

# Chapter: Distributed Platforms and Cloud Services Enabling Machine Learning for Big Data. An Overview<sup>1</sup>

Daniel Pop, Gabriel Iuhasz and Dana Petcu

Institute e-Austria Timisoara,  
West University of Timisoara, Romania

**Abstract:** Applying popular machine learning algorithms to large amounts of data raised new challenges for machine learning practitioners. Traditional libraries does not support properly the processing of huge data sets, so that new approaches are needed. Using modern distributed computing paradigms, such as MapReduce, or in-memory processing novel machine learning libraries have been devised. In parallel, the advance of Cloud computing in the past ten years could not be ignored by machine learning community, thus a rise of Cloud-based platforms have been put in place as well. This chapter aims at presenting an overview of novel platforms, libraries and Cloud services that can be used by data scientists to extract knowledge from un-/semi-structured, large data sets. The overview covers several popular approaches, such as packages enabling distributed computing in popular machine learning environments, distributed platforms for machine learning and Cloud services for machine learning, known as Machine-Learning-as-a-Service approach. We also provide few recommendations for data scientists when considering machine learning approach for their problem.

**Keywords:** machine learning, data mining, cloud computing, big data, distributed computing

## Introduction

Analysing the large amount of collected data in companies, industry and sciences is becoming increasingly important for all impacted domains. The data to be analysed is no longer restricted to sensor data and classical databases, but it often includes textual documents and webpages (text mining, Web mining), spatial data, multimedia data, or graph-like data (molecules, social networks).

---

<sup>1</sup> This manuscript was submitted for review as contributed book chapter to “Data Science and Big Data Computing: Frameworks and Methodologies” book, edited by Zaigham Mahmood, to appear in Springer-Verlag, 2016

Although for more than two decades, parallel database products, such as Teradata, Oracle, Netezza have provided means to realize a parallel implementation of machine learning algorithms, expressing these algorithms in SQL code is a complex and difficult to maintain task. On the other side, large-scale installations of these products are expensive. Another reason for moving away from relational databases is the exponential growth of the unstructured data (e.g. audio and video) and semi-structured data (e.g. Web traffic data, social media content, sensor data) in recent years. The needs of data science practitioners with respect to data analysis tools vary quite largely across different domains, from medical statistics, bio-informatics, social networks analysis or physics. This diversity is equally important for the advancement of machine learning tools and platforms. Consequently, in the past decade researchers moved from the parallelization of machine learning algorithms and support in relational databases towards the design and implementation on top of novel distributed storage (e.g. NoSQL datastores, distributed file systems) and processing paradigms (e.g. MapReduce). From the business perspective, Software-as-a-Service (SaaS) model opened up new opportunities for machine learning providers, who moved the standalone tools towards Cloud-based machine learning services.

In this chapter we survey how distributed storage and processing platforms help data scientists to process large, heterogeneous sets of data. The tools, frameworks and services included in this chapter share a common characteristic: all run on top of distributed platforms. Thus, parallelization of machine learning algorithms, either using multiple cores CPU or GPU were not included here. The reader is referred to [37], a recent, comprehensive study covering that topic. We also avoided commercial solution providers, small or big players, since their offerings are either based on distributed open-source packages, or they do not disclose the implementation details.

In the first section, we briefly introduce the reader to the machine learning field, describing and classifying the types of problems and overviewing the challenges of applying traditional algorithms to large, unstructured datasets. The first category of tools considered in this survey covers tool, packages and libraries that enable data scientists to use traditional environments for data analysis, such as R system, Python, or statistics applications, in order to deal with large data sets. We survey next the distributed platforms for big data processing, either based on Apache Hadoop or Spark, as well as platforms specifically designed for distributed machine learning. We also include a section on scalable machine learning services delivered using Software-as-a-Service business model since they offer easy to use, user-friendly graphical interfaces supporting users in quickly getting and deploying models. The last section of the chapter summarizes our findings and provides readers with a collection of best practices in applying machine learning algorithms.

## Machine Learning for Data Science

The broadest and simplest definition of machine learning is that is a collection of computational methods that use experience, i.e. past information available to the system, to improve performance or to make predictions [25]. This information usually takes the form of electronic records collected and made available for analytical purposes. These records can take the form of pre-labelled training sets (usually by a human operator although this is not always the case). Another important source of data is that resulting from direct interactions with a given environment, either virtual, such as software interactions, network data etc., or relying on real-world natural scenarios, such as weather phenomena, water level etc. Data quality and quantity is extremely important in order to obtain an acceptable learned model. Machine learning relies on data-driven methods that combine fundamental concepts in the field of computer science with optimization, probability and statistics [25].

There is a wide array of applications to which machine learning can and is being applied, such as taming (text mining and document classification), spam detection, keyword extraction, emotion extraction, natural language processing (NLP), unstructured text understanding, morphological analysis, speech synthesis and recognition, optical character recognition (OCR), computational biology, face detection, image segmentation, image recognition, fraud detection, network intrusion detection, board and video games, navigation in self-driving vehicles, planning, medical diagnosis, recommendation systems or search engines. In all these applications, we can identify several types of learning problem, which are:

- *Classification* - assign each item from a data set to a specific category (e.g. given a document, to which domain (history, biology, mathematics) does it belong)
- *Regression and Time Series analysis* - predict a real value for each item, such future stock market values, rainfall runoff etc.
- *Ranking* - returns an ordered set of features based on some user defined criterion (e.g. Web search).
- *Dimensionality reduction (feature selection)* is used for transforming initial large feature spaces into a lower-dimensional representation so that it preserves the properties of the initial representation.
- *Clustering* is used in grouping items based on some predefined distance measure. It is usually used on very large data sets. In sociology it can be used to group individuals into communities [29].
- *Anomaly Detection* is an observation or series of observations which do not resemble any pattern or data item in a data set [7, 38].

In machine learning there are different types of training scenarios [25]. Arguably the most widely used type of training is called *supervised* learning. In this scenario the learner receives a set of labelled data for training and validation. The learned prediction model can be then applied to a larger data set and identify

all unseen data points. This type of learning is used for classification, regression (time series analysis). Supervised methods rely on the availability and accuracy of labelled data sets.

