# A Survey on Big Data Analytics: Methods, Opportunities and Challenges

**Ms. Kiran R. Borade[1], Mrs. Dhanashree S. Joshi[2], Mrs. Tejal S. Sonawane[3]**

Lecturer, Department of Computer Engineering[1,2,3]

Guru Gobind Singh Polytechnic, Nashik, Maharashtra, India

kiranborade@ggsf.edu.in[1], dhanashree.joshi@ggsf.edu.in[2], tejal.sonawane@ggsf.edu.in[3]

**Abstract:** *In this paper, we briefly introduce the basic concepts and features of big data. We are surrounded by a huge amount of data but in demand of information. In the Big Data era, how to quickly get high quality and valuable information from large amounts of data has become an important research tool. Big Data Analytics has become very popular, not only academically, but also in industry and government programs, which can be attributed to the fact that Big Data Analytics offers great promises and poses significant global challenges. The extracted data can be useful to the organization in various aspects. Many decisions have to be made by business organizations based on this big data. Data is growing at a tremendous rate these days and it is an important issue for data management and management to analyze the information needed to save time and cost. The main purpose of this paper is to present a detailed analysis of the various forums/techniques suitable for Big Data processing. In this paper, the various software technologies, methods and tools available for Big Data analysis are explored and analyzed in detail with their strengths and weaknesses. In this paper various data mining and feature extraction methodologies are studied and finally comparison between technologies or tools is done based on performance of various systems.*

**Keywords:** Big data technology, Data Mining, Big Data, Big Data Analytics, Feature extraction

## I. INTRODUCTION

Big data is one of the "hottest" phrases used today. Everyone is talking about big data, and it is believed that science, business, industry, government, society, etc. it will completely change with the influence of big data. Speaking of technology, the process of managing big data involves the collection, storage, transport and exploitation. There is no doubt that the stages of collection, storage and transportation are a necessary precursor to the ultimate goal of data exploitation, which is the core of big data processing. Turning to data analysis, we note that the ─big data ‖ is defined by four Vs ─Volume, Velocity, True, and variability. It is assumed that all or one of them needs to be met in order to configure the problem as a Big Data problem. The volume indicates the size of the data, which is probably too large to be controlled by the current state of algorithms and / or systems. Speed means data is distributed at a faster rate than conventional algorithms can be handled by systems. Sensors read quickly and communicate with data streams. We are approaching the world of quantified self, which introduces data that is not yet available. The fact is that unless data is available, data quality remains a major problem. That is, we cannot assume that with big data it comes with high quality. In fact, in size comes a quality problem, which needs to be addressed in the pre-processing phase of the data or in a learning algorithm. Variety is the most compelling of all Vs as it presents data of different types and methods of a particular object being considered. Each V is not new. Machine learning and data mining researchers have been dealing with these problems for decades. However, the emergence of online-based companies has challenged many traditional process-focused companies — now they require more knowledge based companies driven data rather than process. The purpose of this article is to share the authors' views on big data in their data analysis ideas. The four authors present a wide range of ideas with a wide range of research experience, including computer literacy, machine learning, data mining and science, and multidisciplinary research. The authors represent academics and industry in all four different continents. This diversity brings together an interesting and inclusive perspective in data analysis in the context of today's big data.

It is worth noting that this article is not intended to provide a comprehensive overview of the modern "Big Data Opportunities and Challenges: Discussions from Data Analytics Perspectives November 2014 for big data research, or to provide future big data. Research agenda. Aim to express authors' personal opinions and present their views on the future based on their own views. Therefore there will be limited arguments for evidence or supporting literature, given the rapidly changing world conditions and major shortcomings in academic research reporting. Although all the authors contributed to the paper as a whole, each writer focused on their specific issues in the following discussions. Zhou covers machine reading, while Chawla brings up the concept of excavation data and data science perspective. Jin provides insights from computer intelligence and global meta-heuristic development, and Williams draws on machine learning and data mining technology used as an active data scientist and consultant in the international industry.

## II. MACHINE LEARNING WITH BIG DATA

Machine learning is among the core techniques for data analytics. In this section we will first clarify three common but unfortunately misleading arguments about learning systems in the big data era. Then we will discuss some issues that demand attention.

### 2.1. Three Misconceptions

### 2.1.1. "Models Are Not Important Any More"

Many people today talk about changing complex models with big data, when we have a large amount of data available. The argument is that in the "small data era" models were important but in the "big data era" this may not be the case. Such arguments are said to be based on the strong observation of the type shown in Fig. 1. For smaller data (e.g., 10 data size), the best model is approximately x% than the worst model in the image, and performance enhancements for large data (e.g., 104 data size) are y% & x% . Such awareness can be traced back to many years, as [7], and precede the use of ―great knowledge. ‖ It is interesting to note that in the "big data era", many people take such a number (or similar statistics) to say that having big data is enough to get better performance. Such superficial recognition, however, ignores the fact that even big data (e.g., 104 data sizes in a photo), still has significant differences between different models — still important models. Also, we often hear such controversies: As the statistics show, a simple model with small data achieves excellent performance with big data, and thus, one does not need to have a complex model because a simple model is enough. Unfortunately, this argument is also incorrect. First, there is no reason to conclude that the model that performs poorly on small data is really "simple", and vice versa. Second, even if we thought that the worst-performing model for small data was the simplest, there is no general support for the argument that a simpler model would definitely achieve better data performance in tasks without the current empirical study. Looking at [7] we can find that Algo_1 in Fig. 1 with memorybased approach, Algo_2 with Perceptron or Naïve Bayes, and Algo_3 with ink. It's hard to conclude that Winnow is simpler than a memory-based approach; at the very least, the "simplistic argument" will not be able to explain why Perception performance is better than that of a memory-based approach to big data. A more logical explanation of why a memory-based approach is better for smaller data than for data but even worse for large data may be due to its requirement to load data to memory. This is a memory problem not whether the model is complex or not. Recent interest in in-depth learning [10], [6] provides strong evidence that in large data, complex models are able to achieve much better performance than simple models. We want to emphasize that deep learning methods are not really new, and many ideas can be found from the 1990's [5], [3]. However, there were two major problems that hindered development at that time. First, the calculation tools available at the time could not handle models with thousands of parameters to be tuned. Current in-depth learning models include millions or millions of parameters. Second, the data scale at that time was small, so models with high complexity are likely to be overweight. We see that with the rapid increase in computing power, the training of high-end models becomes even more feasible, while large data sizes significantly reduce the risk of overuse of complex models. From this perspective, one can conclude that in the big data age, complex models are more popular as simple models are often unable to fully exploit data.

