

A Survey on Big Data

¹NIHARIKA SAHU,

Gandhi Institute of Excellent Technocrats, Bhubaneswar, India

²SUBHASHMITA BEHERA,

Black Diamond College of Engineering & Technology, Jharsuguda, Odisha, India

ABSTRACT –Big data is a term used to describe data or data sets that are so massive or complicated that distributed databases are required instead of standard data processing software. Big data has been the foundation of companies like Google, eBay, LinkedIn, and Facebook from the start. It consists of a collection of enormous and intricate data sets, including vast amounts of data, social media analytics, data management tools, real-time data, etc. The design of sensors, data collection, data duration, sharing, storage, analysis, visualization, and information privacy are among the difficulties. Big data refers to datasets that are rapidly growing and have a large degree of variability, making them challenging to manage with conventional tools and methods. Big data analytics is the study of enormous amounts of data to uncover hidden correlations. Big Data is a type of data that is so complex that managing it and gleaning value and untapped knowledge from it calls for new management strategies, algorithms, and analytics. In order to structure Big Data and address the issue of making it relevant for analytics, a different platform, called Hadoop, is required.

Key Words: Big Data, Parallel programming, MapReduce technique.

I. INTRODUCTION

Every digital process and social media exchange produces Big data. The Systems, sensors and mobile devices transmit. The arrival of big data is from multiple sources at a frightening velocity, volume and variety. We need optimal processing power, analytics capabilities and skills to extract meaningful value from big data. More confident decision making can do with accurate big data. Good decisions lead to greater operational efficiency, cost reduction and reduced risk. Analysis of data sets can find new correlations, to spot business trends, prevent

diseases, and combat crime and soon [1]. Scientists, business executives, practitioners of media, and advertising and governments alike regularly meet difficulties with large data sets in areas including Internet search, finance and business informatics. Data sets grow in size in part because they are increasingly being gathered by cheap and numerous information-sensing mobile devices, aerial (remote sensing), software logs, cameras, microphones, radio-frequency identification (RFID) readers, and wireless sensor networks [2][3][4]. Big data "size" is a constantly moving target, ranging from a few dozen terabytes to many petabyte of data. (1 petabyte is 1000 terabytes)



Fig-1: An image of Big data

Here are some real-world examples of Big Data in action:

- Consumer product companies and retail organizations are monitoring social media like Facebook and Twitter to get an unprecedented view into customer behavior, preferences, and product perception.
- Manufacturers are able to monitor minute vibration data from their equipment, which changes slightly as it wears down, to predict the optimal time to replace or maintain. Replacing it too soon wastes money and replacing it too late triggers an expensive work stoppage.
- Manufacturers are also monitoring social networks, but with a different goal than marketers: They are using it to detect aftermarket support issues before a warranty failure becomes publicly detrimental.
- The government is making data public at the national level, state level, and city level for users to develop new applications that can generate public better.
- Financial Services organizations are taking data mined from customer interaction to slice and dice their users into finely tuned segments and

enable these financial institutions to create increasingly relevant and sophisticated offers.

- Advertising and marketing agencies are tracking social media to see which home insurance applications can be immediately processed, and which ones need a validating in-person visit.
- Retail organizations are engaging brands and advocates, changing the perception of brand antagonists, and even enabling enthusiastic customers to sell their products. All these things are doing by embracing social media.
- Hospitals predict those patients that are likely to seek readmission within a few months of discharge by analyzing medical data and patient records. The hospital can then prevent another costly hospital stay.
- To offer more appealing recommendations and more successful coupon programs, the Web-based businesses are developing information products that combined data gathered from customers
- Sport teams are using data for tracking tickets sales and are using big data for tracking team strategies also.

a. Three Vs of big data: volume, velocity and variety⁵. (Big Data Parameters)



Fig-2: An image of Big data

Volume. Volume of data stored in enterprise repositories have grown from gigabytes to petabytes. Many factors contribute to the increase in data volume like transaction-based data stored through the years, unstructured data streaming in from social media etc. Huge amounts of sensor and machine-to-machine data being collected. In the past days, excessive data volume was a storage issue. But with decreasing storage costs, other issues emerge, including how to determine relevance within large data

volumes and how to use analytics to create value from relevant data. Volume referred as amount of data.

