

IT UNIVERSITY OF COPENHAGEN

Clustering Player Behaviors in Data Streams with K-Means in Map-Reduce

by

Sigurdur Karl Magnusson

A thesis submitted in partial fulfillment for the
degree of Master of Science in IT

in the
Computer Science
Software Development and Technology

Supervisors
Julian Togelius, Rasmus Pagh

April 2013

“Any intelligent fool can make things bigger, more complex, and more violent. It takes a touch of genius – and a lot of courage – to move in the opposite direction.”

Albert Einstein

Abstract

The abstract is a short summary of the thesis. It announces in a brief and concise way the scientific goals, methods, and most important results. The chapter “conclusions” is not equivalent to the abstract! Nevertheless, the abstract may contain concluding remarks. The abstract should not be discursive. Hence, it cannot summarize all aspects of the thesis in very detail. Nothing should appear in an abstract that is not also covered in the body of the thesis itself. Hence, the abstract should be the last part of the thesis to be compiled by the author.

A good abstract has the following properties: *Comprehensive*: All major parts of the main text must also appear in the abstract. *Precise*: Results, interpretations, and opinions must not differ from the ones in the main text. Avoid even subtle shifts in emphasis. *Objective*: It may contain evaluative components, but it must not seem judgemental, even if the thesis topic raises controversial issues. *Concise*: It should only contain the most important results. It should not exceed 300–500 words or about one page. *Intelligible*: It should only contain widely-used terms. It should not contain equations and citations. Try to avoid symbols and acronyms (or at least explain them). *Informative*: The reader should be able to quickly evaluate, whether or not the thesis is relevant for his/her work.

An Example: The objective was to determine whether ... (*question/goal*). For this purpose, ... was ... (*methodology*). It was found that ... (*results*). The results demonstrate that ... (*answer*).

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Contents

Abstract	ii
Acknowledgements	iii
List of Figures	vi
List of Tables	vii
Abbreviations	viii
1 Chapter Title Here	1
1.1 Motivation	1
1.2 Problem Statement	2
1.3 Contributions	2
1.4 Project Outline	2
2 Background Theory	3
2.1 Player Behaviour	3
2.1.1 Game Metric	3
2.1.2 Features	3
2.2 Clustering	3
2.2.1 K-Means	3
2.2.2 Streaming and Data Streams	3
2.3 Map-Reduce and Hadoop	3
3 Related Work	4
3.1 Clustering Player Behaviors	4
3.2 Clustering Large Data	4
3.2.0.1 Map-Reduce	7
3.2.1 Data Streams	7
4 Methodology	8
4.1 Game Events Data	8
4.1.1 Building Features	8
4.2 K-Means in Map-Reduce	8
4.2.1 Euclidean distance	8
4.2.2 Normalized Histograms	8
4.3 Experimental Set-Up	8

5	Results and Discussion	9
5.1	Results	9
5.2	Discussion	9
6	Conclusions	10
A	Appendix Title Here	11
	Bibliography	12

List of Figures

List of Tables

Abbreviations

LAH List Abbreviations **Here**

Chapter 1

Chapter Title Here

1.1 Motivation

We live in a time of data, where data is growing at an amazing rate. Data is everywhere around us and companies are generating more and more data. Storing data is cheaper than ever and many smaller and mid-sized companies are buying storage subscriptions at other bigger companies that offer "Storage as a Service" (SaaS). A SaaS company offers e.g. data reliability and durability by storing critical data in multiple facilities and on multiple devices. Having all this data does not give a meaning by itself, it needs to be analysed and interpreted to give a business value.

"Information is the oil of the 21st century, and analytics is the combustion engine."

Peter Sondergaard, senior VP at Gartner

The need to analyse all this data have caused many new companies to rise up that are specialised in analysing large quantities of data, extracting knowledge and converting it to a business value. Many of those companies offer "Analytics as a Service" (AaaS) that offer an analytical software to discover trends and unknown patterns in the data.

1.2 Problem Statement

GameAnalytics.com is a cloud hosted service for collecting, analyzing and reporting game metrics. Working with large quantities of game metric data that needs to be analysed and processed efficiently.

A valuable application for GA is to design and implement a scalable streaming version of a clustering algorithm. That can read a large data in mini-batches into memory and cluster the data incrementally.

The goal would be to incrementally find clusters efficiently in a streaming data, showing predominant characteristics of human behavior, e.g. the hardcore players, casual players, people who didn't understand the game, etc.

1.3 Contributions

- Clustering Player Behaviours using Map-Reduce framework
- Software and knowledge to GameAnalytics for further development

1.4 Project Outline

Chapter 2

Background Theory

2.1 Player Behaviour

2.1.1 Game Metric

2.1.2 Features

2.2 Clustering

2.2.1 K-Means

2.2.2 Streaming and Data Streams

2.3 Map-Reduce and Hadoop

Chapter 3

Related Work

In subsequent sections we will give a overview of some of the related and recent work. We start with looking into work related to Clustering Player Behaviors and then we dive into numerous work related to clustering large set of data using the k-means algorithm, where we can have a finite stream of data, an endless and evolving data stream and finally some parallel clustering implementations in the Map-Reduce framework. In our work we focus on processing large amount of data of finite length in parallel using Map-Reduce that arrives each day. The behavior of the data can change from day to day like when dealing with an endless stream of data but is not in the scope and is discussed in future work in the Conclusions section. [TODO].

3.1 Clustering Player Behaviors

Many researches have been done on clustering and predicting player behaviors over the years to get a better understanding which kind of groups of behaviors are playing a game by exploring user telemetry and player events that is being logged as a the game is being played.

