

# CA4022 Analysis of MovieLens Dataset Using Hadoop's PIG and HIVE.

Seán Cummins

18404744

sean.cummins26@mail.dcu.ie

[Github](#)

October 2021

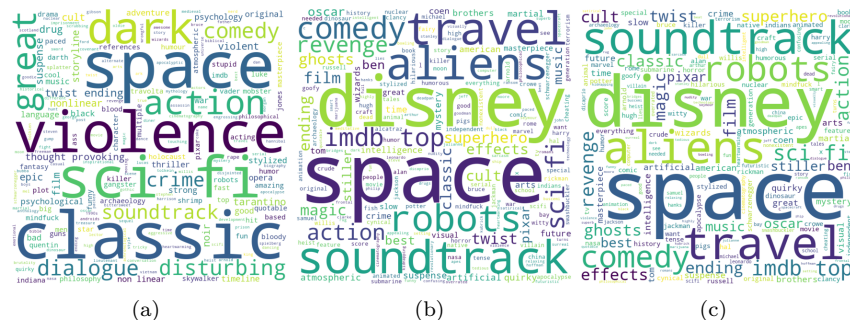
## Abstract

[Apache Hadoop](#) is a software framework that allows for the distributed processing of large data across clusters of computers. Rather than relying on single hardware architectures to perform large amounts of computation at constant high-availability, Hadoop clusters divide the task amongst individual hardware setups called nodes. These nodes are often consumer-grade and in-expensive. Along with being relatively cheap, Hadoop clusters ensure that data operations are not compromised by nodes that become unavailable. In a situation where a cluster node is down, Hadoop will redistribute the workload amongst other nodes ensuring no data loss.

The Hadoop infrastructure deploys the MapReduce algorithm (Dean and Ghemawat, 2008), coined by Google, to divide or 'Map' the task into identical, independent parts amongst nodes in a cluster. The outputs from clusters are then combined and aggregated, or 'Reduced' where possible. [Apache Hive](#) and [PIG](#) are frameworks built on top of the Hadoop infrastructure designed to refine the process of writing custom MapReduce jobs. MapReduce scripts in native Java are long, complex, and difficult to comprehend. Because of this, writing a custom MapReduce script is complicated. To this end, both Hive and PIG are capable of writing equivalent code with considerably less lines and complexity, overall increasing user efficiency and eliminating the need for programmers to write MapReduce scripts in Java.

## 1 Introduction

The following text aims to foster the analysis of the MovieLens dataset which consists of: movies, users, ratings, and tag information. In the following subsection, we will briefly introduce two frameworks commonly used to write and execute custom MapReduce scripts, Apache Hive, and Apache Pig. These frameworks are also used to perform data pre-processing and wrangling. In Section 2,



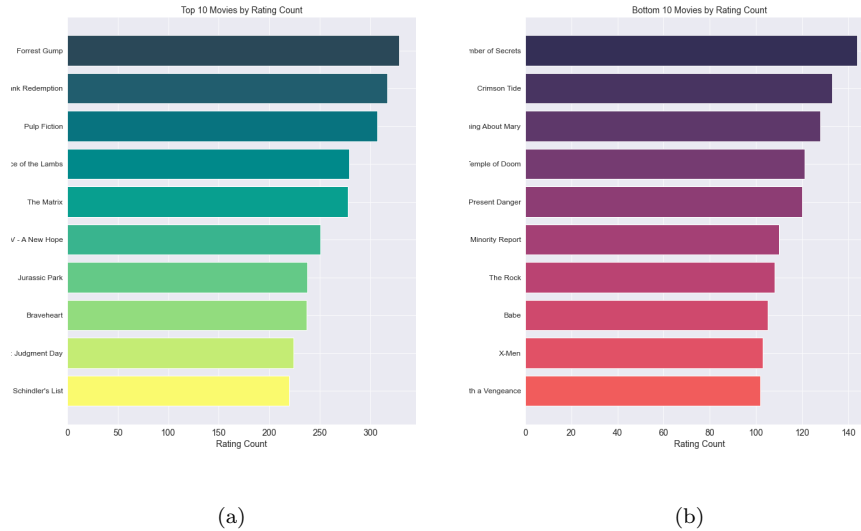


Figure 2: (a) Best Rated Movies (b) Worst Rated Movies

It is interesting to note that the tags 'violence', and 'dark' appear as a common word only in the set of highly rated movies. It is also no surprise that the highly rated movies are often tagged with 'classic'. 'Space' is common amongst all sub figures including (b) All Movies. This could indicate that the dataset is disproportionately representing space movies with respect to other genres. Or possibly, there has been a considerable amount of space movies produced over time. A final hypothesis could be that space movies are difficult to make and consequently, producers take on a lot of risk when they create a space film. This would imply that poorly produced space movies would receive negative criticism. However, when a space movie is well done, it receives high praise, resulting in its strong presence in both the sub figures (a), and (c).

In Figure 2, you will find bar charts representative of the movies included in Figure 1 (a) and (b). These movies are the films with the best and worst average ratings according to the MovieLens dataset. As well as this, they also receive the highest amount of ratings overall. The list of movies in (a) consists of the exact movies you would expect to find here. So called 'cult-classics'. It's interesting to note that movies from the 'worst movies' receive a comparable amount of tags to the 'best movies.' This shows that tags and descriptions are not reserved solely for movies that people enjoy.

### 3 Conclusion

We will now conclude this short review and analysis of Assignment one for CA4022. In this assignment, we were introduced to the Hadoop ecosystem, a powerful framework designed under the "Divide and Conquer" principles. The Hadoop architecture incorporates Google's MapReduce algorithm to divide big data tasks amongst clusters of machines rather than a single machine. This not only increases efficiency, but creates various safety mechanisms in the case of a system failure. Hadoop allows large difficult computations to be performed on a distributed system of machines that can be of consumer-grade. This eliminates the need for large, expensive, complex hardware.

We were also able to perform some light analysis on the MovieLens dataset. This analysis was shallow in the essence of keeping this review short. However, there is room to build further on this. Understanding the emotions portrayed by the tags users write to describe movies is a common research task known as Sentiment Analysis. It would be interesting to better understand the sentiment behind movie reviews in the future.

Please find the link to my Github containing all the PIG and Hive cleaning and analysis scripts below the title of this document. Otherwise, access my Github [here](#).

### References

- [1] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: Simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, jan 2008. ISSN 0001-0782. doi: 10.1145/1327452.1327492. URL <https://doi.org/10.1145/1327452.1327492>.
- [2] Ammar Fuad, Alva Erwin, and Heru Purnomo Ipung. Processing performance on apache pig, apache hive and mysql cluster. In *Proceedings of International Conference on Information, Communication Technology and System (ICTS) 2014*, pages 297–302. IEEE, 2014.
- [3] Dayong Du. *Apache Hive Essentials*. Packt Publishing Ltd, 2015.