In *unsupervised* learning the learner receives unlabelled data that it has to group based on a distance measurement. In some scenarios labelled data is extremely hard to come by thus training a classification model is unfeasible. This type of learning is used for clustering, anomaly detections (a type of clustering) and dimensionality reduction.

In some cases where labelled data is only a small fraction of the overall training data set. This is called *semi-supervised* learning. The idea is that the distribution of unlabelled data can help the learner achieve a much better performance [12].

In *reinforcement* learning the training is done using an evaluation function. This means that training and testing are much more interlaced than in other learning scenarios. The performance of an algorithm in a problem environment is continuously evaluated through the monitoring and evaluation of its performance. Favourable outcomes are rewarded while unfavourable ones are punished. Reinforcement learning is used in genetic algorithm, neural networks etc. *On-line* learning is used when data is available in a sequential way. This means that the mapping between data sets and labels each time a new data point is received.

Due to the popularity of data analytics, machine learning techniques are being investigated by teams with complementary skills across very different business (finance, telecommunications, life sciences etc.).

*Statisticians, data scientists* are now facing data sets size explosion, thus coping with large size data sets is a must. These are users with strong mathematical background, proficient in statistics and mathematical software applications, such as R, Octave, Matlab, Mathematica, Python, SAS Studio or IBM's SPSS, but less experienced in coping with data sets of large dimensions, distributed computing or software development. Their expectations is to easily re-use their algorithms already available in their preferred language, and being able to run them against large data sets on distributed architectures (on-premise, or Cloud based). The section "Distributed and Cloud-based Execution Support in Popular Machine Learning Tools" overviews packages and tools crafted to this purpose.

Teams of *software engineers* often face client requirements asking for the transition from available (large) data warehouse to actionable knowledge. These are users with a vast experience in software development, skilled programmers in general-purpose programming languages and they 'speak' parallel and distributed computing. Deep mathematics and statistics is not necessarily their preferred playground, as they expect tools and libraries to enable them to integrate advanced ML algorithms in their systems and thus quickly get actionable results. They need fast, easy to customize (less number of parameters) and easy to integrate algorithms that run on distributed architectures and are able to fetch data from large data repositories. Tools addressing the requirements of system and software engineers are discussed in section "Distributed Machine Learning Platforms".

*Domain experts* know their data at best, but they are less experienced in ML algorithms and software tools. Ideally they need off-the-shelf software applications, easy to install and use, or Cloud-based Software-as-a-Service solutions allowing them to get insights on their data and produce reports and executable models for further usage. Domain experts mainly find helpful accessing various machine learning techniques using “Machine Learning as a Service” model (see the section with the same name).

Irrespective of the application field, this explosion of data, available in multiple formats, produced by different devices, electronic systems or human users raises specific challenges, which are:

- *Massive data sets.* Data sets are growing faster, being common now to hot numbers of 100TB or more. The Sloan Digital Sky Survey is 5 TB, the Common Crawl web corpus is 81 TB, and the 1000 Genomes Project is 200 TB, just to name a few.
- *Large models.* Massive data sets need large models to be learnt. Some deep neural networks are comprised of more than ten layers with more than a billion parameters [19, 31], collaborative filtering for video recommendation on Netflix comprises 1-10 billion parameters, and multi-task regression model for simplest whole-genome analysis may reach 1 billion parameters as well.
- *Inadequate ML tools and libraries.* Traditional ML algorithms used for decades (K-means, logistic regression, decision trees, Naïve Bayes) were not designed for handling large data sets and huge models; they were not developed for parallel/distributed environments.
- *“Operationalization” of predictive models.* “Operationalize” refers to integrate predictive models into automated decision-making systems and processes on a large scale in order to deliver predictions to end users, who will ultimately benefit from them. Integrating these models into multiple platforms (Web, standalone, mobile) across different business units requires a high degree of customization, which slows deployment, drives up costs and limits scalability.
- *Lack of clear contracts.* More recently, terms such as Analytics as a Service (AaaS) and big data as a Service (BDaaS) are becoming popular. They comprise services for data analysis similarly as IaaS offers computing resources. Unfortunately, the analytics services still lack well defined Service Level Agreements available for IaaS because it is difficult to measure quality and reliability of results and input data, to provide promises on execution times and guarantees on methods for analysing the data. Therefore, there are fundamental gaps on tools to assist service providers and clients to perform these tasks and facilitate the definition of contracts for both parties [23].
- *Inadequate staffing.* Market research shows that inadequate staffing and skills, lack of business support, and problems with analytics software are some of the barriers faced by corporations when performing analytics [28].

In the next three sections a selection of tools, libraries, packages and platforms designed to address these challenges is presented, organized based on their main target audience.

## Distributed and Cloud-based Execution Support in Popular Machine Learning Tools

As annual KD Nuggets survey shows [18], R, Python, SQL and SAS have been rated the preferred languages of choice for past 3 years. One of the early trends matching Cloud computing and data analysis was, around 2010s, the provision of virtual machine images (VMI) for these popular systems (R, Octave or Maple) integrated within public cloud service providers, such as Amazon Web Services, or Rackspace. After several Proof-of-Concept were successfully built, such as Cloudnumbers<sup>2</sup>, CloudStat<sup>3</sup>, Opani<sup>4</sup>, or Revolution R Enterprise<sup>5</sup>, the practice today is to provide VMI through the public cloud providers' marketplaces, such as Amazon Marketplace. One can find Amazon Machine Images (AMI), via the marketplace, for all the popular mathematical and statistics environments. Examples include Predictive Analytics Framework and Data Science Toolbox<sup>6</sup> that support both Python and R, BF Accelerated Scientific Compute for R with accelerated math libraries for boosted performance, or SAS University Edition for SAS Studio.

Much more effort has been invested in the development of plugins for the most popular machine learning platforms to allow data scientists to easily create and run time-consuming jobs over clusters of computers. This approach allows ML practitioners to reuse their existing code and adapt it for large-data sets processing, into the same environment they used for prototyping. It also leverages existing infrastructure (grids, clusters) for large-scale distributed computation and data storage.

