

DATA SCIENCE

Bharathy A¹, Dilnisha K², Haneena Yousaf³

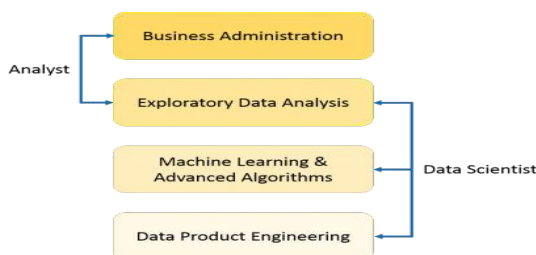
B.Tech 111 Year, Dept. of Computer Science and Engineering

Abstract-Data science has attracted heaps of attention, promising to show immense amounts of information into helpful predictions and insights. during this article, we have a tendency to raise why scientists ought to care regarding data science. To answer, we have a tendency to discuss data science from 3 perspectives: statistical, computational, and human. although each of the 3 is a critical component of data science, we have a tendency to argue that the effective combination of all 3 elements is that the essence of what data science is regarding.

Keywords-Predictive casual analytics, prescriptive analytics ,Business intelligence, artificial intelligence, Decision making, machine learning, clustering, deep learning, exploratory data analysis.

I. INTRODUCTION

Data Science may be a term that escapes any single complete definition, that makes it tough to use, especially if the goal is to use it properly. Most articles and publications use the term freely, with the belief that it's universally understood. However, data science – its strategies, goals, and applications – evolve with time and technology. Data science twenty five years past referred to gathering and improvement datasets then applying statistical strategies to that data. In 2018, data science has grown up to a field that encompasses data analysis, predictive analytics, data mining, business intelligence, machine learning, and so way more. Traditionally, the data that we have a tendency to had was mostly structured and little in size, that may well be analyzed by using the easy BI tools. In contrast to data within the ancient systems that was largely structured, these days most of the data is unstructured or semi-structured. Let's have a glance at the data trends in the image given below which shows that by 2020, over eighty you look after the data are unstructured. This data is generated from totally different sources like monetary logs, text files, multimedia forms, sensors, and instruments. easy BI tools don't seem to be capable of processing this immense volume and sort of data. this is why we'd like a lot of complex and advanced analytical tools and algorithms for processing, analyzing and drawing significant insights out of it.



As you'll see from the above image Data Analyst sometimes explains what's happening by process history of As you'll see from the higher than image, a knowledge Analyst sometimes explains what's happening by process history of the data. On the opposite hand, data scientist not solely does the preliminary analysis to get insights from it, however conjointly uses numerous advanced machine learning algorithms to spot the occurrence of a selected event within the future. a data scientist can inspect the information from several angles, typically angles not identified earlier. So, data Science is primarily accustomed make selections and predictions creating use of predictive casual analytics, prescriptive analytics (predictive plus decision science) and machine learning.

1.1 Predictive casual analytics

If you want a model which can predict the possibilities of a particular event in the future, you need to apply predictive causal analytics. Say, if you are providing money on credit, then the probability of customers making future credit payments on time is a matter of concern for you. Here, you can build a model which can perform predictive analytics on the payment history of the customer to predict if the future payments will be on time or not.

1.2 Prescriptive analytics

If you want a model which has the intelligence of taking its own decisions and the ability to modify it with dynamic parameters, you certainly need prescriptive analytics for it. This relatively new field is all about providing advice. In other terms, it not only predicts but suggests a range of prescribed actions and associated outcomes.

The best example for this is Google's self-driving car which I had discussed earlier too. The data gathered by vehicles can be used to train self-driving cars. You can run algorithms on this data to bring intelligence to it. This will enable your car to take decisions like when to turn, which path to take, when to slow down or speed up.

1.3 Business Intelligence (BI)

BI is the process of analyzing and reporting historical data to guide future decision-making. BI helps leaders make better strategic decisions moving forward by determining what happened in the past using data, like sales statistics and operational metrics.

1.4 Artificial Intelligence (AI)

AI computer systems can perform tasks that normally require human intelligence. This doesn't necessarily mean replicating the human mind, but instead involves using human reasoning as a model to provide better services or create better products, such as speech recognition, decision-making and language translation.