3.2 Clustering Large Data

K-means is one of the most studied clustering algorithm out there and is still actively researched. It's a simple algorithm that partition the data into k partitions by minimizing

its objective function sum of squared error. From its appearance in a standard version by Stuart Lloyd in 1982 [?] , it has been one of the most popular clustering algorithm to research because of its simplicity. There are many different research areas regarding k-means. The problem with manually set the number k of partitions, initializing the centers for the partitions and dealing with outliers in data. In recent years k-means has also been very popular algorithm to study in the Map-Reduce framework where the k-means algorithm can easily be applied to cluster large amount of data sets in parallel. We start with going through work relates to both working with a stream of data of length n that doesn't fit in memory and the incrementally evolving characteristics of clustering endless data streams.

We start with a famous work by Guha et al. [?] where they introduce the STREAM algorithm that is based on the divide-and-conquer strategy and achieves a constant-factor approximation solving the k-median problem (a k-means variant). The algorithm divides the dataset into m pieces of similar sizes. Each of the pieces are independently clustered sequentially and all the centers from all the pieces are then clustered further. They show a new k-median algorithm called LSEARCH that is used by the stream algorithm and is based on local search algorithm solving the facility location problem [?] to solve the k-median problem. Results show that LSEARCH produced better quality clusters than k-means and the hierarchical algorithm BIRCH [?] but took longer to run.

In 2009 Ailon et al. [?] extended the the work of Guha et al. mentioned above and showed an one pass streaming algorithm for k-means with approximation guarantees. Achieving that they introduced a new algorithm called k-means# that builds on the non-streaming algorithm k-means++ by Arthur et al. [CITE] that is a combined algorithm of seeding the initial centers and running k-means. k-means# provides a bi-criterion (α, β) approximation algorithm by choosing $\alpha * k$ centers with approximation factor β . In the divide-and-conquer strategy they run the k-means# independently on each piece of the data to achieve $O(k \log k)$ random centers non-uniformly and use the k-means++ algorithm introduced in work of Arthur et al. [CITE] to find k centers from the intermediate centers found in all of the pieces of the data set.

Another approach using a coresets by selecting a weighted subset from the original data set such that by running any k-means algorithm on the subset will give near similar

results to running k-means on the original data set. Ackermann et al. [CITE] showed a new algorithm called StreamKM++ that uses k-means++ algorithm from Arthur et al. to solve k-means on the subset and also introduce a new data structure called coresets tree to speed up the time for the sampling in the center initialization. StreamKM++ is a streaming version of k-means++ to cluster large data sets. Their approach was shown to be on par with LSEARCH algorithm in cluster quality but outperformed BIRCH by factor of 2. A recent work by Shindler et al. [CITE] proposed an algorithm called *Fast streaming k-means* based on the online facility location algorithm [CITE Adam Meyerson] and extends the work of TODO [CITE] by proving faster running time and a better approximation factor. Their algorithm outperformed both work of Ailon et al. [CITE] and Ackermann et al. [CITE].

Adam Meyerson Online facility location. *in* FOCS, 2001

A different approach was introduced by Sculley in 2010 [CITE]. Using sampling and a gradient descent, Sculley [CITE] introduced a new algorithm called Mini-batch k-means that yields excellent clustering results with low computational cost for large datasets. The algorithm initializes the k centers like the normal k-means algorithm but then for each iteration it picks b examples from the dataset and for each point in the example the center which it is closest to is updated by taking a gradient descent step towards the point, with a learning rate of that center. The approach scales when datasets grow large with redundant examples and allowing convergence to better solutions compared to the on-line stochastic gradient descent (SGD) variant proposed by Bottou et al. [CITE].

Turning to clustering a data stream where data arrives continuously and the behavior of the stream can change over time. This area of research is very popular in the recent years where generating data is rapid, has an unknown length and not possible to access historic data points seen before because of the amount of data.

MapReduce. A parallel version of k-means++ initialization algorithm can be found in the work of Bahmani et al from 2012 [CITE]. The algorithm is called k-means|| and is implemented in the Map-Reduce framework and they show it outperforms the k-means++ algorithm in both sequential and parallel settings.

3.2.0.1 Map-Reduce

3.2.1 Data Streams

Chapter 4

Methodology

4.1 Game Events Data

4.1.1 Building Features

4.2 K-Means in Map-Reduce

4.2.1 Euclidean distance

4.2.2 Normalized Histograms

4.3 Experimental Set-Up

Chapter 5

Results and Discussion

5.1 Results

5.2 Discussion

Chapter 6

Conclusions

Appendix A

Appendix Title Here

Write your Appendix content here.

Bibliography

- [1] Nir Ailon, Ragesh Jaiswal, and Claire Monteleoni. Streaming k-means approximation. *Advances in Neural Information Processing Systems*, 22:10–18, 2009.
- [2] M. Charikar and S. Guha. Improved combinatorial algorithms for the facility location and k-median problems. In *Foundations of Computer Science, 1999. 40th Annual Symposium on*, pages 378–388, 1999.
- [3] S. Guha, A. Meyerson, N. Mishra, R. Motwani, and L. O’Callaghan. Clustering data streams: Theory and practice. *Knowledge and Data Engineering, IEEE Transactions on*, 15(3):515–528, 2003.
- [4] S. Lloyd. Least squares quantization in pcm. *Information Theory, IEEE Transactions on*, 28(2):129–137, 1982.
- [5] Tian Zhang, Raghu Ramakrishnan, and Miron Livny. Birch: an efficient data clustering method for very large databases. *SIGMOD Rec.*, 25(2):103–114, June 1996.