Since R is the preferred option among Machine Learning practitioners, several packages were developed in order to enable big data processing within R, most of them being available under CRAN<sup>7</sup> packages page. These R extensions make possible to distribute the computational workload on different type of clusters, while accessing data from distributed file systems. First example is the RHadoop [4], a collection of five R packages that enables R users to run MapReduce jobs on Hadoop by writing R functions for mapping and reducing. Similarly, RHIFE<sup>8</sup> is

---

<sup>2</sup> <http://cloudnumbers.com>

<sup>3</sup> <http://cs.croakun.com>

<sup>4</sup> <http://opani.com>

<sup>5</sup> <http://www.revolutionanalytics.com>

<sup>6</sup> <http://datasciencetoolbox.org>

<sup>7</sup> [http://cran.r-project.org/web/packages/available\\_packages\\_by\\_name.html](http://cran.r-project.org/web/packages/available_packages_by_name.html)

<sup>8</sup> <http://www.stat.purdue.edu/~sguha/rhipe/doc/html/index.html>

another R package that brings MapReduce framework to R practitioners, providing seamless access to Hadoop cluster from within R environment. Using specific R functions, programmers are able to launch MapReduce jobs on the Hadoop cluster, results being easily retrieved from HDFS. Segue<sup>9</sup> for R project makes it easier to execute MapReduce jobs from within R environment on elastic clusters at Amazon Elastic MapReduce<sup>10</sup>, but lacks support for handling large data sets. RHive is an extension enabling distributed computing via HIVE in R, by a seamless integration between HQL (Hive Query Language) and R objects and functions. Snow (Simple Network of Workstations) [20] and its variants (snowfall, snowFT, doSnow) implement a framework that is able to express an important class of parallel computations and is easy to use within an interactive environment like R. It supports three types of clusters: socket-based, MPI, and PVM. Support for manipulating large data sets in R is available in H5 plugin, which provides an interface to the HDF5 API through S4-objects, supporting fast storage and retrieval of R-objects to/from binary files in a language independent format. The pbd\* (pbdBASE, pbbMPI, pbdNCDF4, pbdSLAP etc.) is a collection of R packages for programming with big data, enabling MPI distributed execution, NetCDF file system access, or tools for scalable linear algebra.

As far as Python is concerned, we should start by mentioning pyDoop<sup>11</sup>, a Python MapReduce and HDFS API for Hadoop [36]. Anaconda<sup>12</sup> is a free, scalable Python distribution for large-scale data analytics and scientific computing. It is a collection of Python packages (NumPy, SciPy, Pandas, IPython, Matplotlib, Numba, Blaze, Bokeh) that enables fast large data sets access, GPU computation, access to distributed implementations of ML algorithms and more. IPython.parallel<sup>13</sup> provides a sophisticated and powerful architecture for parallel and distributed computing [5] that enables IPython to support many different styles of parallelism including single program multiple data (SPMD), multiple program multiple data (MPMD), message passing using MPI, data parallel and others. In a tutorial at PYCON 2013, O. Grisel [26] presented how scikit-learn [30], a popular open-source library for machine learning in Python, can be used to perform distributed machine learning algorithms on a cheap Amazon EC2 cluster using IPython.parallel and StarCluster<sup>14</sup>. We should note as well that most of the libraries and frameworks considered in the next sections offer Python language bindings, but we choose not to include them in this section.

Other mathematical and statistics environments have seen similar interest in embracing big data processing. For example, HadoopLink<sup>15</sup> is a package that allows MapReduce programs being implemented in Mathematica and run them on

---

<sup>9</sup> <http://code.google.com/p/segue>

<sup>10</sup> <http://aws.amazon.com/elasticmapreduce>

<sup>11</sup> <https://github.com/crs4/pydoop>

<sup>12</sup> <https://store.continuum.io/cshop/anaconda>

<sup>13</sup> <http://ipython.org/ipython-doc/dev/parallel/>

<sup>14</sup> <http://star.mit.edu/cluster/>

<sup>15</sup> <https://github.com/shadanan/HadoopLink>

a Hadoop cluster. It looks more like a proof of concept (PoC), being stalled since 2013. Matlab has its Parallel Computing Toolbox which extends the capabilities of Matlab MapReduce and Datastore<sup>16</sup> in order to run big data application. Matlab Distributed Computing Server also supports running parallel MapReduce programs on Hadoop clusters<sup>17</sup>.

There are extensions to traditional machine learning libraries that enable execution on top of Hadoop, or Spark clusters. Weka [13], one of the most popular libraries for data mining, supports both Hadoop and Spark execution through Weka Hadoop integration [24]. There is also a commercial distribution, Pentaho [33], that offers a complete solution for big data analytics, supporting all phases of an analytics process - from pre-processing to advanced data exploration and visualization, which uses distributed Weka execution for analytics. Another example is the KNIME's [6] big data extension<sup>18</sup>, which enables the access to Hadoop via Hive. RapidMiner [15] has Radoop<sup>19</sup> that enables the deployment of workflows on Hadoop.

## Distributed Machine Learning Platforms

After distributed processing and storage environments (Hadoop, Dryad, MPI) reached an acceptable level of maturity, they became an increasingly appealing foundation for the design and implementation of new platforms for machine learning algorithms. These provide users out-of-the-box algorithms, which are run in parallel mode over a cluster of (commodity) computers. These solutions does not use statistics, or mathematics software packages, rather they offer self-contained, optimised implementations in general purpose programming languages (C/C++, Java) of state-of-the-art ML methods and algorithms. This section focuses on ML platforms specifically designed for distributed and scalable computing. The Table 1 summarizes most popular recent platforms, detailing for each the product license, type of ML problems supported, distributed environment supported for deployment, the size of users' community (evaluated based on number of downloads and forks) and the programming language used for implementation.