1.5 Machine Learning

A subset of AI, machine learning refers to the process by which a system learns from inputted data by identifying patterns in

that data, and then applying those patterns to new problems or requests. It allows data scientists to teach a computer to carry out tasks, rather than programming it to carry out each task step-by-step. It's used, for example, to learn a consumer's preferences and buying patterns to recommend products on Amazon or sift through resumes to identify the highest-potential job candidates based on key words and phrases.

1.6 Clustering

Clustering is classification but without the supervised learning aspect. With clustering, the algorithm receives inputted data and finds similarities in the data itself by grouping data points together that are alike.

1.7 Deep Learning

A more advanced form of machine learning, deep learning refers to systems with multiple input/output layers, as opposed to shallow systems with one input/output layer. In deep learning, there are several rounds of data input/output required to assist computers to solve complex, real-world problems. A deep dive can be found.

1.8 Exploratory Data Analysis (EDA)

EDA is often the first step when analysing datasets. With EDA techniques, data scientists can summarize a dataset's main characteristics and inform the development of more complex models or logical next steps.

1.9 Decision Science

Under the umbrella of data science, decision scientists apply math and technology to solve business problems and add in behavioural science and design thinking (a process that aims to better understand the end user).

2. RELATED WORK

Our study was inspired by prior work in end-user programming, teaching data science, practitioners as instructors, and broadening computing education to learners who do not self-identify as programmers. Data Science and End-user Programming Data science is a broad term that encompasses a wide variety of activities related to acquiring, cleaning, processing, modelling, visualizing, and presenting data [35, 41]. Although data visualization is a highly active area of HCI research, what is more relevant to our study is prior HCI research on programming as performed by non-professional programmers. Kandel et al. found great variation in levels of programming ability amongst data scientists [41]. Many of them write code in languages such as Python and R [21, 25, 35, 37], but they are not professional software engineers; moreover, many do not even have formal training in computer science. Much of data scientists' coding activities can be considered end-user programming [46] since they often write code for themselves as a means to gain insights from data rather than intending to produce reusable software artifacts. Related terms for this type of insight-driven coding activity include exploratory programming (from Kery et al. [42, 43]) and research programming (from Guo's dissertation [32]). However, as we discovered in our study, modern data scientists are not merely writing ad-hoc prototype code. They are now developing increasingly mature technology stacks for writing modular and reusable software (e.g., Figure 1). In the terminology of Ko et al., they are now engaging in end-user

software engineering [46] with more of an emphasis on code quality and reuse; in Segal's related terminology, data scientists are now becoming professional end-user developers [65]. Along these lines, software engineering researchers such as Kim et al. have studied the role of data scientists within industry engineering teams [44]. In contrast to prior HCI work that focuses on what data science practitioners do on the job, our study instead focuses on how they pass on those skills to novices via teaching.

3. DRAWBACKS

While Data Science is a very lucrative career option, there are also various disadvantages to this field. In order to understand the full picture of Data Science, we must also know the limitations of Data Science. Some of them are as follows:

1. Data Science is Blurry Term

Data Science could be a terribly general term and doesn't have a particular definition. whereas it's become a bunk, it's terribly arduous to write down the precise that means of a data scientist. a data Scientist's specific role depends on the sector that the company is specializing in. whereas some individuals have represented data Science to be the fourth paradigm of Science, few critics have referred to as it a mere rebranding of Statistics.

2. Mastering Data Science is near to impossible

Being a mixture of many fields, data Science stems from Statistics, computer science and mathematics. It's far away from attainable to master every field and be equivalently knowledgeable altogether of them. Whereas several on-line courses are attempting to fill the skill-gap that the data science industry is facing, it's still unacceptable to be skilled at it considering the immensity of the sphere. Someone with a background in Statistics might not be ready to master computer science on short notice so as to become a skilled data scientist. Therefore, it's a associate degree dynamical, field that needs the person to stay learning the various avenues of data Science.

