

Real-Time Big Data Stream Analytics

Albert Bifet

Université Paris-Saclay

Télécom ParisTech

Département Informatique et Réseaux

46 rue Barrault

75634 Paris Cedex 13, FRANCE

albert.bifet@telecom-paristech.fr

Abstract

Big Data is a new term used to identify datasets that we cannot manage with current methodologies or data mining software tools due to their large size and complexity. Big Data mining is the capability of extracting useful information from these large datasets or streams of data. New mining techniques are necessary due to the volume, variability, and velocity, of such data. MOA is a software framework with classification, regression, and frequent pattern methods, and the new APACHE SAMOA is a distributed streaming software for mining data streams.

1 Introduction

Big Data is a new term used to identify the datasets that due to their large size, we can not manage them with the typical data mining software tools. Instead of defining “Big Data” as datasets of a concrete large size, for example in the order of magnitude of petabytes, the definition is related to the fact that the dataset is too big to be managed without using new algorithms or technologies. There is need for new algorithms, and new tools to deal with all of this data. Doug Laney (Laney, 2001) was the first to mention the 3 V’s of Big Data management:

- Volume: there is more data than ever before, its size continues increasing, but not the percent of data that our tools can process
- Variety: there are many different types of data, as text, sensor data, audio, video, graph, and more
- Velocity: data is arriving continuously as streams of data, and we are interested in obtaining useful information from it in real time

Nowadays, there are two more V’s:

- Variability: there are changes in the structure of the data and how users want to interpret that data
- Value: business value that gives organizations a competitive advantage, due to the ability of making decisions based in answering questions that were previously considered beyond reach

For velocity, data stream real time analytics are needed to manage the data currently generated, at an ever increasing rate, from such applications as: sensor networks, measurements in network monitoring and traffic management, log records or click-streams in web exploring, manufacturing processes, call detail records, email, blogging, twitter posts and others. In fact, all data generated can be considered as streaming data or as a snapshot of streaming data, since it is obtained from an interval of time.

In the data stream model, data arrive at high speed, and algorithms that process them must do so under very strict constraints of space and time. Consequently, data streams pose several challenges for data mining algorithm design. First, algorithms must make use of limited resources (time and memory). Second, they must deal with data whose nature or distribution changes over time.

We need to deal with resources in an efficient and low-cost way. In data stream mining, we are interested in three main dimensions:

- accuracy
- amount of space necessary
- the time required to learn from training examples and to predict

These dimensions are typically interdependent: adjusting the time and space used by an algorithm can influence accuracy. By storing more pre-computed information, such as look up tables, an algorithm can run faster at the expense of space. An algorithm can also run faster by processing less information, either by stopping early or storing less, thus having less data to process.

2 MOA

Massive Online Analysis (MOA) (Bifet et al., 2010) is a software environment for implementing algorithms and running experiments for online learning from evolving data streams. MOA includes a collection of offline and online methods as well as tools for evaluation. In particular, it implements boosting, bagging, and Hoeffding Trees, all with and without Naïve Bayes classifiers at the leaves. Also it implements regression, and frequent pattern methods. MOA supports bi-directional interaction with WEKA, the Waikato Environment for Knowledge Analysis, and is released under the GNU GPL license.

3 APACHE SAMOA

APACHE SAMOA (SCALABLE ADVANCED MASSIVE ONLINE ANALYSIS) is a platform for mining big data streams (Morales and Bifet, 2015). As most of the rest of the big data ecosystem, it is written in Java.

APACHE SAMOA is both a framework and a library. As a framework, it allows the algorithm developer to abstract from the underlying execution engine, and therefore reuse their code on different engines. It features a pluggable architecture that allows it to run on several distributed stream processing engines such as Storm, S4, and Samza. This capability is achieved by designing a minimal API that captures the essence of modern DSPEs. This API also allows to easily write new bindings to port APACHE SAMOA to new execution engines. APACHE SAMOA takes care of hiding the differences of the underlying DSPEs in terms of API and deployment.

As a library, APACHE SAMOA contains implementations of state-of-the-art algorithms for distributed machine learning on streams. For classification, APACHE SAMOA provides a Vertical Hoeffding Tree (VHT), a distributed streaming version of a decision tree. For clustering, it includes an algorithm based on CluStream. For regression,

HAMR, a distributed implementation of Adaptive Model Rules. The library also includes meta-algorithms such as bagging and boosting.

The platform is intended to be useful for both research and real world deployments.

3.1 High Level Architecture

We identify three types of APACHE SAMOA users:

1. Platform users, who use available ML algorithms without implementing new ones.
2. ML developers, who develop new ML algorithms on top of APACHE SAMOA and want to be isolated from changes in the underlying SPEs.
3. Platform developers, who extend APACHE SAMOA to integrate more DSPEs into APACHE SAMOA.

4 Conclusions

Big Data Mining is a challenging task, that needs new tools to perform the most common machine learning algorithms such as classification, clustering, and regression.

APACHE SAMOA is a platform for mining big data streams, and it is already available and can be found online at <http://www.samoa-project.net>. The website includes a wiki, an API reference, and a developer's manual. Several examples of how the software can be used are also available.

Acknowledgments

The presented work has been done in collaboration with Gianmarco De Francisci Morales, Bernhard Pfahringer, Geoff Holmes, Richard Kirkby, and all the contributors to MOA and APACHE SAMOA.

References

- Gianmarco De Francisci Morales and Albert Bifet. 2015. SAMOA: Scalable Advanced Massive Online Analysis. *Journal of Machine Learning Research*, 16:149–153.
- Albert Bifet, Geoff Holmes, Richard Kirkby, and Bernhard Pfahringer. 2010. MOA: Massive Online Analysis. *Journal of Machine Learning Research*, 11:1601–1604, August.
- Doug Laney. 2001. 3-D Data Management: Controlling Data Volume, Velocity and Variety. *META Group Research Note*, February 6.