Name	License	<sup>20</sup> ML	Distributed	Comm.	Lang.
------	---------	------------------	-------------	-------	-------

<sup>16</sup> <http://www.mathworks.com/help/matlab/large-files-and-big-data.html>

<sup>17</sup> <http://www.mathworks.com/help/distcomp/big-data.html>

<sup>18</sup> <https://www.knime.org/knime-big-data-extension>

<sup>19</sup> <https://rapidminer.com/products/radoop/>

<sup>20</sup> ANO = Anomaly detection; CLS = Classification; CLU = Clustering; DL = Deep learning; DR = Dimensionality reduction; FRQ = Frequent pattern; MET =



		<b>Problem</b>	<b>Environment</b>		
Petuum	Open source (Sailing Lab)	DL, CLS, CLU, RGR, MET, TOP	Clusters or Amazon EC2, Google GCE	Medium	C++
Jubatus	LGPL v2.1	CLS, RGR, ANO, CLU, REC, Graph	Zookeeper	Medium	C++
Mllib (MLBase)	Apache 2.0	RGR, CLS, REC, CLU	Spark	Large	Scala, Java
Mahout	Apache 2.0	Collaborative Filtering, CLS, CLU, DR, TOP	Hadoop, Spark, H2O	Medium	Java
Oryx	Apache 2.0	REC, CLS, RGR, CLU	Hadoop, Spark	Low	Java
Trident-ML	Apache 2.0	CLS, RGR, CLU, DR	Storm	Low	Java
H2O	Apache 2.0	DL, RGR, CLS, CLU, DR	Hadoop	Medium	Java
GraphLab Create	Apache 2.0	CLU, CLS, RGR, DL, REC	Hadoop, Spark, MPI	High	C++
Vowpal Wabbit	Ms-PL	CLS, RGR, CLU	Hadoop	Medium	C++
Deeplearning 4J	Apache 2.0	DL	Hadoop, Spark, AWS, Akka	Medium	Java, Scala
Julia's MLBase	MIT License	CLS	Julia		Julia

---

Metrics learning; REC = Recommendation; RGR = Regression; TOP = Topic modelling

Flink-ML	Apache 2.0	CLS, RGR, CLU, REC	Flink Hadoop	Low	Scala
DryadLINQ	Ms Academic		Dryad, Hadoop YARN	None	C#, LINQ
NIMBLE	NA	CLU, FRQ, ANO,	Hadoop	None	Java
SystemML	NA	RGR, PageRank	Hadoop	None	DML

Table 1: Distributed ML frameworks

The IBM Research Lab has been one of the pioneers who invested in distributed machine learning frameworks. NIMBLE [2] and SystemML [1] are two high-level conceptual frameworks supporting the definition of ML algorithms and their execution on Hadoop clusters. NIMBLE, a sequel to IBM’s Parallel Machine Learning Toolbox [10], features a multi-layered framework enabling developers to express their ML algorithms as tasks, which are then passed to the next layer, an architecture independent layer, composed of one queue of DAGs of tasks, plus worker threads pool that unfold this queue. The bottom layer is an architecture dependent layer that translates the generic entities from upper layer into various runtimes, the only distributed environment supported within the proof of concept being Hadoop alone. The layered architecture of the system hides the low-level control and choreography details of most of the distributed and parallel programming paradigms (MR, MPI etc), it allows developers to compose parallel ML algorithms using reusable (serial and parallel) building blocks, but also it enables portability and scalability. SystemML proposes an R-like language (Declarative Machine Learning language) that includes linear algebra primitives and shows how it can be optimized and compiled down to MapReduce. Authors report an extensive performance evaluation on three (Group Nonnegative Matrix Factorization, Linear regression, PageRank) ML algorithms on varying data and Hadoop cluster sizes. These two systems are purely research endeavours, and they are not available to the community.

Most of the frameworks rely on Hadoop’s MapReduce paradigm and the underlying distributed file storage system (HDFS) because it simplifies the design and implementation of large-scale data processing systems . Only few frameworks (e.g. Jubatus, Petuum, GraphLab Create) have tried to propose novel distributed paradigms, customised to machine learning for big data, in order to optimize the complex, time-consuming ML algorithms.

Recognizing the limitations and difficulties of adapting general-purpose distributed frameworks (Hadoop, MPI, Dryad etc.) to ML problems, a team at CMU under E. P. Xing lead designed a new framework for distributed machine

learning able to handle massive data sets and cope with big models. Petuum<sup>21</sup> (from Perpetuum Mobile) [9, 41] takes advantage of data correlation, staleness, and other statistical properties to maximize the performance for ML algorithms, realized through core features such as a distributed Parameter Server and a distributed Scheduler (STRADS). It may run either on-premises clusters or on cloud computing resources like Amazon EC2 or Google GCE

Jubatus<sup>22</sup> is a distributed computing framework specifically designed for online machine learning on big data. A loose model sharing architecture allows it to efficiently train and share machine learning models by defining three fundamental operation: update, mix and analyze [14]. Comparing to Apache Mahout, Jubatus offers stream processing and online learning, which means that the model is continuously updated with each data sample that is coming in, by fast, not memory-intensive algorithms. It requires no data storage, nor sharing; only model mixing. In order to efficiently support online learning, Jubatus operates updates on local models and then each server transmits its model difference that are merged and distributed back to all servers. The mixed model improves gradually thanks to all servers' work.

GraphLab Create<sup>23</sup>, former GraphLab project [43], is a framework for machine learning that expresses asynchronous, dynamic, graph-parallel computation while ensuring data consistency and achieving a high degree of parallel performance, in both shared-memory and distributed settings. It is an end-to-end platform enabling data scientists to easily create intelligent apps at scale, from cleaning the data, developing features, training a model, and creating and maintaining a predictive service. It runs on distributed Hadoop/YARN clusters, as well on local machine or on EC2 and it exposes a Python interface for an easy accessibility.

Apache Mahout [27] is a scalable machine learning framework built on top of Hadoop that features a rich collection of distributed implementations of machine learning and data mining algorithms. Although initially created on top of Hadoop, starting with version 0.10 it supports additional execution engines, such as Spark and H2O, while Flink<sup>24</sup> is a work in progress. The same release introduces Mahout-Samsara, a new math environment created to enable users to develop their own extensions, using Scala language, based on general linear algebra and statistical operations. Mahout-Samsara comes with an interactive shell that runs distributed operations on a Spark cluster. This make prototyping or task submission much easier and allows users to customize algorithms with a whole new degree of freedom.

H2O<sup>25</sup> and MLLib [11] are two of the most actively developed projects. Both feature distributed, in-memory computations and are certified for Apache Spark (MLLib being part of Spark), as well as for Hadoop platforms. This in-memory

---

<sup>21</sup> <http://petuum.org>

<sup>22</sup> <http://jubat.us>

<sup>23</sup> <https://dato.com/products/create>

<sup>24</sup> <https://flink.apache.org/>

<sup>25</sup> <http://0xdata.com/product/>

capability means that in some instances these frameworks outperform Hadoop-based frameworks [44]. MLLib has been shown to be more scalable than Vowpal Wabbit. One important distinction when comparing H2O with other MapReduce applications is that each H2O node (which is a single JVM process) runs as a mapper in Hadoop. There are no combiners, nor reducers. Also, H2O has more built-in analytical features and a more mature REST API for R, Python and JavaScript than MLLib.