### 2.1.2. "Correlation Is Enough"

Some popular big data books, including [23], say that it is enough to find ―connections ‖ in big data. The importance of "causality" will be eliminated from the "connection", with some urging that we enter a "relationship". Trying to find the cause represents the main purpose of searching for a deeper understanding of data. This is often a challenge for many real domains [24]. However, we must emphasize that the correlation is far from sufficient, and the role of causality can be 10 102 103 104 70 75 80 85 90 95 100 Predictability Accuracy Data Size (%) Algo_1 Algo_2 Algo_3 will never be replaced by a relationship. The reason is that a person invests in data analysis in order to gain useful information in making wise decisions and / or to take appropriate action, while communication abuse will be misleading or catastrophic. We can easily find many examples, even in ancient mathematical literature, that show that relationships cannot replace the role of causality. For example, a robust data analysis on public safety conditions in many cities has revealed that the number of hospitals and the number of car thefts are closely related. Indeed, car theft is increasing at almost the same rate as the construction of new hospitals. Having noted such a link, how can the mayor react to reducing car theft? The ―obvious ‖ solution is to halt the construction of new hospitals. Unfortunately, this is a violation of communication information. It will only reduce the chances of patients getting health care on time, and they are less likely to have anything to do with car theft. Instead, the increase in both car theft cases and the number of hospitals is actually affected by subtle fluctuations, i.e., the people living there. If a person believes that the relationship is adequate and does not delve too deeply into the details, he may act as the mayor who plans to reduce car theft by limiting hospital construction. Sometimes computer challenges may obscure the cause, and in such cases, finding the right link, will be able to provide useful information. However, overemphasizing the importance of ―correlation ‖ and taking the causality shift by merging it as a feature of the ―big data era ‖ can be dangerous, and lead to unnecessary, negative consequences.

### 2.1.3. "Previous Methodologies Do Not Work Any More"

Another popular argument is that previous research methods were designed for small data and could not work well on large data. This debate is often held by people who are passionate about newly proposed strategies, thus seeking "completely new" paradigms. We appreciate the search for new paradigms as this is one of the ways to conduct new research. However, we emphasize the importance of ―old-fashioned ‖ methods. First, we must emphasize that researchers have been trying to work with "big" data, such as what is considered big data today may not be considered big data in the future (e.g., ten years). For example, in a popular article [3], the author stated that "By learning activities with 10,000 training examples and more becomes impossible ...". The title of the paper ―Making SVM Great Training Work ‖ suggests that the purpose of the article was ―great ‖ ―a test data set on paper typically contains thousands of samples, and the largest contains 49,749 samples. This was considered ―amazing data ‖ at the time. Today, few people will consider 50,000 samples as big data. Second, many methods of previous research still hold large numbers. We may consider [6], the continuation of the KDD 1999. workshop on page 2 emphasizes that ―use ... on a distributed computer ... becomes increasingly important in ensuring system balance and collaboration as data continues to grow inexorably in size and complexity ‖ . Indeed, the "current" assistant for managing big data, such as the most efficient computer, the most compatible and distributed computer, the most efficient storage, etc., has been used in data analysis for many years and will remain popular in the future.

### 2.2. Opportunities and Challenges

It is difficult to identify "completely new" stories that bring big data. However, there are always some important factors one hopes to see the most attention and effort put into them. First, although we have been trying to manage (increasingly) large data, we have often assumed that basic calculations can be done in memory without sewing. Although the current data size reaches such a scale that the data becomes difficult to maintain and difficult even for most scans. However, many important learning objectives or practical steps are out of line, smooth, flexible and do not rot over samples. For example, AUC (ROC Zone) [4], and your preparation, naturally requires repeated scanning of the entire database. Is it readable by scanning the data only once, and if it needs to store something, the storage requirement is small and

independent of the data size? We call this ―single-pass reading ‖ and it is important because in many big data systems, the data is not only larger but also collected over time, which is why it is impossible to know the final size of the database. Fortunately, there are recent efforts toward this direction, including [2]. On the other hand, even though we have big data, is all data important? The answer is very likely that they are not. Then, the question is whether we can identify important data sets in the main database? Second, the advantage of large data in machine learning lies in the fact that with more samples available for reading, the risk of overheating becomes less. We all understand that controlling overload is one of the most important concerns in the development of machine learning algorithms and in the application of machine learning techniques in practice. Concerns about overuse have led to a natural tendency for simpler models to have a smaller tuning parameter. However, parameter tuning parameters may change with big data. Now we can try to train a model with billions of parameters, because we have big enough data, driven by powerful computational areas that allow the training of such models. The great success of in-depth learning [10] over the past few years serves as a good showcase. However, much of the work of indepth learning relies heavily on engineering techniques that are difficult to replicate and read by others, other than the authors themselves. It is important to learn the mysteries after a thorough study; for example, why and when are some of the ingredients of current in-depth reading strategies, e.g., advance training and quitting, helpful and helpful? There have been recent efforts in this regard. Additionally, we might ask if it is possible to improve the parameter tuning index instead of the nearly complete current search? Third, we need to realize that big data often contains too many "interests", and from such data we can get "whatever we want"; in other words, we can find supporting evidence of any controversy on which we agree. So, how can we judge / evaluate ―findings ‖ ? One important solution is to turn to a mathematical hypothesis test. The use of mathematical tests can help in at least two areas: First, we need to make sure that what we did is exactly what we wanted to do. Second, we need to ensure that our gains are not caused by minor data breaches, especially due to misuse of all data. Although mathematical experiments have been studied for centuries and have been used in machine learning for decades, the design and delivery of adequate mathematical tests is no small feat, and has actually been a misuse of mathematical tests [17]. In addition, a statistical test that deserves large-scale data analysis, not only for statistical efficiency but also for the use of only part of the data, remains an interesting but less explored area of study. Another way to assess the validity of analytical results is to find translated models. Although most learning models in black box machines, there have been studies to improve the understanding of models such as law enforcement [6]. Visualization is another important technique, although it is often difficult with more than three dimensions. In addition, big data tends to exist in a distributed way; that is, different parts of the data may be stored by different owners, and no one owns all the data. It is often the case that certain sources are very important in a particular analysis objective, while other sources are less important. Given the fact that different data owners may authorize an analyst with different access rights, can we use resources without access to all data? What information should we have about this purpose? Even if the owners agree to provide certain data, it may be very difficult to move the data due to its large size. So, can we exploit the data without moving it? In addition, data in different locations may have different label quality, and may contain significant label sound, possibly due to crowdsourcing. Can we read low quality and / or even contradictory label information? In addition, we often assume that data is distributed uniformly and independently; however, the basic i.i.d. guesses cannot hold on different data sources. Can we learn more effectively and efficiently beyond i.i.d? thought? There are a few first studies on these important big data issues, including [4], [8], [6]. In addition, if the same data is provided, different users may have different needs. For example, in a product recommendation, some users may want the recommended items to be good, and some users may want all the recommended items to be good, while other users may want all the good things returned. Calculations and large data storage may be barriers to the development of each model for different needs separately.

