email and username of banned users cannot be used again.

The **updater "bot"** is not a real user but an entity used to periodically update the database in order to add new movies and update external ratings.

1.3.2 Non-Functional Requirements

- **Availability:** the Database must be replicated in order to be always available. Write operations on the Database can be eventually consistent.
- **Scalability:** the application must be able to scale to an arbitrary number of servers.
- **Security:** passwords must be stored in a secure way.
- **Responsive UI:** Client-side application must provide a responsive view both for pc, laptops and mobile devices.

2 Design

2.1 Use-case diagram

The use-case diagram is shown in Figure 1. Different colors are used to highlight cases that are exclusive of some actors: white cases are referred to all users; blue cases are referred to registered user and admin; yellow cases are exclusive of the admin.

2.2 Class diagram

The class diagram is shown in Figure 2. It was decided to show *Country*, *Year* and *Genre* as separate entities as they are some of the fields over which aggregate statistics are calculated.

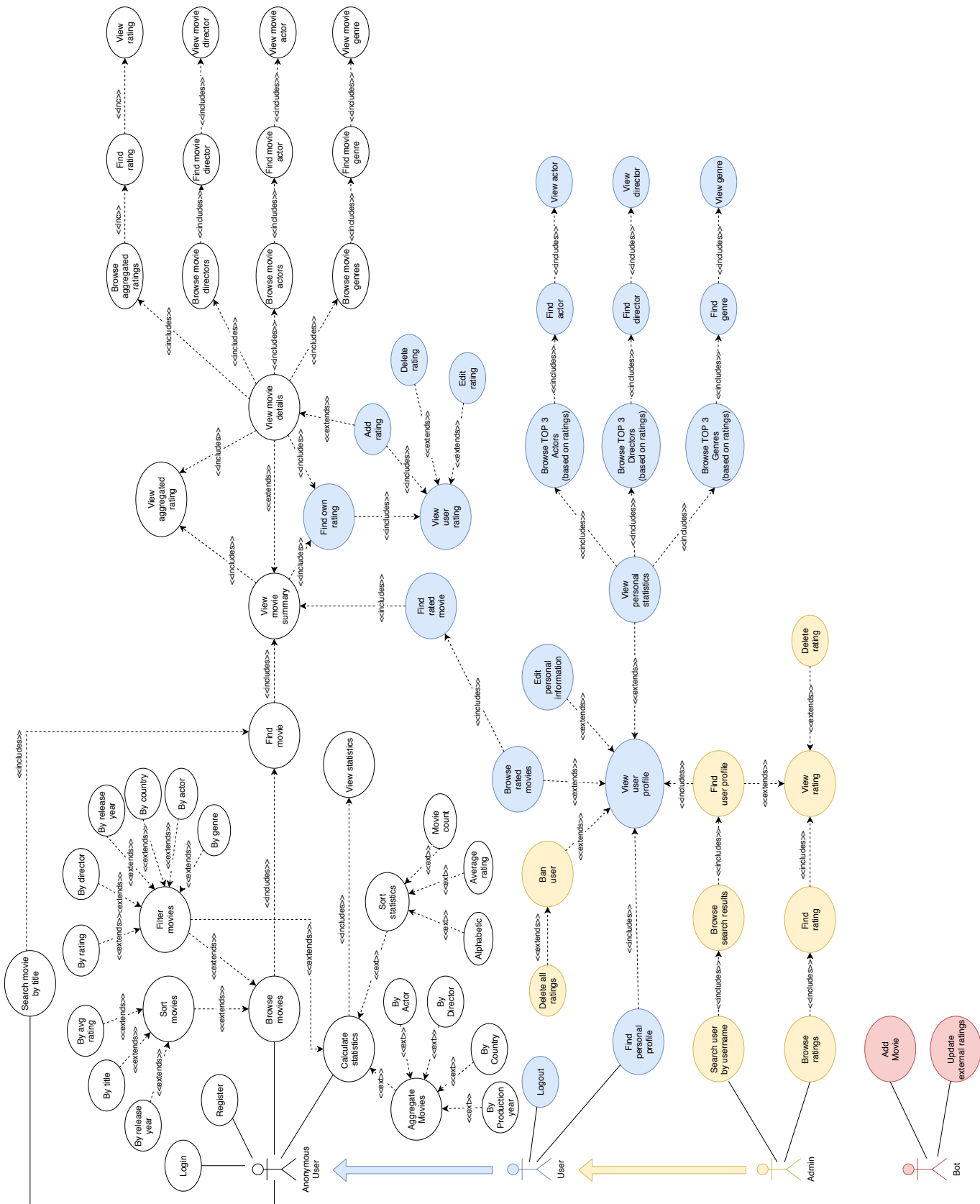


Figure 1: Use-case diagram

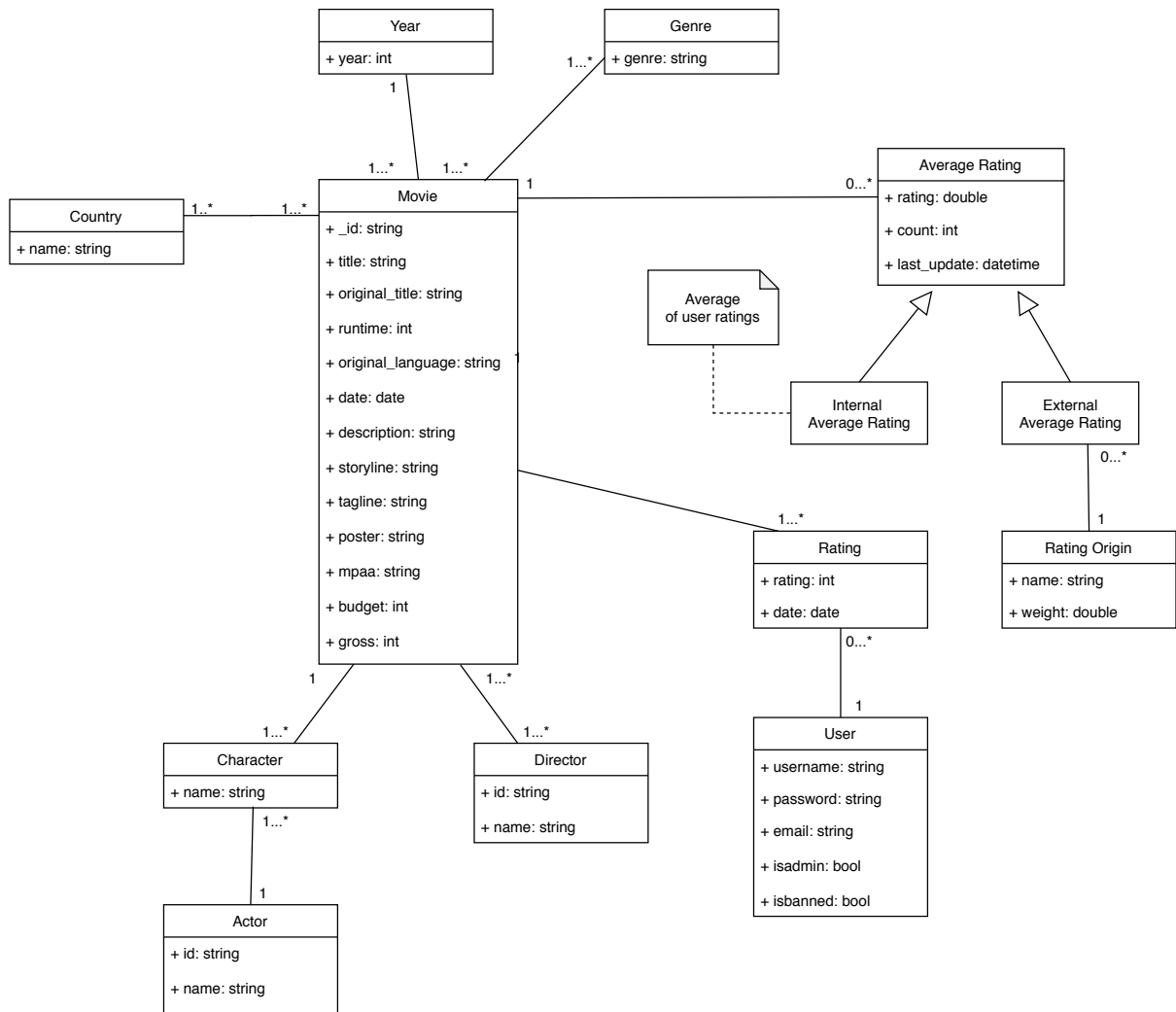


Figure 2: Class diagram for the identified entities

2.3 Data model

The data model is split in 3 different collections:

- Movies (Section 2.3.1)
- Users (Section 2.3.3)
- Ratings (Section 2.3.2)

In each following subsection, an example document is shown for every collection.