Vowpal Wabbit<sup>26</sup> [39] is an open source, fast, out-of core learning system, currently sponsored by Microsoft research. It has an efficient implementation of online machine learning, using the so called “hash trick” [34] as the core data representation, which results in significant storage compression for parameter vectors. VW reduces regression, multi-class, multi-label, or structured prediction problems to a weighted binary classification problem. A Hadoop-compatible computational model called AllReduce [3] has been implemented in order to eliminate MPI and MapReduce drawbacks which relate to machine learning. Using this model a 1000 node cluster was able to learn a tera-feature data-set in one hour [3].

Julia<sup>27</sup> is high-level, high-performance and dynamic programming language. It is designed for computing and provides a sophisticated compiler, distributed parallel execution, numerical accuracy and has an extensive mathematical function library. Comparing to traditional MPI, Julia’s implementation of message passing is “one sided”, simplifying thus the process management. Furthermore, these operations typically do not look like “message send” and “message receive” but rather resemble higher-level operations like calls to user functions. It also provides a powerful browser based notebook using IPython. It also possesses a built in package manager and it is able to call C functions directly. It is specially designed for parallelism and distributed computation. It also provides a variety of classification, clustering and regression analysis packages<sup>28</sup> implemented in Julia.

Another framework focusing on real-time online machine learning is Trident-ML [16], built on top of Apache Storm, a distributed stream processing framework. It processes batches of tuples in a distributed way which means that it can scale horizontally. However, Storm does not allow state updates to append simultaneously, a shortage that hinders distributed model learning.

The Apache Oryx 2<sup>29</sup> framework is a realization of the lambda architecture built on top of Spark and Apache Kafka. It is a specialized framework that provides real-time, large scale machine learning. It consists of three tiers: lambda, machine learning and application. The lambda tier is further split up into batch, speed and serving tier respectively. Currently it has only 3 end-to-end implementation that implement the batch, speed and serving layer (collaboration filtering, k-Means clustering, classification and regression based on random

---

<sup>26</sup> [https://github.com/JohnLangford/vowpal\\_wabbit/](https://github.com/JohnLangford/vowpal_wabbit/)

<sup>27</sup> <http://julialang.org/>

<sup>28</sup> <http://mlbasejl.readthedocs.org/en/latest/>

<sup>29</sup> <http://oryxproject.github.io/oryx/>

forest). Although it has only these three complete implementations its main design goal is not that of a traditional machine learning library but more a lambda architecture based platform for MLLib and Mahout. At this point it is important to note several key differences between Oryx 1<sup>30</sup> and Oryx 2. Firstly Oryx 1 has a monolithic tier for lambda architecture while Oryx 2 has three as mentioned in the previous paragraph. The streaming based batch layer in Oryx 2 is based in Spark while in the first version it was a custom MapReduce implementation in the Computational Layer. Two of the most important differences relate to the deployment of these frameworks. Oryx 2 is faster yet more memory hungry than the previous version because of its reliance on Spark. Second, the first version supported local (non-Hadoop) deployment while the second version does not.

DryadLINQ<sup>31</sup> [21] is LINQ<sup>32</sup> (Language INtegrated Query) subsystem developed at Microsoft Research on top of Dryad [22], a general purpose architecture for execution of data parallel applications. A DryadLINQ program is a sequential program composed of LINQ expressions performing arbitrary side-effect-free transformations on datasets, and can be written and debugged using standard .NET development tools. The system transparently translates the data-parallel portions of the program into a distributed execution plan which is passed to the Dryad execution platform that ensures efficient and reliable execution of this plan. Following Microsoft's decision to focus on bringing Apache Hadoop to Windows systems, this platform has been abandoned and Daytona project took off, which has recently become Windows Azure Machine Learning platform.

Deeplearning4J<sup>33</sup> is a open source distributed deep learning library written in Java and Scala. It is largely based on ND4J library for scientific computation that enables GPU, as well as native code integration. It is also deployable on Hadoop, Spark and Mesos. The main differences between this library and the others mentioned is that it is mainly focused on business use cases not on research. This means that some features such as parallelism is automatic, meaning that worker nodes are set up automatically.

Apache SAMOA<sup>34</sup> is a distributed streaming machine learning framework enabling developers to focus on implementing distributed algorithms and not worry about the underlying complexities of the stream processing engines (Storm, S4, Samza, etc.).

---

<sup>30</sup> <https://github.com/cloudera/oryx>

<sup>31</sup> <http://research.microsoft.com/en-us/projects/dryad/>

<sup>32</sup> <http://msdn.microsoft.com/netframework/future/linq/>

<sup>33</sup> <http://deeplearning4j.org>

<sup>34</sup> <https://samoa.incubator.apache.org/>

## Machine Learning as a Service

This section focuses on Software-as-a-Service providers for machine learning services (MLaaS). These services are accessible via RESTful interfaces, and in some cases, the solution may also be installed on-premise (e.g. ersatz). The favourite class of machine learning problems addressed by these services is predictive modelling (e.g. BigML, Google Prediction API, Eigendog), while clustering and anomaly detection receive less attention. We did not include in this category the fair number of SQL over Hadoop processing solutions (e.g. Cloudera Impala, Hadapt, Hive), because their main target is not machine learning problems, rather fast, elastic and scalable SQL processing of relational data using the distributed architecture of Hadoop.

