

Big Data: Data Science Applications and Present Scenario

Shubhankar Chaturvedi and Shwetank Kanava
Arya College of Engineering and Information Technology, Jaipur

Abstract

In this paper we are presenting some simple study of data science which has been discussed very frequently in scientific community. We are also giving some recent trends and techniques and their impact on scientific as well as social community.

I. INTRODUCTION

Data science is not a new concept or term for statistical scientists, but this is a simply computerization of statistical old methods as per the need of present time. Big Data are large data sets. However, there is a lot of informative value hidden in this data, so many companies have desire to access this data for maximize their profit. It is very innovative and competitive research area in present time.

II. DATA SCIENCE & CHALLENGES

During the last few years, the most challenging problem the world developed was data science. The data science problem means that data is growing at a much faster rate than computational solution techniques, and it is the result of the fact that storage cost is getting cheaper gradually day by day, therefore, keeping data safe and secure for further use becomes cheaper with time. Social activities, biological explorations, scientific experiments, along with the sensor devices are great data contributors. Data science is beneficial to the society and business but at the same time, it brings challenges to the scientific communities. The existing traditional tools, machine learning algorithms and techniques are not capable of handling, managing and analyzing big data. Although various scalable machine learning algorithms, techniques and tools are prevalent. In this paper we have identified the most relevant issues and challenges related to data science and point out a comprehensive comparison of various technical Theories and techniques from many fields and disciplines are used to investigate and analyze a large amount of data to help decision makers in many industries such as science, engineering, economics, politics, finance, and education; Computer Science- Pattern recognition, visualization, data warehousing, High performance computing, Databases, AI; Mathematics - Mathematical Modeling in various

scientific methods like Statistics, Statistical and Stochastic modeling, Probability and queuing for handling data science problem.

III. DATA SCIENCE AND ACADEMICS

In the words of Alex Szalay, these sorts of researchers must be "Pi-shaped" as opposed to the more traditional "T-shaped" researcher. In Szalay's view, a classic PhD program generates T-shaped researchers: scientists with widebut-shallow general knowledge, but deep skill and expertise in one particular area. The new breed of scientific researchers, the data scientists, must be Pishaped: that is, they maintain the same wide breadth, but push deeper both in their own subject area and in the statistical or computational methods that help drive modern research.

IV. DATA MINING/SCIENCE WITH BIG DATA

Aspects of big data have been studied and considered by a number of data mining researchers over the past decade and beyond. Mining massive data by scalable algorithms leveraging parallel and distributed architectures has been a focus topic of numerous workshops and conferences. However, the embrace of the Volume aspect of data is coming to a realization now, largely through the rapid availability of datasets that exceed terabytes and now penta bytes—whether through scientific simulations and experiments, business transactional data or digital footprints of individuals. Astronomy, for example, is a fantastic application of big data driven by the advances in the astronomical instruments. Each pixel captured by the new instruments can have a few thousand attributes and translate quickly to a penta scale problem. This rapid growth in data is creating a new field called Astro-informatics, which is forging partnerships between computer scientists, statisticians and astronomers. The emergence of big data from various domains, whether in business or science or humanities or engineering, is presenting novel challenges in scale and provenance of data, requiring a new rigor and interest among the data mining community to translate their algorithms and frameworks for data-driven discoveries.

V. ESSENTIAL POINTS

Big Data has given rise to Data Science; Data science is rooted in solid foundations of mathematics and statistics, computer science, and domain knowledge, not everything with data or science is Data Science, The use cases for Data Science are compelling.

VI. DATA VARIETY

Data presents itself in varied forms for a given concept. It is presenting a new notion to learning systems and computational intelligence algorithm for classification, where the feature vector is multi-modal, with structured and unstructured text, and still the notion is to classify one concept from another. How do we create a feature vector, and then a learning algorithm with an appropriate objective function to learn from such varied data?

A. Data size

On one hand, we develop “one-pass learning” algorithms that require only one scan of the data with limited storage irrelevant to data size; on the other hand, we try to identify smaller partitions of the really valuable data from the original big data.

B. Data trust

While data is rapidly and increasingly available, it is also important to consider the data source and if the data can be trusted. More data is not necessarily correct data, and more data is not necessarily valuable data. A keen filter for the data is a key.

VII. DISTRIBUTED EXISTENCE

Owners of different parts of the data might warrant different access rights. We must aim to leverage data sources without access to the whole data, and exploit them without transporting the data. Extreme distribution: Taking this idea even further, the unit-level data may be what we see as the level of data distribution, as we deal with issues of privacy and security. New approaches to modeling big data will be required to work with extreme distributed data.

VIII. DIVERSE DEMANDS

People may have diverse demands whereas the high cost of big data processing may disable construction of a separate model for each demand. Can we build one model to satisfy the various demands? We also need to note that, with big data, it is possible to find supporting evidence to any argument we want; then, how to judge/evaluate our “findings”?

IX. SUB-MODELS

Diverse demands might also relate to diversity of the behaviors that we are modeling within our application domains. Rather than one single model to cover it all, the model will consist of ensembles of a large number of smaller models that together deliver better understandings and predictions than the single, complex model.

X. INTUITION IMPORTANCE

Data is going to power novel discoveries and action oriented business insights. It is important to still attach intuition, curiosity and domain knowledge without which one may become myopic and fall in the chasm of “correlation is enough”. Computational intelligence should be tied with human intuition.

XI. RAPID MODEL

As the world continues to “speed up”, decisions need to be made more quickly because fraudsters can more quickly find new methods in an agile environment, model building must become more agile and real-time. Algorithms such as meta-heuristics have achieved great success in academic research, but have rarely been employed in industry. One major obstacle is the huge computational cost required for evaluating the quality of candidate designs of complex engineering systems. The emerging big data analytic technologies will remove the obstacle to a certain degree by reusing knowledge extracted from the huge amount of high-dimensional, heterogeneous and noisy data. Such knowledge can also be acquired with new visualization techniques. Big data driven optimization will also play a key role in reconstruction of large-scale biological systems.

XII. COMPLEX OPTIMIZATION

Definition of decision variables, setup of the objectives and articulation of the constraints are three main steps in formulating optimization problems before solving them. For optimization of complex systems, formulation of the optimization problem itself becomes a complex optimization problem. The big data approach might provide us new insights and methodologies for formulating optimization problems, thus leading to a more efficient solution. In closing the discussion, we emphasize that the opportunities and challenges brought by big data are very broad and diverse, and it is clear that no single technique can meet all demands. In this sense, big data also brings a chance of “big combination” of techniques and of research.

