

Lab 5 – Counting Stars

CC5212-1

April 6, 2016

No more counting words!

Today we count the stars.

More specifically, we will count co-stars in movies that are listed in IMDb.¹ Your mission ... should you choose to accept it ... is to create a set of MapReduce jobs that count the pairs of actors or actresses that have co-starred in the most movies together. Hopefully it shouldn't be impossible.²

The goal of the lab is to demonstrate that you can now independently code MapReduce jobs to process large amounts of data. Since we did something very similar on the board in the lecture on Monday and since last Wednesday, we already saw the syntax and steps needed to code and run a MapReduce job, in this lab I'm not going to give you detailed step-by-step instructions. Instead, I'll give you the instructions to get you started and will point you to the material that will help you complete the lab. The rest will be up to you.

- The instructions for logging into the server, for accessing HDFS, for running Hadoop code, etc. are in last weeks' instructions: <http://aidanhogan.com/teaching/cc5212-1-2016/lab/04/mdp-1ab04.pdf>.
- The data you need are on HDFS in the `/uhadoop/shared/imdb/` folder. The `imdb-stars.tsv` file is about 1GB and contains 13 million roles played by different actors in IMDb (pretty much the full database). Have a look inside the start of the file.³ The header (not included in the file) is as follows:

Star Name	Movie Name	Year	Movie Number	Movie Type	Episode Name	Starring As	Role
...

Columns are tab delimited. Some values may not apply and may be `null`.

Star Name is the name of the star. The file is sorted alphabetically so you might see stars with weird names in there at the start of the file.

Movie Name is the name of the movie or tv series, etc., that the star appears in.

Year is ... well yes, the year the movie was released.

Movie Number is used when a movie with the same name appears in the same year (e.g., <http://www.imdb.com/title/tt0801505/> is the second movie called “Crash” listed in 2004 so it will have `II` here). Often this will be `null` (if there was only one movie with that name in that year).

Movie Type is the type of movie ... if it's a theatrical movie, a TV movie, a TV series, etc.⁴

Episode Name is the name of an episode if it was a TV series.

Starring As is the name of the actor/actress in the credits.

Role is the character they played.

¹The Internet Movie Database ... <http://imdb.com/> ... I'm sure many of you know it?

²Okay, that “joke” was terrible.

³`hdfs dfs -cat [file] | more` pressing spacebar to scan, Ctrl+C to end.

⁴In the code later, the options will appear in `org.mdp.imdb.ActorMovieParser.MovieRole.MovieType`

We will mainly focus on the first five columns. The Star Name is unique for a person. **However, for movies, Movie Name is not unique ... only the combination of Movie Name, Year and Movie Number are enough together to uniquely identify a movie.** To make a unique key for a movie, you'll need to concatenate the values of these three columns. We only want to consider movies: specifically you should only consider input lines where the Movie Type is equal to THEATRICAL_MOVIE; if it is something else, skip it. Let's look for info about a specific actor: type `hdfs dfs -cat /uhadoop/shared/imdb/imdb-stars.tsv | grep -e "^Pacino, Al" | more`, flick through some results with spacebar, Ctrl+C to end.

- Download the code project from <http://aidanhogan.com/teaching/cc5212-1-2016/lab/05/mdp-lab05.zip> and open it in Eclipse. In this project there is a Main class, an example of the CitationCount for you to adapt. You do not need to do anything with Main ... feel free to ignore it.
- The exercise for today is to adapt the CitationCount example provided in the package to create a sequence of Hadoop classes to count the number of movies that each pair of actors/actresses have co-starred in. The output should look something like:

```
De Niro, Robert##Pacino, Al      x
Cage, Nicholas##Kidman, Nicole  y
...                               ...
```

... where x is the number of theatrical movies that both Robert De Niro and Al Pacino star in, etc. (## is just a delimiter; you can use a tab or anything unambiguous). Some tips:

- You will have to create your own classes today. Follow the CitationCount example.
- As mentioned previously, only include theatrical movies and make sure the movie key you use is unique (see above).
- You will need multiple jobs. It's okay to run them manually one-by-one. Check the output of one before moving onto the next. The output of one task becomes the input of the next.
- /uhadoop/shared/imdb/imdb-stars-100k.tsv contains a subset of 100,000 roles from the input file (which contains 13 million). You should use this for testing your code. When it works for this, you can try run the process over the full file.
- When you're calling your jar from the command line, the first argument is the name of the class with the main method you need. The second argument should be the location of the input file. The third argument should be the location of the output file.
- Do not copy the input data into your personal folder ... it's too big to make 30–35 copies. Just pass the input location /uhadoop/shared/imdb/imdb-stars.tsv. But do write the output to your personal folder on HDFS (e.g., /uhadoop/USERNAME/imdb/out1/).
- OPTIONAL: try using MapReduce to get the top-10 most frequent co-stars. You can do this by sorting everything and taking the head of the results, for example. Or you can create a reduce job that keeps a sorted list of the top 10 results it has seen on each machine and outputs only them.
- Add the new class files that you created to a ZIP and submit to u-cursos before Monday night.