Name	ML Problems	Data Source	Model Export	Deployment
Azure ML	CLS, RGR, CLU, ANO	Upload, Azure	None	Cloud
PredictionIO	RGR, CLS, REC, CLU	Upload, Hbase	None	Local, Cloud
Ersatzlabs	DL	Upload	None	Cloud, Local
ScienceOps <sup>35</sup> (ScienceBox)	RGR, CLS, REC, CLU	S3, Upload	PMML	Cloud, Local
Skymind	DL	Upload	None	Cloud, Local
BigML <sup>36</sup>	CLS, RGR, CLU	Upload, S3, Azure, odata	PMML	Cloud
Amazon ML	CLS, RGR, CLU	S3, Redshift Upload	None	Cloud
BitYota <sup>37</sup>	CLS, RGR, CLU	S3,	None	Cloud

<sup>35</sup> <https://yhathq.com/products/scienceops>

<sup>36</sup> <http://bigml.com>

		Azure		
Google Prediction API	CLS, RGR, CLU, ANO	Upload, Google Cloud Storage	PMML	Cloud
EigenDog <sup>38</sup>	CLU, RGR	Upload, S3	None	Local, Cloud
Metamarkets <sup>39</sup>	CLU, ANO	Upload, HDFS	None	Local (Druid), Cloud
Zementis ADAPA <sup>40</sup>	CLS, RGR, CLU	S3, Azure, Upload, SAP HANA	PMML	Local, Cloud
Predictobot <sup>41</sup>	CLS, RGR	Upload	None	Cloud

Table 2: MLaaS Offerings

Table 2 presents a selection of most popular MLaaS solutions as well as some of their key characteristics: supported machine learning problems, supported data sources, whether the built model can be exported and model deployment. As for supported data acquisition options, all platforms support data upload (.csv, .arff are widely accepted), whereas several support integration with Cloud and distributed storage solutions (S3, HDFS etc.). Predictive model training, verification and visualization is supported by all the platforms, however only few support predictive model exporting in PMML<sup>42</sup>. In the remaining of this section, we will discuss several services from Table 2.

Windows Azure Machine Learning, formerly project Daytona, was officially launched in February 2015 as a Cloud based platform for big data processing. It comes with a rich set of predefined templates for data mining workflows, as well with a visual workflow designer that allows end-users to compose complex

---

<sup>37</sup> <http://bityota.com>

<sup>38</sup> <https://eigendog.com/#home>

<sup>39</sup> <http://metamarkets.com/>

<sup>40</sup> <http://zementis.com/products/adapa/amazon-cloud/>

<sup>41</sup> <http://predictobot.com>

<sup>42</sup> <http://www.dmg.org/v4-1/GeneralStructure.html>

machine learning workflows. In addition, it supports the integration of R and Python scripts within workflows, and is able to run the jobs on Hadoop and Spark platforms. The built models are deployed in a highly scalable cloud environment and can easily be accessed via Web services [40].

PredictionIO<sup>43</sup> is based on open source software, such as Spark. Thus the solution can be deployed and hosted on any infrastructure. This contrasts to Azure Machine Learning that requires data to be uploaded into Windows Azure. Also, it is possible to write custom distributed data processing tasks in Scala while on Azure custom scripts can only be run on a single node. There are no restrictions on the size of the training data or on the number of concurrent request. It can be deployed on AWS, using Vagrant, Docker or even starting from source code.

The recent popularity of deep learning has resulted in the creation of several services wrapping well-known deep learning libraries (Tehano, pylearn2, deeplearning4j etc.) into a machine learning as a service. Some good examples are Ersazlabs<sup>44</sup> and Skymind<sup>45</sup>. These provide similar services and support distributed, as well as GPU deployment.

Amazon Machine Learning service<sup>46</sup> allows users to train predictive models in the cloud. It targets a similar use case as Azure Machine learning from Microsoft and Google's Predictive API. It has similar features to many large scale learning applications including visualization and basic data statistics. The exact learning algorithm it uses is not known however, it is similar to Vowpal Wabbit. There are some limitation such as the inability to export the learned model or to access data which is not stored inside amazon (Amazon S3 or Redshift).

Google Prediction API<sup>47</sup> is Google's cloud-based machine learning tools that can help analyse your data. It is closely connected to Google Cloud Storage<sup>48</sup> where training data is stored and offers its services using a REST interface, client libraries allowing programmers to connect from Java, JavaScript, .NET, Ruby, Python etc. In the first step, the model need to be trained from data, supported models being classification and regression for now. After the model is built, one can query this model to obtain predictions on new instances. Adding new data to a trained model is called Streaming Training and it is also nicely supported. Recently, PMML pre-processing feature has been added, i.e. Prediction API .supports pre-processing your data against a PMML transform specified using

---

<sup>43</sup> <https://prediction.io/>

<sup>44</sup> <http://www.ersatzlabs.com/>

<sup>45</sup> <http://www.skymind.io/about/>

<sup>46</sup> <http://docs.aws.amazon.com/machine-learning/latest/mlconcepts/mlconcepts.html>

<sup>47</sup> <https://developers.google.com/prediction/>

<sup>48</sup> <https://developers.google.com/storage/>



PMML 4.0 syntax; does not support importing of a complete PMML model that includes data. Created models can be shared as hosted models in the marketplace.

## Related studies

Since 1995, when K. Thearling [17] presented a massively parallel architecture and the algorithms for analyzing time series data, allegedly, one of the first approaches to parallelization of ML algorithms, many implementations were proposed for ML algorithms parallelization for both shared and distributed systems. Consequently, many studies tried to summarize, classify and compare these approaches. We will address in this section only the most recent ones.

Upadhyaya [37] presents an overview of machine learning efforts since 1995 onwards grouping the approaches based on prominent underlying technologies: those employed on GPUs (2000–2005 and beyond), those using MapReduce technique (2005 onward), the ones that did not consider neither MapReduce, nor GPUs (1999–2000 and beyond) and, finally, few efforts discussing the MapReduce technique on GPU.

The book “Scaling up machine learning: parallel and distributed approaches” by Bekkerman et al. [32] presents an integrated collection of representative approaches, emerged in both academic (Berkeley, NYU, University of California etc.) and industrial (Google, HP, IBM, Microsoft) environments, for scaling up machine learning and data mining methods on parallel and distributed computing platforms. It covers general frameworks for highly scalable ML implementations, such as DryadLINQ and IBM PMLT, as well as specific implementations of ML techniques on these platforms, like ensemble decision trees, SVM, or k-Means.

A broader study is conducted by Assunção et al. [23], who discusses approaches and environments for carrying out analytics on Clouds for big data applications. Model development and scoring, i.e. machine learning, is one of the areas they considered, alongside other three: data management and supporting architectures, visualisation and user interaction, and business models. Through a detailed survey, they identify possible gaps in technology and provide recommendations for the research community on future directions on Cloud-supported big data computing and analytics solutions.