## III. 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, including [1], [14], [3], [5], [6]. 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 petabytes—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 petascale 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. A similar caveat also plays with the concept of Veracity of data. The issue of data quality or veracity has been considered by a number of researchers [39], including data complexity [9], missing values [19], noise [58], imbalance [13], and dataset shift [39]. The latter, dataset shift, is most profound in the case of big data as the unseen data may present a distribution that is not seen in the training data. This problem is tied with the problem of Velocity, which presents the challenge of developing streaming algorithms that are able to cope with shocks in the distributions of the data. Again, this is an established area of research in the data mining community in the form of learning from streaming data [3], [48]. The key opportunity here is to take the academic literature for a testdrive in real industry settings where issues of scale and delivery often supersede the desire for accuracy. Depending on the application domain, a simpler model might be preferred, even if slightly less accurate. However, as demonstrated by the success of deep learning, computational advances are opening new doors to new opportunities. The issue with Variety is, undoubtedly, unique and interesting. A rapid influx of unstructured and multimodal data, such as social media, images, audio, video, in addition to the structured data, is providing novel opportunities for data mining researchers. We are seeing such data rapidly being collected into organizational data hubs, where the unstructured and structured cohabit and provide the source for all data mining. A fundamental question is related to integrating these varied streams or inputs of data into a singular feature vector presentation for the traditional learning algorithms. An example of big data that has the elements of the four V's is the social media and network data. The last decade has witnessed the boom of social media/network websites, such as Facebook, LinkedIn, and Twitter. Together they facilitate an increasingly wide range of human interactions that also provide the modicums of big data. The ubiquity of social networks manifests as complex relationships among individuals. It is generally believed that the research in this field will enhance our understandings of the topology of social networks and the patterns of human interactions [8], [18], [33], [36], [41], [54]. However, such data present numerous challenges from provenance (individual created data), veracity of data as the data is effectively crowd-sourced, volume as millions of individual are contributing content or making connections, and variety as the data not only comprises of the social network structure but also content such as text and images. Our call to the community is to reconvene some of the traditional methods and identify their performance benchmarks on ―big data‖, and identify novel directions for ground-breaking research built on the foundations we have already developed.

**3.1. From Data to Knowledge to Discovery to Action**

Recent times have greatly enhanced our ability to collect large amounts of data, giving us the opportunity to make dramatic changes in the way we analyze and understand data. This data outlines many features that can complement not only hypothesis-driven research but also enable the discovery of new ideas or events from rich data, which may include location data, temporary data, view data, various data sources. , text data, random data, etc. Such level data and character lengths present new scientific challenges driven by data charting the path from data to information to comprehension. This data-driven acquisition process will include an integrated system of descriptive analysis and predictable modeling for useful information or hypothetical ideas. These ideas not only relate to each other but also help to explain the origin or help to confirm something that has been noted. These observations or predictable analyzes may be helpful in informing decisions, including specific actions that cannot be properly measured in terms of cost and impact of an action. A set of alternating hypotheses leads to conditions that cannot be measured by status. Brynjolfsson et al. [11] studied 179 large companies and found that companies that received data-driven decisionmaking had a 5 to 6 percent high productivity rate. The main difference was that these companies relied on data and statistics instead of relying on experience and understanding. Health care is another area that testifies to the important use of big data. United Healthcare, for example,

is wasting effort on the behavior of mining clients as found in recorded voice files. The company uses natural language processing and text data to identify customer feelings and satisfaction. It is a vivid example of capturing large amounts of data, developing analytical models, and obtaining measurable and possible data. Big data presents unparalleled opportunities: accelerating scientific discovery and innovation; to improve health and well-being; constructing sections of a novel reading that may not have been available to date; improve decision-making by empowering data analysis; understanding the dynamics of human behavior; and the impact of trade on the global economy.

### 3.2. Opportunities and Challenges

Big data clearly presents us with exciting opportunities and challenges in data mining research. First, data-driven science and discoveries should try to find action-focused information that leads to drawing new discoveries or implications. Without understanding the nuances of personal data and background, one can fall into the trap of easy and misleading relationships, sometimes leading to false discovery and understanding. It is important to fully understand and appreciate the domain in which the person works, and all awareness and detail should be properly constructed in that domain. It requires a person immersion in the domain to perform feature engineering, data analysis, machine learning, and knowledge of system configuration and website design, as well as performing that-if analysis. This does not mean that a data scientist will be an expert in all aspects. Instead a data scientist may be the inventor of machine learning but who is well versed in the design of a system or website or visualize or perform rapid prototyping. But a data scientist cannot be separated from a domain without risking failure. Leading data acquisition into an application domain requires deep curiosity, the ability to ask big questions, and to deliver a variety of data sources, in addition to data science chopsticks. Second, algorithms are important, but before we jump into the novel algorithms journey to tackle the four Vs of big data, it is important for the public to consider the progress made so far, do a thorough research on it and identify potential. the challenges, challenges and obstacles of the present. In addition, it can be very beneficial to consider additional data source while considering a simple algorithm with an interest in answering questions about the domain. It's about completing the picture, after all. Third, any improvements in measuring algorithms should be accompanied by improvements in construction, systems, and the creation of new data. We are seeing a shift to technologies such as NoSQL and Hadoop, given the schemes-free environment of new data types and the spread of random data. It is an opportunity for algorithmic researchers to work with system researchers to integrate machine learning or data mining algorithms as part of a pipeline to naturally use low data storage structures and computational fabric. Fourth, the basic paradigm that exists before us is data driven. A data scientist should be an inquisitive outsider who can ask data questions, highlight the limitations of available data and identify additional data that can improve the performance of algorithms in a particular task. The hypothesis here is that there may be external data in the data taken by a particular company, which can provide a significant value. For example, consider the problem of predicting a patient's recovery when exiting. This problem of reduced learning can be worthwhile by considering lifestyle data, which is outside the patient's Electronic Medical Record (EMR). We see this as one of the key opportunities and challenges presented directly within the data mining environment of the vast data research context.