Velocity. Data is streaming in at extraordinary

speed and must be dealt with in a timely manner. RFID sensors and smart metering are driving the need to deal with fast-moving of data in near-real time. It is a challenge for most organizations to react quickly enough to deal with data velocity. Velocity referred the speed of data processing. For time-sensitive processes such as catching fraud, big data must be used. It streams into your enterprise in order to maximize its value.

Variety. Today data comes in different types of formats. Structured and numeric data in traditional data bases. Information created from line-of-

business applications. Unstructured text documents, e-mail, video, audio and financial transactions. Managing, merging and governing different varieties of data are something many organizations still struggle with. Different types and sources of data are there. Data variety exploded from structured and legacy data stored in enterprise storage to unstructured, semistructured, audio, video etc.

We consider two additional dimensions when thinking about big data:

Variability. With the increasing velocities and varieties of data, data flows can be highly inconsistent with periodic peaks. It is trending in social media. Everyday seasonal and event-triggered peak data loads cannot be able to manage. Even more unstructured data involved. The inconsistency the data can show at times—which can hamper the process of handling and managing the data properly. The inconsistency the data can show at times can hamper the process of handling and managing the data properly.

Complexity. Today's data comes from different types of sources. It is still an undertaking to link, match and transform data across systems. Anyway, it is necessary to connect and correlate relationships, hierarchies and multiple data linkages. Otherwise your data can quickly spiral out of control. When large volumes of data come from multiple sources, the data management is very complex. Especially data must be linked, connected, and correlated the users can grasp the information or messes the data is supposed to convey.

Veracity The quality of captured data, which varies so high. The accuracy analysis of data depends on the veracity of source data.

II. PARALLEL PROGRAMMING & MAP REDUCE

Data analysis software parallelizes fairly naturally. Many programmers are interested to building programs on the parallel model. The parallel research had the most success in the field of parallel databases. Rather than requiring the programmer to unknot an algorithm into separate threads to be run on separate cores, parallel databases let them break up the input data tables into pieces, and pump each piece through the same

single-machine program on each processor. This “parallel dataflow” model makes parallel programming as easy as programming a single machine. And it works on “shared-nothing” clusters of computers in a data center: The machines involved can communicate via simple streams of data messages, without a need for an expensive shared RAM or disk infrastructure. [6]

Famous big data analysis tool is Hadoop. Apache Hadoop is an open-source software framework. It is written in Java for distributed storage and distributed processing of big data on computer clusters built from commodity hardware. All the modules in Hadoop are designed with a fundamental assumption that hardware failures (of individual machines or racks of machines) are common place and thus should be automatically handled in software by the framework. [7]

The heart of Hadoop is **MapReduce**. It is this programming paradigm that allows for massive scalability across thousands of servers in a Hadoop cluster. It is useful for batch processing on petabytes or zeta bytes of data stored in Apache Hadoop. If we are familiar with clustered scale-out data processing solutions. Then the MapReduce concept is simple to understand. MapReduce programming model has twisted a new page in the parallelism story. The MapReduce framework is a parallel data flow system that works by dividing data across machines. Each of which runs the same single-node logic. MapReduce asks programmers to write traditional code, in languages like C, Java, Python and Perl. In addition to its familiar syntax, MapReduce allows programs to be written to and read from traditional files in a file system, rather than requiring database schema definitions.

MapReduce refers to two separate and distinct tasks. The first is the job of map, which takes a set of data and converts it into another set of data. Individual elements are broken down into value pairs. The reduce job takes the output from a map as input and combines those data values into a smaller set of values. The reduce job is always performed after the map job. So the sequence of the name MapReduce.