2.3.1 Movies

```
1 {
2   "_id": "tt7286456",
3   "title": "Joker",
4   "original_title": "Joker",
5   "runtime": 122,
6   "countries": ["USA", "Canada"],
7   "original_language": "English",
8   "year": 2019,
9   "date": "2019-10-04",
10  "description": "In Gotham City, mentally troubled comedian [...]",
11  "storyline": "Joker centers around an origin of the iconic arch [...]",
12  "tagline": "Put on a happy face.",
13  "poster": "https://m.media-amazon.com/images/M/[...].jpg",
14  "mpaa": "Rated R for strong bloody violence, disturbing behavior, [...]",
15  "budget": 55000000,
16  "gross": 1074251311,
17  "characters": [
18    {
19      "name": "Joker",
20      "actor_name": "Joaquin Phoenix",
21      "actor_id": "nm0001618"
22    },
23    ...
24  ],
25  "directors": [
26    {
27      "id": "nm0680846",
28      "name": "Todd Phillips",
29    }
30  ],
31  "genres": ["Crime", "Drama", "Thriller"],
32  "ratings": [
33    {
34      "source": "internal",
35      "avgrating": 9,
36      "count": 100,
37      "weight": 2,
38      "last_update": "2019-11-28 12:34:56"
39    },
40    {
41      "source": "IMDb",
42      "avgrating": 8.6,
43      "count": 628981,
```

```

44         "weight": 1,
45         "last_update": "2019-11-28 12:34:56"
46     },
47     ...
48 ],
49 "total_rating": 8.87,
50 "last_scraped": "2019-11-28 12:34:56"
51 }

```

Notes Character and director nested documents contain redundant data (person's name) that is introduced to reduce the number of seeks necessary to return the movie details to the user.

This collection is mainly read-heavy since most fields will be written only once. The only fields subject to change are the ones related to ratings. For these reasons, we can afford many indices to speed-up the aggregations.

Indices Indices on the following fields will be set for this collection:

- title: speed-up sorting by title
- title, original_title (text index): movie lookup by name
- countries (multikey index): filter and aggregation by country
- genres (multikey index): filter and aggregation by genre
- year: filter and aggregation by year
- characters.actor_id (multikey index): filter and aggregation by actor
- directors.id (multikey index): filter and aggregation by director
- total_rating: filter and sorting by rating
- last_scraped: fast lookup for next-to-be-scraped movie (oldest movie).

Among this fields, *total_rating* is the only one that will be updated frequently (all other fields with index are updated with the periodic scraper whose period is of some days). However it is necessary to efficiently fulfill many queries.

2.3.2 Ratings

```

1 {
2   "_id": {
3     "user_id": <ObjectId>,
4     "movie_id": "tt7286456"
5   },
6   "date": "2020-01-20",
7   "rating": 10
8 }

```

Notes The raw ratings are kept separated from the movie they refer for three reasons:

1. single ratings are never shown to the user, except when the user looks his own ratings or when the admin looks all ratings

2. they may be accessed per-user, per-movie (for calculating statistics) and globally (by the administrator). The optimization for one use-case would influence the performance of another.
3. new ratings are expected to be frequently added and, thus, the use of a nested document inside the *movies* collection would be infeasible because it would grow indefinitely.

However, this choice mandates the introduction of indices for efficiently obtaining ratings of a user or ratings of a movies are required.

Indices

- `_id.movie_id`: speed up aggregations on movie.
- `_id.user_id`: speed up user ratings browse.
- `date`: speed up browsing user and global ratings.

2.3.3 Users

```
1 {
2   "_id": <ObjectId>,
3   "username": "joker",
4   "password": "<HASHED PASSWORD>",
5   "email": "joker@dccomics.com",
6   "isadmin": True (optional),
7   "isbanned": False (optional),
8   "sessions": [
9     {
10      "session_id": "<session_id>",
11      "expiry": "2020-01-28"
12    },
13    ...
14  ]
15 }
```

Notes *Username, email and session_id* must be unique across the whole collection.