## IV. GLOBAL OPTIMIZATION WITH BIG DATA

Another key area where big data offers opportunity and challenges is global optimization. Here we aim to optimize decision variables over specific objectives. Meta-heuristic global search methods such as evolutionary algorithms have been successfully applied to optimize a wide range of complex, largescale systems, ranging from engineering design to reconstruction of biological networks. Typically, optimization of such complex systems needs to handle a variety of challenges as identified here.

### 4.1. Global Optimization of Complex Systems

Complex systems tend to have a large number of dynamic decisions and involve a large number of objectives, where the correlation between dynamic decisions may be less linear and objectives often conflict. Development problems with a large number of dynamic decisions, known as major development problems, are extremely challenging. For example,

the performance of many search algorithms will be greatly reduced as the number of decisions increases, especially if there is a complex correlation between dynamic decisions. Divide and conquer the widely accepted strategy for dealing with greater efficiency where the main problem is to find a correlation between dynamic decisions so that the related relationships are grouped by the same minorities and independent relationships are grouped into different categories. - people. Over the past two decades, metaheuristics have been shown to be effective in solving multi-objective development problems, where objectives often clash on their own. The main reason is that in a humanbased search method, different people can capture different trade relationships between conflicting targets, e.g., in improving the structure of a complex structure [2]. As a result, it is possible to gain a small set that represents a complete Pareto solution by doing one run, especially problems to improve two or three objectives. The multi-objective metaheuristics developed to date can be broadly divided into three categories, namely weight-based methods [8], Pareto-based methods [6] and index-based algorithms [5]. Unfortunately, none of these methods can work well if the number of objectives is much higher than three. This is because the number of Pareto-optimal solutions is huge and gaining a smaller set of them is no longer feasible. In weight-loss methods, it may be difficult to build a limited number of weight-bearing compounds to represent the best Pareto solutions of the highest magnitude. On Pareto-based approaches, many solutions for people of limited size are incomparable. Thus, few people dominate others and the pressure to choose the best solutions is lost. An additional complication is the growing calculation costs of establishing governance relationships as the number of objectives increases. Performance-based methods also have a major problem with calculations, e.g., in hyper-volume calculations. A second major challenge associated with the efficiency of complex systems is the expensive computer systems for evaluating the quality of solutions. In many complex preparation problems, time-consuming numerical simulations or costly tests need to be performed to assess eligibility. The costly calculation costs in a prohibited manner make it difficult to use search based search algorithms worldwide for such complex problems. Another promising approach is the use of computer-assisted models, known as surrogates, to replace the cost-effectiveness component of the test [2]. However, building surrogates can be very challenging in large problems with limited data samples that are too expensive to collect. Complex improvement issues are often subject to a large amount of uncertainty, such as different environmental conditions, system damage, or changing customer demand [9]. Two basic ideas can be adopted to deal with the uncertainty of performance. One is to find insensitive solutions to small dynamic variations of decision decisions or eligibility tasks, known as the right solid solutions [9]. However, if the changes are large and continuous, the meta-heuristics of the good moving average will often be improved, known as dynamic efficiency [4]. Different from the rigid approach to dealing with uncertainty, a strong improvement aims to keep track of the positive whenever it changes. Theoretically this sounds perfect, but it really isn't necessary for two reasons. First, good tracking is highly computer-based, especially if competency testing is expensive. Second, a change in design or solution may be costly and conventional changes are not allowed in most cases. Considering these two factors, a new approach to uncertainty has been proposed, called resilience over time [3]. The key is to achieve real-world transactions between finding a solid solution and good tracking. That is, the design or solution will only change if the current solution is not acceptable, and a suitable new solution will gradually change over time, which is not really the best solution at that moment.

### 4.2. Big Data in Optimization

Meta-heuristic global optimization of complex systems will not be possible without data generated in numerical measurements and physical examination. For example, improving the design of a race car is a major challenge as it involves many sub-systems such as the front wing, rear wing, chassis and tires. A large number of dynamic decisions are involved, which may slow down the search functionality of meta-heuristics. To alleviate this difficulty, the data generated by aerodynamic engineers in their daily work will greatly assist in determining which subsystem, or as a continuous step which sub-system component, is essential for aerodynamic development and driving. Analyzing and extracting such data, however, is a challenging task, as the amount of data is large, and the data may be stored in different ways and contaminated with noise. In other words, this data is fully reflected by the four Vs of the big data. In addition, as testing of the suitability of race car designs is time consuming, surrogates are important for the development of race cars. Another example is the reconstruction of computer networks for the control of biological elements. Reconstruction of gene control