Interesting analyses have been made available in online press and blogs [35, 42, 8], who have reviewed open-source, or commercial players for big data analytics and prediction.

## Summary and Guidelines

Analysing big data sets gives users the power to identify new revenue sources, develop loyal and profitable customer relationships, and run the organization more

efficiently and cost effectively, all in all giving them competitive advantage over competition. Big data analytics is still a challenging and time demanding task that requires important resources, in terms of large e-infrastructure, complex software, skilled people, significant effort, without any guarantee on ROI. After reviewing more than 40 solutions, our key findings are summarized below.

Existing programming paradigms for expressing large-scale parallelism (MapReduce, MPI) are the *de facto* choices for implementing distributed machine learning algorithms. The initial enthusiastic interest devoted to MapReduce has been balanced in recent years by novel distributed architectures specifically designed for machine learning problems. Nevertheless, Hadoop remains the state-of-the-art platform for processing large data sets stored on HDFS, either in MapReduce jobs or using higher-level languages and abstractions.

Both, research and industry, invested significant efforts in developing “as a Service” solutions for big data problems (Analytics as a Service, Data as a Service and Machine Learning as a Service), in order to benefit of the advantages Cloud computing provides: resources on-demand (with costs proportional to the actual usage), scalability and reliability.

Although state-of-the-art tools and platforms provide intuitive graphical user interfaces, current environments lack an interactive process, and techniques should be developed to facilitate interactivity in order to include analysts in the loop, by providing means to reduce time to insight. Systems and techniques that iteratively refine answers to queries and give users more control of processing are desired.

In the end of this chapter, several best practices, readily available in the literature, are being outlined:

(1) Understand the business problem. Having a well-defined problem, knowing specific constraints available for the problem under investigation, can greatly improve in performance of your ML algorithms.

(2) Understand the ML task. Is it supervised or unsupervised? What activities are required to get the data labelled? The same features (attributes, domains, labels) need to be available at both times, training and testing. Pick a machine learning method appropriate to the problem and the data set. This is the most difficult task and here are some questions you should consider: Do human users need to understand the model? Is the training time a constraint for your problem? What is an acceptable trade-off between having an accurate answer versus having the answer quickly? Keep in mind that there is no single best algorithm; experiment with several algorithms and see which one gives better results for your problem.

(3) In case of predictive modelling, carefully select and partition the data at hand in training and validation set, which will be used to build your model, versus test set that you will use to test the performance of your model. More data for training the model, better predictive performance. Better data always beats a better algorithm, no matter how advanced it is. Visualize the data with at least univariate histograms. Examine correlations between variables.

(4) Well prepare your data. Deal with missing and invalid values (misspelled words, values out of range, outliers). Take enough time, because no matter how robust a model is, poor data will yield poor results.

(5) Evaluate your model using confusion matrix, ROC (Receiver Operating Characteristic) curve, precision, recall or F1 score. Do not overfit your model, because the power lies in good prediction of unseen examples.

(6) Use proper tools for your problems. Low-level programming environments you might found difficult to use; try first machine learning services offered by Cloud service providers, which are easy to use and powered by state-of-the-art algorithms.

## Acknowledgements

This work was partially supported by the European Commission H2020 co-funded project DICE (GA 644869).

## References

1. A. Ghoting et al. -- SystemML: Declarative machine learning on mapreduce. In Proceedings of the 2011 IEEE 27th International Conference on Data Engineering, ICDE 11, pages 231-242, Washington, DC, USA, 2011
2. A. Ghoting, P. Kambadur, E. Pednault, and R. Kannan – NIMBLE: A Toolkit for the Implementation of Parallel Data Mining and Machine Learning Algorithms on MapReduce, KDD 11
3. Agarwal, A.; Chapelle, O.; Dudik, M. & Langford, J., A Reliable Effective Terascale Linear Learning System, *CoRR*, 2011, *abs/1110.4198*
4. Antonio Piccolboni (2015) RHadoop - <https://github.com/RevolutionAnalytics/RHadoop/wiki>, Accessed: 13.May.2015
5. B. Granger, F. Perez, M. Ragan-Kelley, Using IPython for parallel computing, <http://minrk.github.com/scipy-tutorial-2011>
6. Berthold, M. R., Cebron, N., Dill, F., Gabriel, T. R., Kötter, T., Meinl, T., Ohl, P., Sieb, C., Thiel, K. and Wiswedel, B. *KNIME: The Konstanz Information Miner* Studies in Classification, Data Analysis, and Knowledge Organization (GfKL 2007) Springer, 2007
7. Chandola, V.; Banerjee, A. & Kumar, V. *Anomaly Detection: A Survey* *ACM Comput. Surv.*, *ACM*, 2009, *41*, 15:1-15:58
8. D. Harris, 5 low-profile startups that could change the face of big data, Last accessed on May 2015 from <http://gigaom.com/cloud/5-low-profile-startups-that-could-change-the-face-of-big-data/>
9. Dai, Wei, Jinliang Wei, Xun Zheng, Jin Kyu Kim, Seunghak Lee, Junming Yin, Qirong Ho, and Eric P. Xing. "Petuum: A framework for iterative-convergent distributed ML." *arXiv preprint arXiv:1312.7651* (2013)