3. Large Amount of Domain Knowledge Required

Another disadvantage of data Science is its dependency on Domain knowledge. An individual with a substantial background in Statistics and computer science can realize it troublesome to resolve data Science problem without its background knowledge. Constant holds true for its vice-versa. For example, A health-care business acting on working on analysis of genomic sequences would require an appropriate worker with some knowledge of genetics and molecular biology. This enables the data Scientists to create calculated choices so as to help the company. However, it becomes troublesome for { a data | a data | an information } scientist from

a different background to accumulate specific domain knowledge. This additionally makes it troublesome to migrate from one business to another.

4. Arbitrary Data May Yield Unexpected Results

A Data Scientist analyses the data and makes careful predictions in order to facilitate the decision-making process. Many times, the data provided is arbitrary and does not yield expected results. This can also fail due to weak management and poor utilization of resources.

5. Problem of Data Privacy

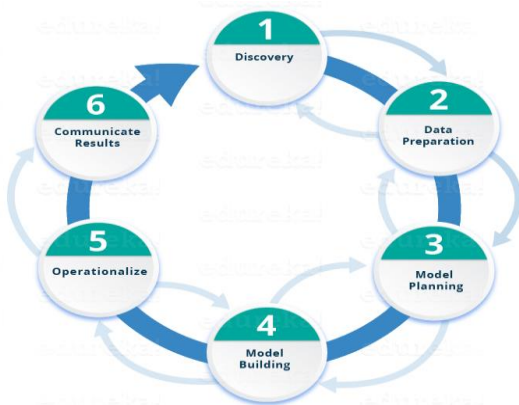
For many industries, data is their fuel. Data Scientists help companies make data-driven decisions. However, the data utilized in the process may breach the privacy of customers. The personal data of clients are visible to the parent company and may at times cause data leaks due to lapse in security. The ethical issues regarding preservation of data-privacy and its usage have been a concern for many industries.

4. BUSINESS INTELLIGENCE (BI) v/s DATA SCIENCE

- BI basically analyses the previous data to search out hindsight and insight to explain the business trends. BI allows you to take data from external and internal sources, prepare it, run queries on that and create dashboards to answer the questions like quarterly revenue analysis or business issues. BI can evaluate the impact of certain events within the near future.
- Data Science may be more innovative approach, associate exploratory method with the main focus on analysing the past or current data and predicting the long run outcomes with the aim of creating enlightened decisions. It answers the open-ended queries on “what” and “how” events occur.

5. LIFE CYCLE OF DATA SCIENCE

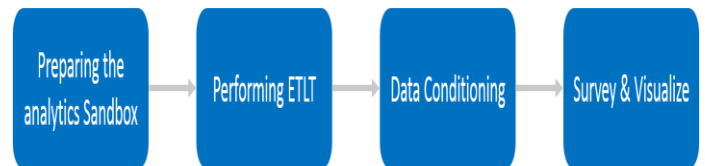
Here is a brief overview of the main phases of the Data Science Lifecycle:



Phase 1—Discovery: Before you start the project, it's vital to know the assorted specifications, needs, priorities and required budget. you need to possess the power to raise the proper queries. Here, you assess if you have got the desired resources present in terms of individuals, technology, time and data to support the project. during this part, you also have to be compelled to frame the business problem and formulate initial hypotheses (IH) to test.



Phase 2—Data preparation: In this phase, you require analytical sandbox in which you can perform analytics for the entire duration of the project. You need to explore, preprocess and condition data prior to modeling. Further, you will perform ETLT (extract, transform, load and transform) to get data into the sandbox. Let's have a look at the Statistical Analysis flow below.

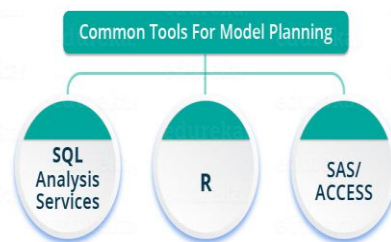


You can use R for data cleaning, transformation, and visualization. This will help you to spot the outliers and establish a relationship between the variables. Once you have cleaned and prepared the data, it's time to do exploratory analytics on it. Let's see how you can achieve that.



Phase 3—Model planning: Here, you will determine the methods and techniques to draw the relationships between variables. These relationships will set the base for the algorithms which you will implement in the next phase. You will apply Exploratory Data Analytics (EDA) using various statistical formulas and visualization tools.