Indices

- `username` (text index): find user by username
- `username`: find user by exact username
- `email`: find user by email
- `sessions.session_id` (multikey index): find user by session id

2.4 Software Architecture

The application will be made of the following 4 components:

- **Mongo DB**: a MongoDB cluster will be deployed with replication.
- **React Frontend**: web-based UI.

- **Java Backend:** using *Spring*, the *Java backend* will provide REST APIs to the *React frontend*.
 - **Updater “bot”:** two Python scripts are needed in order to nightly update the DB with the latest movies and ratings:
 1. the **IMDB parser** periodically parses the IMDb dataset to add the latest movies.
 2. the **scraper** continuously parses the rating sources to update the ratings.
- These scripts will executes asynchronously from the *Java backend*.

3 Database

3.1 Creation of the database

The database is created from scratch starting from the public IMDb dataset available at the link <https://www.imdb.com/interfaces/>.

We developed a Python script, in order to obtain the data in the form described in the Data Model. However, some information where missing. So we integrated missing information scraping them from other movie’s sites while we were scraping for ratings.

Collections are saved into Json files and then upserted in MongoDB database with another simple python script with the use of the **UpdateOne** module. Update requests are inserted in append into an array and served through the **bulk_write()** function.

3.2 Scrapping

3.3 Updating the database

4 Backend Implementation

The API specification can be found in the docs/api.md file. This section will go through the implementation of the main APIs.

4.1 Java classes overview

TODO - Class diagram - Some words - Each class has an Adapter to conveniently convert from/to Mongo Document

4.2 Authentication

Register In order to register a user, it must first be checked that no user with same username and password exists (snippet 1) and then, if no duplicate entry is found, the user can be added to the database (snippet 2). If creation was successful, a new session for the new user is then added (snippet 5). Any subsequent API call must be accompanied by the session identifier in order to identify the user (snippet 6).

Change Password This query (snippet 3) is simply done by matching the user by user-name and setting the new password.

Login In order to authenticate the user, a matching user with same username and password is found in the database (snippet 4). If a user matches, a new session for the given user is then added (snippet 5).

Logout The session is removed from the DB (snippet 7).

```
1 // 1 - check user uniqueness
2 usersCollection.find(
3     or(
4         eq("username", u.getUsername()),
5         eq("email", u.getEmail())
6     )
7 ).first();
8
9 // 2 - insert user
10 usersCollection.insertOne(User.Adapter.toDBObject(u));
11
12 // 3 - edit password
13 usersCollection.updateOne(
14     eq("username", u.getUsername()),
15     set("password", u.getPassword())
16 );
17
18 // 4 - authenticate user
19 usersCollection.find(
20     and(eq("username", u.getUsername()),
21         eq("password", u.getPassword())
22     ))
23 .first();
24
25 // 5 - add session
26 usersCollection.updateOne(
27     eq("_id", u.getId()),
28     push("sessions", Session.Adapter.toDBObject(s))
29 );
30
31 // 6 - finds the user whose session is s
32 usersCollection.find(
33     and(
34         eq("sessions._id", s.getId())
35     ))
36 .first();
37
38 // 7 - remove session
39 usersCollection.updateOne(
40     eq("_id", u.getId()),
41     pull("sessions", eq("_id", s.getId()))
42 );
```

4.3 Movies

4.3.1 Browse

Returns a list of filtered movies with the given sorting. The result is paged.

1. build a list of filters based on user specifications. People names are matched if all words in the query string are substrings of the full name (this is done using a regex).
2. Make the filter to pass to the *find* function by and-ing all conditions.
3. Make the query:

find all movies matching the conditions

sort by the user-defined sort order

project only movie important details

skip the first *n* pages

limit the results to the page size

If user is logged in, a separate query will be done to fetch the user's ratings of the returned movies (4.4.2).

```
1 // 1 - define filters
2 List<Bson> conditions = new ArrayList<>();
3 if (minRating != -1)
4     conditions.add(gte("total_rating", minRating));
5
6 if (maxRating != -1)
7     conditions.add(lte("total_rating", maxRating));
8
9 if (director != null && !director.isEmpty()){
10     List<Bson> directorConditions = new ArrayList<>();
11     for (String s:director.split(" ")){
12         directorConditions.add(regex("directors.name", s, "i"));
13     }
14     conditions.add(and(directorConditions.toArray(new Bson[]{})));
15 }
16
17 if (actor != null && !actor.isEmpty()){
18     List<Bson> actorConditions = new ArrayList<>();
19     for (String s:actor.split(" ")){
20         actorConditions.add(regex("characters.actor_name", s, "i"));
21     }
22     conditions.add(and(actorConditions.toArray(new Bson[]{})));
23 }
24
25 if (country != null && !country.isEmpty())
26     conditions.add(eq("countries", country));
27
28 if (fromYear != -1){
29     conditions.add(gte("year", fromYear));
30 }
31
32 if (toYear != -1)
33     conditions.add(lte("year", toYear));
34
35 if (genre != null && !genre.isEmpty())
36     conditions.add(eq("genres", genre));
37
38 // 2 - and conditions
```

```

39 Bson filters;
40
41 if (!conditions.isEmpty())
42     filters = and(conditions.toArray(new Bson[]{}));
43 else
44     filters = new BsonDocument();
45
46 // 3 - make the query
47 FindIterable<Document> movieIterable = moviesCollection
48     .find(filters)
49     .sort(sorting)
50     .projection(include("title", "year", "poster", "genres", "total_rating", "
description"))
51     .skip(n*(page-1))
52     .limit(n);

```

4.3.2 Search

During a search, movies that (fuzzy) match the entire query are returned sorted by “likely-hood” (meta text score, in MongoDB). Results are paged as before.

If user is logged in, a separate query will be done to fetch the user’s ratings of the returned movies (refer to 4.4.2).

```

1 FindIterable<Document> movieIterable = moviesCollection
2     .find(text("\"+ query + "\""))
3     .projection(Projections.metaTextScore("score"))
4     .sort(Sorts.metaTextScore("score"))
5     .skip(n*(page-1))
6     .limit(n);

```

4.3.3 Get Details

Given a movie id, return all movie details. Code snippet not reported for brevity’s sake.

4.3.4 Browse Statistics

This function is used to return aggregated statistics on movies based on one of the following information attributes:

- genre
- year
- country
- director
- character

Before performing the aggregation, movies can be filtered as shown in the *GetMovieList* API. The statistics returned show the following informations about the aggregation attribute:

- aggregated attribute

- average rating calculated on the aggregation
- number of movies present in the aggregation (calculated after filtering)

Results can also be sorted in ascending or descending order based on each of the previous attributes.

The aggregation pipeline is performed with the following code:

```

1 AggregateIterable<Document> iterable = moviesCollection.aggregate(
2     Arrays.asList(
3         Aggregates.match(filters),
4         Aggregates.unwind("$" + realUnwindBy),
5         Aggregates.group(
6             "$" + realGroupBy,
7             Accumulators.first("name", "$" + realGroupName),
8             Accumulators.avg("avg_rating", "$total_rating"),
9             Accumulators.sum("movie_count", 1)
10        ),
11        Aggregates.sort(sorting),
12        Aggregates.skip(n*(page-1)),
13        Aggregates.limit(n)
14    );
15

```

The **match** function uses the array of conditions "*filters*" to perform all the filters in a single step. **Unwind** deconstructs the array field and returns a document for each array element (this is necessary for country, director and character fields). Then, the **group** function groups documents by the value expressed in the "*realGroupBy*" and then applies accumulator expression to each group:

- **Accumulators.first** sets "*realGroupName*" as name of the group
- **Accumulators.avg** calculates the average rating of all the "*total_rating*" and saves them in the "*avg_rating*" attribute
- **Accumulators.sum** calculates the number of movies in the aggregation, by adding 1 to the variable "*movie_count*" for each movie

After that, the results are sorted with the **sort** function, where the argument "*sorting*" is calculated as follow:

```

1 Bson sorting;
2 if (sortOrder == 1) {
3     sorting = ascending(realSortBy);
4 } else if (sortOrder == -1) {
5     sorting = descending(realSortBy);
6 } else {
7     throw new RuntimeException("sortOrder must be 1 or -1.");
8 }

```

Here, "*realSortBy*" is the value to sort.