networks can be seen as a complex development problem, in which a large number of parameters and network connections need to be determined. Although meta-heuristic development algorithms have been shown to be very promising, genetic reconstruction data is the largest data by nature [5]. Data obtained from genetic expression increase with exposure level [9]. Data volume continues to grow with the development of next generation strategic strategies such as high-impact testing. In addition, data from experimental biology, such as microarray data, are noisy, and genetic expression tests rarely have the same growth conditions and thus produce different data sets. Data diversity is also greatly enhanced by the use of deletion data, where the gene is removed to determine its control parameters. Disruption testing contributes to the reconstruction of genetic control networks, which, however, are another source of biological data variability. Data collected from different gene labs on the same biological network often differ. It is also very important to develop development algorithms to be able to obtain problem-specific information during development. Acquisition of problem-solving information can assist in capturing problem-solving to make the search more efficient. In large multi-purpose problems, such information can be used to direct search in a highly promising search environment, and to specify preferences over objectives so that the search will focus on the most important trade. Unfortunately, sometimes only limited information is available for the problem to be resolved. It is therefore interesting to obtain information from similar use efficiency problems or goals that have been previously resolved [20]. In this case, reusing information can be a real challenge. The relationship between the challenges in the development of complex systems and the big data environment is shown in Fig. 2.
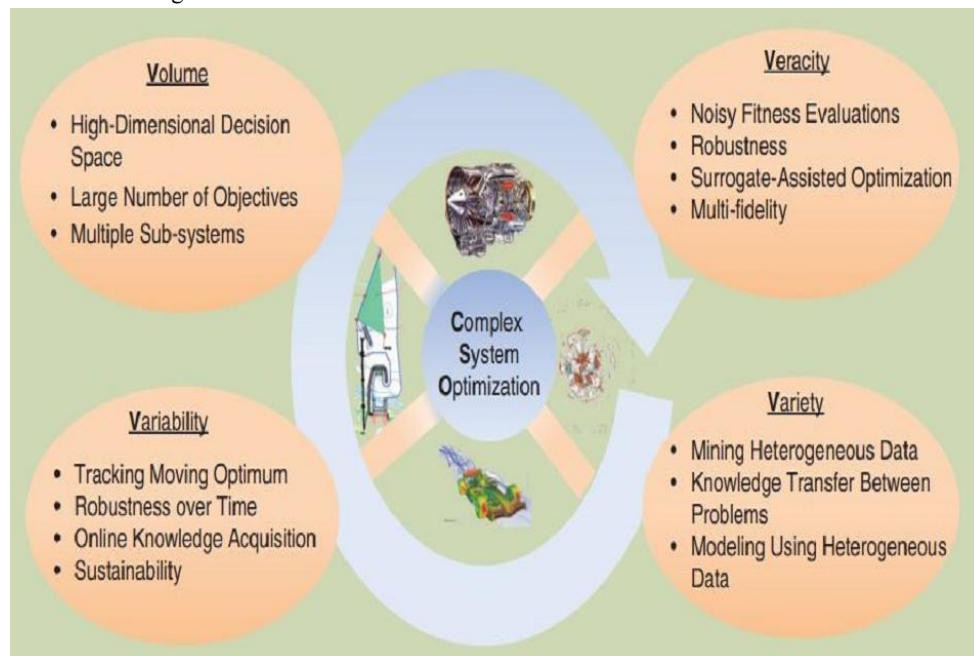


**Figure 2:** Relationship between the challenges in complex engineering optimization and the nature of big data.

## 4.3. Opportunities and Challenges

As discussed above, big data is widely regarded as critical to the success of complex design processes. Great effort has been devoted to the use of data to improve the performance of meta-heuristic optimization algorithms to solve major problems in the presence of large amounts of uncertainty. It is believed that growth in big data research could create new opportunities and pose new challenges to data-driven efficiency. Answering the following questions can be important in turning the challenges posed by big data into opportunities. First, how can we seamlessly integrate modern learning methods with fully functional techniques? Many learning strategies have developed, such as low-level reading [23], advanced reading [15], active reading [7] and n-depth reading [10] have been developed over the past decade. However, these techniques are rarely used profitably within development without a few exceptions, and are essential for acquiring

domain information with a large amount of diverse and noisy data. For optimal development using metaheuristics, such information is the ultimate decision in setting flexible and integrated problem representation, designing successful search operators, building high quality surrogates, and refining user preferences for multiple purposes. Second, how can we create a development problem so that new technologies developed in big data research can be used effectively? The traditional formulation of development problems consists of defining objective activities, changes in decisions. This applies perfectly to small, well-defined problems. Unfortunately, the construction of complex problems in itself is a recurring learning process. New ways of analyzing data on big data can be of interest in simplifying the formation of complex performance problems. For example, in order for people of other children to be used to predict quality in human-based development, predicting specific fitness is less important than finding related sequences for candidate designs. It is also possible to find variables of meta decisions that may be more effective in directing the search process than using the flexibility of the original decision? Third, how do we visualize a high decision space and a high resolution space to understand the solutions reached and make choices [7], [3]? How can advanced techniques in the analysis of big data be used in development? Overcoming these challenges with a large data framework will bring significant improvements to global overall performance in the years to come.

## V. INDUSTRY, GOVERNMENT AND SOCIETY WITH BIG DATA

We have presented above some of the technical challenges for research around the disciplines impacted and challenged by big data. In the end, we must also focus on the delivery of benefit and outcomes within industry, business and government. Over the decades, many of the technologies we have covered above, in machine learning, data mining, and global optimization, have found their way into a variety of large-scale applications. We now ask what are the impacts we see today in industry and government of big data, how is this affecting and changing society, and how might these changes affect our research across all these disciplines? In this section we present a perspective on data analytics from experiences in industry and government. The discussion is purposefully presented as a point of view of future practice rather than presenting a scientifically rigorous argument. We identify areas where a focus from research might deliver impact to industry, government and society.

### 5.1. Decentralizing Big Data

It is useful to point out that over the past two decades we have seen a time when the public has seen the collection of personal data in the interests of business and government. As users, we are enticed by the significant benefits of providing our data to these organizations, and these organizations now store big data that we have understood as a current marketing concept or term. Google, Apple and Facebook, as well as many other Internet companies that exist today, provide services ranging from finding old friends to the ability to share our thoughts, personal information and daily activities in public. With the abundance of our email, our diary and calendars, our photos and thoughts and personal activities, now hosted by Google, for example, there is a great opportunity to identify and address the entire customer experience on a large scale. Combine that with our web logs, updates from our location and archive of our documents in Google Drive, and we begin to understand the vast scope of the data collected about each of us, individually. This data can be used to better target the services advertised to us, using amazing algorithmic technology to deliver new information and information. Together, these multi-source data stores entice us to bring our personal data to data collectors in return for the sophisticated services they provide. Attraction, of course, looks amazing, evidenced by the large number of users in each growing online ecosystem. Customers who drive this data collection by these online companies are not service users but, for example, commercial advertisers and sophisticated System Development • High Demensional Decision Space • Large Number of Objectives • Multiple Lower Systems • Tracking Moving Optimum Internet • Sustainability • Various Mining Data • Information Transmission Between Problems • Modeling Using Different Data Volume Variations • Sound Efficiency Testing • Strength • Individual Assistance • Enhanced Enhancement • Multiple Reliability and Assurance in New Image government services. Data is also made available to other organizations, intentionally and / or improperly. As data is gradually being collected by cloud services over time, there is a growing need for broader discussion and understanding of the privacy, security, and public issues that arise from such collections of personal data. As a