Let's have a look at various model planning tools.



1.**R** has a complete set of modelling capabilities and provides a good environment for building interpretive models.

2.**SQL Analysis services** can perform in-database analytics using common data mining functions and basic predictive models.

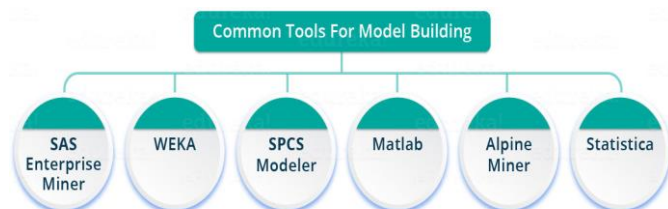
3. **SAS/ACCESS** can be used to access data from Hadoop and is used for creating repeatable and reusable model flow diagrams.

Although, many tools are present in the market but R is the most commonly used tool.

Now that you have got insights into the nature of your data and have decided the algorithms to be used. In the next stage, you will apply the algorithm and build up a model.

Phase 4—Model building: In this phase, you will develop datasets for training and testing purposes. You will consider whether your existing tools will suffice for running the models or it will need a more robust environment (like fast and parallel processing). You will analyse various learning techniques like classification, association and clustering to build the model.

You can achieve model building through the following tools.



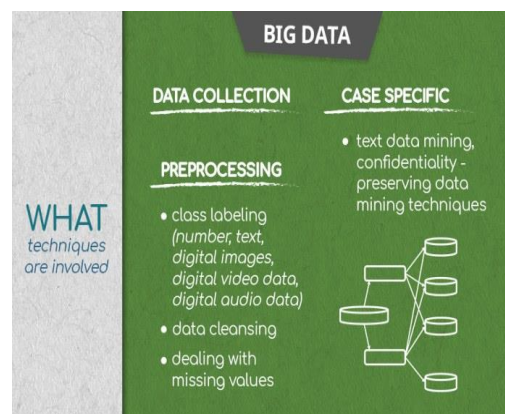
Phase 5—Operationalize : In this phase, you deliver final reports, briefings, code and technical documents. In addition, sometimes a pilot project is also implemented in a real-time production environment. This will provide you a clear picture of the performance and other related constraints on a small scale before full deployment.

Phase 6—Communicate results: Now it is important to evaluate if you have been able to achieve your goal that you had planned in the first phase. So, in the last phase, you identify all the key findings, communicate to the stakeholders and determine if the results of the project are a success or a failure based on the criteria developed in Phase 1.

6. BIG DATA IN DATA SCIENCE

When it comes to big data and data science, there is some overlap of the approaches used in traditional data handling, but there are also a lot of differences.

First of all, big data is stored on many servers and is infinitely more complex.



In order to do data science with massive data, pre-processing is even a lot of crucial, as the quality of the data may be a ton larger. you may notice that conceptually, a number of the steps are just like traditional data pre-processing, however that's inherent to working with data.

- Collect the data
- Class-label the data

Keep in mind that massive data is extremely varied, instead of rather than 'numerical' v/s 'categorical', the labels are 'text', 'digital image data', 'digital video data', digital audio data', and so on.

- data cleansing

The strategies here are massively varied, too; as an example, you'll be able to verify that a digital image observation is prepared for processing; or a digital video, or...

- data masking

When collecting data on a mass scale, this aims to make sure that any confidential information within the data remains private, without hindering the analysis and extraction of insight. the process involves concealing the original data with random and false data, allowing the scientist to conduct their analyses without compromising personal details. Naturally, the scientist will do this to traditional data too, and sometimes is, however with massive data the data are often much more sensitive, which masking a lot more urgent.

Where will data come from?

Traditional data could come back from basic client records, or historical stock price information.

Big data, however, is all-round us. A systematically growing variety of firms and industries use and generate massive data. think about on-line communities, as an example, Facebook, Google, and LinkedIn; or money

commerce data. Temperature measuring grids in numerous geographical locations also amount to massive data, in addition as machine data from sensors in industrial instrumentality. And, of course, wearable tech.