Finally, the **skip** and the **limit** functions are used to manage the display of the results.

4.4 Ratings

4.4.1 CRUD operations

A rating can be added, read, updated or deleted. Related code snippets are not shown for brevity's sake.

After a rating is updated, the corresponding movie rating information are asynchronously updated (refer to 4.4.5).

4.4.2 Get all ratings for a user and list of movies

This particular operation is done to add the user's ratings to a query result set. First the list of movie ids of the movies in the query result is built and then all ratings of the given user to any of those movies is returned.

By building the array of identifiers, just one query to the database is necessary.

NB: the number of movies will always be small due to paging (20).

```
1 // 1 - build list of ids of the movies I'm interested in
2 List<String> ids = new ArrayList<>();
3 for (Movie m:movies){
4     ids.add(m.getId());
5 }
6
7 // 2 - fetch corresponding ratings from database
8 FindIterable<Document> ratingIterable = ratingsCollection.find(
9     and(
10         eq("_id.user_id", u.getId()),
11         in("_id.movie_id", ids)
12     )
13 );
```

4.4.3 Browse All

Admins can browse all ratings in the collection in descending date order. Results are paged as seen in previous sections. Related code snippet is not shown for brevity's sake.

4.4.4 Browse User Ratings

Users can browse all their ratings in descending date order. Results are paged as seen in previous sections. Related code snippet is not shown for brevity's sake.

4.4.5 Update Internal Rating

After a rating is changed, the average internal rating and the total rating of the movie should be updated. This is done asynchronously after the update in order not to make the client wait for this additional operation.

The update takes in input the old and the new rating (one of them can be null in case of insertion or deletion) and consists in three parts:

1. fetch the movie from the DB
2. calculate new statistics

3. update the database. Both internal rating and total rating are updated. For the former there are three cases:

- the internal rating was not present before and it needs to be added.
- the internal rating existed and needs to be updated
- the internal rating existed and needs to be removed

For the total rating there are two cases: update or unset in case the rating we removed was the last one.

Furthermore, this operation is also protected from possible concurrent bakends running in other machines since it modifies the database only if it detected no updated between the fetch of the movie and the updated. It does so by making sure that the *last_updated* field remains the same by adding a filter in the update query (if no internal rating was present before it just checks that it was not added in the meanwhile). In case of failure, it retries up to 5 times. It is very unlikely to fail but in case it will be corrected by the scraper that always builds the internal rating using the Ratings collection.

```
1 // 1 - fetch movie
2 Movie m = dba.getMovieDetails(movieId);
3
4 // 2 - calculate new aggregated ratings
5 // [...]
6
7 // 3 - update database
8 if (oldUpdate == null) { // there is no previous internal rating entry
9     match = and( // make sure no one added it in the meanwhile
10         eq("_id", movieId),
11         not(elemMatch("ratings", new Document().append("source", "internal")))
12     );
13     if (internalRating.getCount() != 0) { // push it if needs to be added
14         updates.add(push("ratings", AggregatedRating.Adapter.toDBObject(
15             internalRating)));
16     }
17 } else { // there is a previous entry: replace or pull
18     match = and( // make sure no one updated it in the meanwhile
19         eq("_id", movieId),
20         elemMatch("ratings", new Document()
21             .append("source", "internal")
22             .append("last_update", oldUpdate))
23     );
24     if (internalRating.getCount() != 0) { // replace it if it exists
25         updates.add(set("ratings.$", AggregatedRating.Adapter.toDBObject(
26             internalRating)));
27     } else { // else pull it
28         updates.add(pull("ratings", eq("source", "internal")));
29     }
30 }
31
32 if (m.getRatings().size() > 0) { // update total_rating
33     updates.add(set("total_rating", m.getTotalRating()));
34 } else { // no rating available: unset it
35     updates.add(unset("total_rating"));
36 }
```

```
35
36 UpdateResult result = moviesCollection.updateOne(match, combine(updates));
37 if (result.getModifiedCount() == 1)
38     // ok
39 else
40     // retry: someone modified it in the meanwhile
```

4.5 Users

4.5.1 Get Profile

4.5.2 Ban

4.5.3 Search