community, we have reduced our focus on these issues and have come to realize that the central data store is the only way we can deliver these desirable services. This notion of a centralized collection of personal and private data should be challenged. The one-position model primarily serves the interests of data-gathering organizations -more than people. We see a growing interest and opportunity to transform this data collection approach into its own, and to give personal interests first, then second organizations. Appearance, for example, OwnClowd1 and SpiderOak, 2 as customized or encrypted changes to Google Drive and DropBox, reflect this trend. The gradual discovery of the importance of privacy leads to changes in how we store data. We will see this lead to some governments introducing new management practices to online companies over time, and companies themselves start developing and marketing the importance of protecting our data. We will also begin to see that data transfer can be stored in one place but with another person associated with it — the appropriate data owner. This widespread data distribution presents one of the biggest challenges for data scientists in the near future for big data. We present below three major data challenges related to this emerging paradigm of overarching data. Two challenges to starting a scale and the suitability of our building model http://en.wikipedia.org/wiki/Owncloud http://en.wikipedia.org/wiki/Spideroaking. The main focus however is a challenge for the community in the coming years when data is transported from one place to the most widely distributed data stores. This is one of the most important challenges to be presented in the future of big data, and it has a huge impact on the way we think about machine learning, data mining and global efficiency.

### 5.2. Scaled Down Targeted Sub-Models

There has been a lot of focus on the need to upgrade our machine learning algorithms to build models for the entire population available — and build them faster. Indeed, from the very beginning of data mining [45] the key focus has been to measure algorithms. This is what we used to identify as the distinguishing factor between data mining (or sitebased data acquisition) and the home studies of the many algorithms we used: machine learning and mathematics [56]. Our goal was to use all available data, avoid the need to take samples, and make sure we capture information for everyone. The principle remains with us today, as we continue to be very busy and able to collect data — lots of data — and thus introduce new businesses and refine old ones. However, in the last few decades we have talked about big data, big, big, big, big, funny. The current trend is to refer to it as big data or big data [40]. No matter, it's just a lot of data. In business and government our data sets today contain anything from 100 notifications of 20,000 variables, to 20 million alerts per 1,000 variables, to one billion out of 10,000 variables. Such large data sets challenge any algorithm. Although our research usually focuses on algorithms, the challenges presented in data collection, storage, management, deception, refinement, and conversion are often much larger (i.e., time-consuming) than the actual model structure. You are challenged with big data, what is the job, of course, of a data scientist? In the world to come, learning algorithms will penetrate our vast databases — they will cut and sell data, and they will identify confusing, patterns, and behaviors. But how will this be delivered? The productive approach would be to build on the first successful concept of building an integrated model [55], but taken on a large new scale. In fact, we are seeing the development of mechanisms that severely divide our data into smaller, more subdivided sets, representing more subdivisions, often targeting behavioral archetypes. Within these sub-regions we can better understand and model the many behaviors that businesses are interested in. The concept is not new [57] but has received little attention over the years until now when we see the need to cut and sell big data on logical ways to reveal these vast amounts of behavioral patterns. Today in industry and government this approach brings new models that show amazing results based on collections of thousands of very different small local models. The challenge here is how to identify areas of low morality where we build our models. Instead of building a predictive model on a scale over big data we are building a community of thousands of smaller models, which as a group apply to the entire population. This is done by understanding and dealing with the nuances and idiosyncrasies of different small areas. It is within a very small number of people that speculative models, for example, are built. The final model becomes a combination of the individual models used in each new observation. The most recent successful implementation of this approach, which has been implemented in practice, has analyzed more than two million of the 1000 dynamic variables to identify 20,000 such sub-areas of behavior. The sub-spaces are created using a combination of analysis analysis and the introduction of the decision tree, and each sub-area is illustrated, each identifying and representing the new concepts.

Sub-spaces can be overlapping —each focus may be for multiple sub-areas (i.e., one view may show multiple behaviors). In each of the 20,000 vacancies, smaller speculative models can be developed, and by combining these into a whole — using global development for more complex purpose functions — we bring a more effective global model, which is used, for example, to point the whole risk. Australians who pay taxes, as they connect online to submit their tax returns. We can (and no doubt) continue to explore new computer paradigms such as in-site analytics to greatly measure machine learning and mathematical algorithms in big data. But the big game will be in recognizing and modeling the many facets of the various behaviors we all display individually, and share in small numbers, where the modeling itself will be done.

### 5.3. Right Time, Real Time, Online Analytics

The traditional data miner is, in fact, involved in the development of a batchoriented model, using machine learning and mathematical algorithms. From our historical data store we use algorithms such as retrieval of objects, decision tree formulation, random forests, sensory networks, and vector support machines. Once we have built our model (s) then we continue to produce models. In a business environment we are moving our models into production. The models then work on new tasks as they arrive, perhaps finding each job and deciding on the treatment of that purchase — that is, based on the outcome, how should our systems be considered? The process is typical of how data mining is delivered to many large organizations today, including government, financial institutions, insurance companies, healthcare providers, marketing, and so on. At the Australian Taxation Office, for example, every day a collection of data mining models compromises on all transactions (tax returns) received. The Australian Department of Immigration [4], as another example, has developed a risk model that tests all passengers when boarding their international flight to Australia (known as Advance Passenger Processing) using data mining models. Such examples are pervasive in industry and government. Today's fast-paced situation now requires more than that. Large and old organizations, all over the world, have preferred not to be too quick in their ability to respond to the rapid changes brought about online by our data-rich world. Organizations no longer have the luxury of spending a few months building scoring models in batch mode. We need to be able to evaluate each activity as it happens, and learn flexibly as the model interacts with the large number of tasks we face as it happens. We need to build models in real time to respond in real time and learn and change their behavior in real time. Growth research studies do not really start. Growing learning [15], [6], reading at any time or at any time [9], and data mining mines [21] have all addressed the same problems in different ways over the decades. Now there is a growing opportunity to use this approach. The question remains as to how we can maintain and improve our information store over time, and work to forget old, potentially inaccurate information? The development of flexible, fastmoving students in real time — that is, as they interact with the real world — remains a challenge and will continue to be a major challenge in our big data world.