6. APPLICATIONS OF DATA SCIENCE

Application of data science is numerous. The foremost common use is in finance, genetics, banking, medicine, business and transportation for issues like financial trading, credit scoring, fraud detection, online advertising, web search, recommendations for cross-selling, etc. Several companies have targeted their business on data. They use data to {find |to seek out |to search out} hidden patterns which will facilitate them to find applicable solutions and improve decision-making process. This could facilitate organizations perceive their customers, markets, and therefore the business as a whole by anticipating growth, trends and business insights based on vast amounts of data. Principles and techniques of data science are also applied to general client relationship management to analyse client behaviour so as to cut back attrition and to increase expected client price. Within the finance industry, data science is employed for credit scoring, fraud detection and trading. In the healthcare, classification algorithms may be used to notice cancer and tumours at an early stage using Image Recognition software. In Genetic Industries, data science is employed for analysing and classifying patterns of genomic sequences. Using Machine Learning, data Scientists have developed recommendation systems that recommend completely different merchandise to customers based on their historical habits. In producing, industrial robots use data Science technologies like Reinforcement Learning and Image Recognition to take over repetitive jobs. In transport, Self-Driving Cars are developing based on Reinforcement Learning and Detection algorithms. Another application of data science is in colloquial agents (Siri by Apple). They use Speech Recognition system to know users, to convert human speech into textual data and to produce an acceptable response.

7. CONCLUSION

With the enormous increase in data, there is a constant need for analysing such a large amount of data. Data Science will manage this data and develop useful machine learning models that predict future results.

We are able to conclude that data Science is rising multidisciplinary field with roots in mathematics, statistics, and computer science. Because it engages in extracting, analysing, visualizing, managing and storing giant amounts of data, it's a really wide selection of application from business and finance to healthcare and transportation. The most goal

{of data |of knowledge |of information} Scientists is to acknowledge and use significant insights from data so as to assist organisations in taking smarter selections. Throughout that method, they use completely different tools and strategies to spot redundant patterns and hidden data inside the information. They conjointly use the foremost powerful hardware, most effective algorithms and programming systems to resolve the data related issues. During this paper, we want to introduce data science as a new, powerful field with various applications which will offer a competitive advantage and long-term stability.

REFERENCES

1. Barber, M. (2018). Data science concepts you need to know! Part 1. Retrieved 15. 9. 2019 from <https://towardsdatascience.com/introduction-to-statistics-e9d72d818745>.
2. Dataflair Team. (2019). What is data science?: a complete data science tutorial for beginners [Blog]. Retrieved 8. 10. 2019 from <https://data-flair.training/blogs/what-is-data-science/>.
3. Dhar, V. (2013). Data science and prediction. *Communications of the ACM*, 56(12), 64–73.
4. Foote, K. D. (2016). A brief history of data science. Retrieved 17. 10. 2019 from <https://www.dataversity.net/brief-history-data-science/#>.
5. Grossmann, W. and Rinderle-Ma, S. (2015). *Fundamentals of business intelligence*. Berlin; Heidelberg: Springer.
6. Merritt-Holmes, M. (2016). 10 differences between data science and business intelligence. Retrieved 15. 10. 2019 from <https://www.itproportal.com/2016/08/18/10-differences-between-data-science-and-business-intelligence/>.
7. Provost, F. and Fawcett, T. (2013). *Data science for business: what you need to know about data mining and data-analytic thinking* (1 st ed.). Sebastopol: O'Reilly.
8. Provost, F. and Fawcett, T. (2013). Data science and its relationship to big data and data-driven decision making. *Big data*, 1(1), 51–59.
9. Van der Aalst, W. (2016). Data science in action. In W. van der Aalst, *Process mining* (pp. 3–23). Berlin; Heidelberg: Springer.
10. What is data science? [Blog]. (2019). Retrieved 12. 10. 2019 from <https://intellipaat.com/blog/what-is-data-science/>.
11. Zahavi, J. (1999). Mining data for nuggets of knowledge. Retrieved 7. 10. 2019 from

<https://knowledge.wharton.upenn.edu/article/mining-data-for-nuggets-of-knowledge/>.