10. E. Pednault, E. Yom-Tov, A. Ghoting – IBM Parallel Machine Learning Toolbox, in R. Bekkerman, M. Bilenko and J. Langford (editors) – Scaling up Machine Learning, Cambridge University Press, 2012
11. Franklin, M.; Gonzalez, J.; Jordan, M. I.; Pan, X.; Smith, V.; Sparks, E.; Talwalkar, A.; Venkataraman, S. & Zaharia, M. *MLlib*, 2013
12. Gander, M.; Felderer, M.; Katt, B.; Tolbaru, A.; Breu, R. & Moschitti, A. Moschitti, A. & Plank, B. (Eds.), *Anomaly Detection in the Cloud: Detecting Security Incidents via Machine Learning*, Trustworthy Eternal Systems via Evolving Software, Data and Knowledge, Springer Berlin Heidelberg, 2013, 379, 103-116
13. Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P. and Witten, I. H. *The WEKA Data Mining Software: An Update* SIGKDD Explor. Newsl., ACM, 2009, Vol. 11(1), pp. 10-18
14. Hido, S.; Tokui, S. & Oda, S., Jubauts: An Open Source Platform for Distributed Online Machine Learning, NIPS 2013 Workshop on Big Learning, Lake Tahoe
15. Hofmann, M. & Klinkenberg, R., *RapidMiner: Data Mining Use Cases and Business Analytics Applications*, Chapman & Hall/CRC, 2013
16. Jain, A. & Nalya, A., *Learning Storm*, Packt Publishing, 2014
17. K.K. Thearling, Massively parallel architectures and algorithms for time series analysis, in: L. Nadel, D. Stien (Eds.), *Lectures in Complex Systems*, Addison-Wesley, 1995
18. KD Nuggets (2014) - <http://www.kdnuggets.com/polls/2014/languages-analytics-data-mining-data-science.html>, Accessed on: 15.May.2015
19. Krizhevsky, A., Sutskever, I. and Hinton, G. E., *ImageNet Classification with Deep Convolutional Neural Networks*. NIPS 2012: Neural Information Processing Systems, Lake Tahoe, Nevada.
20. L. Tierney, A. J. Rossini, Na Li, Snow: A parallel computing framework for the R System, *Int J Parallel Prog* (2009) 37:78-90, DOI 10.1007/s10766-008-0077-2
21. M. Budi, D. Fetterly, M. Isard, F. McSherry, and Y. Yu – Large-Scale Machine Learning using DryadLINQ, in R. Bekkerman, M. Bilenko and J. Langford (editors) – Scaling up Machine Learning, Cambridge University Press, 2012
22. M. Isard et al. – Dryad: distributed data-parallel programs from sequential building blocks. In *SIGOPS Operating System Review*, 2007
23. M.D. Assunção, R. N. Calheiros, S. Bianchi, M. A. S. Netto, R. Buyya, *big data computing and clouds: Trends and future directions*, *J. Parallel Distrib. Comput*, 2014, <http://dx.doi.org/10.1016/j.jpdc.2014.08.003>
24. Mark Hall (2013), Weka and Spark - <http://markahall.blogspot.co.nz/>, Accessed on: 13.May.2015
25. Mohri, M.; Rostamizadeh, A. & Talwalkar, A. *Foundations of Machine Learning*, The MIT Press, 2012
26. O. Grisel, Advanced Machine Learning with scikit-learn, PYCON 2013, Tutorial, <https://us.pycon.org/2013/schedule/presentation/23/>

27. Owen, S.; Anil, R.; Dunning, T. & Friedman, E., *Mahout in Action*, Manning Publications Co., 2011
28. P. Russom, big data Analytics, *TDWI best practices report*, The Data Warehousing Institute (TDWI) Research (2011)
29. Patcha, A. & Park, J.-M. *An Overview of Anomaly Detection Techniques: Existing Solutions and Latest Technological Trends*, *Comput. Netw.*, Elsevier North-Holland, Inc., 2007, 51, 3448-3470
30. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, E. *Scikit-learn: Machine Learning in Python*, *Journal of Machine Learning Research*, 2011, Vol. 12, pp. 2825-2830
31. Q. Le, M. A. Ranzato, R. Monga, M. Devin, K. Chen, G. Corrado, J. Dean, A. Ng, *Building high-level features using large scale unsupervised learning*, International Conference in Machine Learning, 2012
32. R. Bekkerman, M. Bilenko, J. Langford (Editors), *Scaling up machine learning: parallel and distributed approaches*, Cambridge University Press, 2012
33. Roldn, M. C., *Pentaho Data Integration Beginner's Guide*, Packt Publishing, 2013
34. Rosen, J.; Polyzotis, N.; Borkar, V. R.; Bu, Y.; Carey, M. J.; Weimer, M.; Condie, T. & Ramakrishnan, R., *Iterative MapReduce for Large Scale Machine Learning*, *CoRR*, 2013, [abs/1303.3517](https://arxiv.org/abs/1303.3517)
35. S. Charrington, Three new tools bring machine learning insights to the masses, February 2012, Read Write Web, <http://www.readwriteweb.com/hack/2012/02/three-new-tools-bring-machine.php>
36. S. Leo and G. Zanetti. Pydoop: a Python MapReduce and HDFS API for Hadoop. In Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing, 819-825, 2010.
37. S. R. Upadhyaya, Parallel approaches to machine learning — A comprehensive survey, *Journal of Parallel and Distributed Computing*, Volume 73, Issue 3, March 2013, Pages 284-292, ISSN 0743-7315, <http://dx.doi.org/10.1016/j.jpdc.2012.11.001>
38. Sagha, H.; Bayati, H.; Millán, J. D. R. & Chavarriaga, R. *On-line Anomaly Detection and Resilience in Classifier Ensembles* *Pattern Recogn. Lett.*, Elsevier Science Inc., 2013, 34, 1916-1927
39. Shi, Q.; Petterson, J.; Dror, G.; Langford, J.; Smola, A. & Vishwanathan, S., *Hash Kernels for Structured Data*, *J. Mach. Learn. Res.*, *JMLR.org*, 2009, 10, 2615-2637
40. Stephen F. Elston, *Data Science in the Cloud with Microsoft Azure Machine Learning and R*, O'Reilly, 2015

41. W. Dai, A. Kumar, J. Wei. Q. Ho, G. Gibson and E. P. Xing, High-Performance Distributed ML at Scale through Parameter Server Consistency Models, AAI 2015
42. W. Eckerson, New technologies for big data, 2012, [http://www.b-eye-network.com/blogs/eckerson/archives/2012/11/new\\_technologie.php](http://www.b-eye-network.com/blogs/eckerson/archives/2012/11/new_technologie.php)
43. Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, J. M. Hellerstein – Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud, Proceedings of the VLDB Endowment, Vol. 5, No. 8, August 2012, Istanbul, Turkey
44. Zaharia, M.; Chowdhury, M.; Franklin, M. J.; Shenker, S. & Stoica, I., *Spark: Cluster Computing with Working Sets*, Proceedings of the 2Nd USENIX Conference on Hot Topics in Cloud Computing, USENIX Association, 2010, 10-10