### 5.4. Extreme Data Distribution: Privacy and Ownership

After considering the two central challenges surrounding big data, we now consider the challenge of game modification. The future holds us hopeful that people will once again manage their data in the hands of big central stores. We expect to see this as a systematic development of our understanding of what is best for society as we develop and govern our society in this age of massive data collection and surveillance and its high risk of privacy. Data ownership has become a challenging problem in our data-rich world. Data collectors in business and government sectors are learning to properly collect large amounts of big data. However, with the help of civil libertarians, philosophers and whistle blowers, the public is gradually becoming more and more aware of the need for better governance in data collection and processing. Recent events such as Wikileaks3 and Edward Snowden4 help raise the level of discussion that provides insight into the dangers of data integration in the middle — regardless of who owns the data. We are well aware of the dangers of a single point of failure — relying on our one-sizefits-all data, such as big data sets, securely stored, and used by data collectors for the benefit of our community, and the individuals of that community. . However, one point of failure will mean that just one mistake or one violation can lead to devastating consequences on a large scale. And with more sophisticated attacks, they happen more often. Even though it is not complicated, Snowden has shown that having all the data stored in one place increases the risk significantly, harming the industry, governments, and society as a whole. After identifying

the risks, we usually work out strategies to reduce those risks. The obvious strategy is to withdraw from the unsafe natural collection of data collection. Personal data needs to go back to people whose data they own. We can and should collect and store such data ourselves, under our control, as individuals, minimizing total social risk. Attractive and reliable services, services provided by Google, Apple, Facebook, and many other online ecosystems, are yet to be introduced. Companies may retain their ability to benefit from data, while data itself is retained by data owners. Under this future scenario, instead of summarizing all the calculations, we can bring the calculations to intelligent operators who use our own resources and liaise with service providers. Business models can still be profitable and users can regain their privacy. Our personal data will be locked after encryption technology, and people will hold the key to unlock the data as they wish. Data will be provided on our smart portable devices where encryption is removed and provides the services (email, photos, music, instant messaging, purchases, searches, web logs, etc.) that we need. Data will be processed on these encrypted personal clouds, which operate on smart mesh-network devices connected to smart devices. We see the beginning of this change happening now with projects like Freedom Box, 5 OwnCloud, and IndieWeb, as well as the widespread and wide acceptance of powerful smartphones! This idea of big distributed data poses a challenge to developing technologies and algorithms that work on more advanced data distribution. How can we spread the word about this extreme level of data distribution and continue to build models that learn in the context of big data? The challenge of overdistributed data and learning is one that will grow rapidly in the next few years. It will require a very different approach to the development of machine learning, data mining and global efficiency. Compare this with how we can look at a group of people working together — we individually own a large data store and share and read and use that data as we interact with the world and other people. The same is true of future learning algorithms.

## 5.5. Opportunities and Challenges

Identifying significant data big challenges leads one to question how the future of data collection could emerge in the years to come. As the pendulum reaches the limit of large data collection in one place, over time, and possibly faster than we expected, we will see the pendulum begin to recede. It should turn back to a state where we return data to the ownership and control of people related data. Data will be widely distributed, and individual records will be distributed far and wide, just as data owners will be distributed far and wide throughout our world. With the overcrowding of data we will be challenged to provide the services we expect with massive central data storage. Those challenges are certainly insurmountable, but it will require a lot of research, innovation, and software development to deliver. The first challenge presented was one of properly classifying big data (ultimately large overdistributed data), identifying the behavioral groups we study in, and even modeling and learning at each level. The second challenge is to re-focus on bringing learning algorithms that you learn (yourself) in real time, or at least in real time, and to do this online. Lastly, how do we bring this to a point where extreme data distribution where database records are now widely distributed and protected by privacy, and how we can bring learning agents who care for the interests of their ﹣owner‖ .

## IV. WRAP UP

From the data analysis ideas we have presented here, there are many new opportunities and challenges brought by big data. Some of this is not really new, but it is news that has not received the attention it deserves. Here we recall some of the most important / exciting things:

- Data size: On the other hand, we create "oneworld reading" algorithms that require only one data scanner with limited storage size unrelated to data size; on the other hand, we are trying to identify small portions of really valuable data in the actual big data.
- Data variability: Data manifests itself in different ways of a particular concept. Introduces a new concept in learning programs and computer intelligence algorithm for differentiation, where the feature vector is multi-modal, has a structured and unstructured text, and the concept is still to distinguish one concept from another. How do we create a feature vector, and then a learning algorithm with the right purpose function to learn from such a wide variety of data?

- Data Reliability: Although data is available quickly and growing, it is also important to consider the source of the data and if the data is unreliable. Additional data is not accurate data, and additional data is not relevant data. A sharp data filter is key.
- Distributed life: Owners of different pieces of data may grant different access rights. We should aim to maximize data sources without access to all data, and exploit them without moving data. We will need to take into account the fact that different sources may come with different label quality, there may be a lot of noise in the data due to mass availability, and i.i.d. guesses may not end at all sources.
- Extreme distribution: If we take this idea even further, unit level data may be what we see as data distribution, as we deal with privacy and security issues. New ways to model big data will be needed to work with over-distributed data.
- Different needs: People may have different needs while the high cost of big data processing can hamper the creation of a different model for each need. Can we build a single model to meet different needs? We also need to note that, with big data, it could happen in finds supporting evidence in any argument we want; so, how do you judge / evaluate "our findings"?
- Sub-models: Different needs may also relate to the diversity of behaviors we emulate within the context of our application. Rather than one model encompassing it all, the model will contain ensembles of a large number of smaller models that will bring better understanding and prediction than a single, more complex model.
- Emotional significance: Data will improve the acquisition of novels and action-specific business details. It is important that you still attach intuition, curiosity and background information otherwise the person may become myopic and fall into the trap of "communication is enough". Computer intelligence should be integrated with human understanding.
- Fast model: As the world continues to ―accelerate‖, decisions need to be made quickly because fraudsters can quickly find new ways in an aging environment, model construction must be fast and real-time.
- Extensive preparation: Global development algorithms such as meta-heuristics have achieved great success in educational research, but are rarely used in the industry. One major obstacle is the huge computer costs required to assess the quality of candidate designs for sophisticated engineering systems. The growing technology of data analysis will remove the barrier to some extent by reusing information extracted from large amounts of high, varied and noisy data. Such information can also be acquired through new visual aids. Big datadriven efficiency will also play an important role in the reconstruction of large biological systems.
- Sophisticated efficiency: The interpretation of dynamic decisions, setting goals and articulating issues are three main steps in creating development problems before solving them. In order to develop sophisticated systems, the construction of the performance problem itself becomes a complex performance problem.

## VII. CONCLUSION

Big data may provide us with new information and ways to create development problems, thus leading to a more efficient solution. In closing the discussion, we emphasize that the opportunities and challenges posed by big data are very broad and varied, and it is clear that there is no one way that can meet all the needs. In this sense, big data also brings the opportunity for a ―great combination ‖ of strategies and research.

## REFERENCES

[1]. J. Abello, P. M. Pardalos, and M. G. Resende, Handbook of Massive Data Sets (Massive Computing). vol. 4, New York: Springer, 2002.

[2]. C. C. Aggarwal and P. S. Yu, Eds., Privacy-Preserving Data Mining: Models and Algorithms. New York: Springer, 2008.

[3]. C. C. Aggarwal, Data Streams: Models and Algorithms. vol. 31, New York: Springer, 2007.

[4]. (2013). Australian Dept. Immigration. Fact sheet 70— Managing the border Internet. [Online]. Available: http:// www.immi.gov.au/media/fact-sheets/70border.htm

**[5].** J. Bader and E. Zitzler, ―HypE: An algorithm for fast hypervolume-based manyobjective optimization,‖ Evol. Comput., vol. 19, no. 1, pp. 45–76, 2011.

**[6].** P. Baldi and P. J. Sadowski, ―Understanding dropout,‖ in Advances in Neural Information Processing Systems 26. Cambridge, MA: MIT Press, 2013, pp. 2814–2822.

**[7].** M. Banko and E. Brill, ―Scaling to very very large corpora for natural language disambiguation,‖ in Proc. 39th Annu. Meeting Association Computational Linguistics, Toulouse, France, 2001, pp. 26–33.

**[8].** A.-L. Barabási, Linked: The New Science of Networks. New York: Basic Books, 2002.

**[9].** M. Basu and T. K. Ho, Data Complexity in Pattern Recognition. London, U.K.: Springer, 2006.

**[10].** Y. Bengio, ―Learning deep architectures for AI,‖ Foundations Trends Mach. Learn., vol. 2, no. 1, pp. 1–127, 2009.

**[11].** E. Brynjolfsson, L. Hitt, and H. Kim. (2011). Strength in numbers: How does data-driven decision making affect firm performance? [Online]. Available: http://ssrn.com/abstract=1819486

**[12].** T. Chai, Y. Jin, and S. Bernhard, ―Evolutionary complex engineering optimization: Opportunities and challenges,‖ IEEE Comput. Intell. Mag., vol. 8, no. 3, pp. 12–15, 2013.

**[13].** N. V. Chawla, ―Data mining for imbalanced datasets: An overview,‖ in Data Mining and Knowledge Discovery Handbook, O. Maimon and L. Rokach, Eds. New York: Springer, 2005, pp. 853–867.

**[14].** N. V. Chawla, Learning on Extremes-Size and Imbalance-of Data. Tampa, FL: Univ. South Florida, 2002.

**[15].** Q. Da, Y. Yu, and Z.-H. Zhou, ―Learning with augmented class by exploiting unlabeled data,‖ in Proc. 28th AAAI Conf. Artificial Intelligence, Quebec City, QC, Canada, 2014.

**[16].** K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, ―A fast and elitist multiobjective genetic algorithm: NSGAII,‖ IEEE Trans. Evol. Comput., vol. 6, no. 2, pp. 182–197, 2002.

**[17].** T. G. Dietterich, ―Approximate statistical tests for comparing supervised classification learning algorithms,‖ Neural Comput., vol. 10, no. 7, pp. 1895–1923, 1998.

**[18].** D. Easley and J. Kleinberg, Networks, Crowds, and Markets. Cambridge, U.K.: Cambridge Univ. Press, 2010.

**[19].** Farhangfar, L. Kurgan, and J. Dy, ―Impact of imputation of missing values on classification error for discrete data,‖ Pattern Recognit., vol. 41, no. 12, pp. 3692– 3705, 2008.

**[20].** L. Feng, Y.-S. Ong, I. Tsang, and A.-H. Tan, ―An evolutionary search paradigm that learns with past experiences,‖ in Proc. IEEE Congr. Evolutionary Computation, Brisbane QLD, Australia, 2012, pp. 1–8.

**[21].** M. M. Gaber, A. Zaslavsky, and S. Krishnaswamy, ―Mining data streams: A review,‖ ACM SIGMOD Record, vol. 34, no. 2, pp. 18–26, 2005.

**[22].** W. Gao, R. Jin, S. Zhu, and Z.-H. Zhou, ―Onepass AUC optimization,‖ in Proc. 30th Int. Conf. Machine Learning, Atlanta, GA, 2013, pp. 906–914.

**[23].** W. Gao and Z.-H. Zhou, ―Dropout rademacher complexity of deep neural networks,‖ CORR abs/1402.3811, 2014.

**[24].** J. A. Hanley and B. J. McNeil, ―A method of comparing the areas under receiver operating characteristic curves derived from the same cases,‖ Radiology, vol. 148, no. 3, pp. 839–843, 1983.

**[25].** G. E. Hinton, P. Dayan, B. J. Frey, and R. M. Neal, ―The ―wake-sleep‖ algorithm for unsupervised neural networks,‖ Science, vol. 268, no. 5214, pp. 1158–1161, 1995.

**[26].** G. E. Hinton and R. R. Salakhutdinov, ―Reducing the dimensionality of data with neural networks,‖ Science, vol. 313, no. 5786, pp. 504–507, 2006.

**[27].** H. Ishibuchi, M. Yamane, N. Akedo, and Y. Nojima, ―Many-objective and many-variable test problems for visual examination of multiobjective search,‖ in Proc. IEEE Congr. Evolutionary Computation, Cancun, QROO, 2013, pp. 1491–1498.

**[28].** H. Ishibuchi and T. Murata, ―Multi-objective genetic local search algorithm,‖ in Proc. IEEE Int. Conf. Evolutionary Computation, Nagoya, Japan, 1996, pp. 119–124.

**[29].** Y. Jin and J. Branke, ―Evolutionary optimization in uncertain environments—A survey,‖ IEEE Trans. Evol. Comput., vol. 9, no. 3, pp. 303–317, 2005.