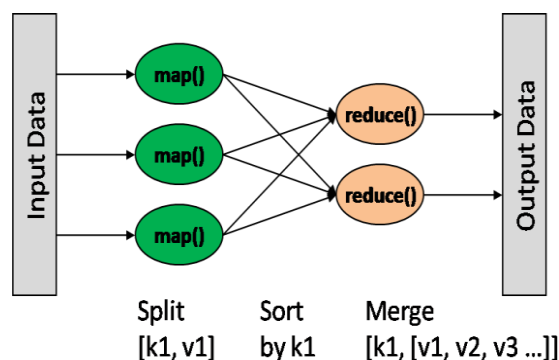


Fig-2:MapReduce

III. BEST BIG DATA ANALYTICS USE CASES

Sentiment Analysis

Sentiment analysis offers powerful business intelligence to enhance the customer experience, revitalize a brand, and gain competitive advantage. The key to successful sentiment analysis lies in the ability to dig for multi-structured data pulled from different sources into a single database.

360-Degree View of Customer

A 360-degree customer view offers a deeper understanding of customer behavior and motivations. Obtaining a 360-degree customer review requires analysis of data from different sources like social media, data collecting sensors, mobile devices etc. From there, more effective micro-segmentation and real-time marketing are getting a result.

Ad Hoc Data Analysis

Ad-hoc analysis only looks at the data requested or needed, providing another layer of analysis for data sets that are becoming larger and more varied. Big data ad-hoc analytics can help in the effort to gain greater insight into customers by analyzing the relevant data from unstructured sources, both external and internal.

Real-Time Analytics

Systems that offer real-time analytics quickly decipher and analyze data sets, providing results even as data is being generated and collected. This high-velocity method of analytics can lead to immediate reaction and changes. It allows for better sentiment analysis, split testing, and improved targeted marketing.

Multi-Channel Marketing

Multi-channel marketing creates a seamless across different types of media like company websites, social media, and physical stores. During all stages of the buying process multi-channel marketing requires an integrated big data approach.

Customer Micro-Segmentation

Customer micro-segmentation provides more tailored and targeted messaging for smaller groups. This personalized approach requires analysis of big data collected through sources like customers' online interactions, social media etc.

Ad Fraud Detection

Ad fraud detection requires data analysis of fraud strategies by recognizing patterns and behaviors. Data that shows irregularity of group behavior make it so ad fraud is found out and blocked before it is spread.

Clickstream Analysis

Click stream analysis helps to grow the user experience by optimizing company websites, and offering better insight into customer segments. Click stream analysis helps to personalize the buying experience, getting an improved return on customer visits with big data.

Data Warehouse Modernization

Integrate big data and data warehouse capabilities to boost operational efficiency. Optimize your data warehouse to enable fresh types of analysis. Use big data technologies to set up a staging area or landing zone for your new data before reformatting what data should be moved to the data warehouse. Invest in frequent access to stored data from warehouse and application database using in-sequence integration software and tools.

Big Data and Predictive Modeling

The most common uses of big data by companies are for tracking business processes and outcomes, and for building a wide array of predictive models. Amazon and Netflix recommendations rely on predictive models of what book or movie an individual might want to purchase. Google's search results and news feed rely on algorithms that predict the significance of particular webpages or articles. Apple's auto-complete function tries to forecast the rest of one's text or email based on past convention patterns. Online advertising and marketing rely greatly on automated predictive models that aim to identify individuals who might be particularly likely to answer to offers.

The application of predictive algorithms extends well beyond the online world. In healthcare, it is now common for insurers to adjust payments and quality measures based on "risk scores," which are resulting from predictive models of human being health expenses and outcomes. An individual's risk score is naturally a weighted sum of health indicators that recognize whether an individual has different persistent conditions, with the weights chosen based on a statistical analysis. Credit card companies use predictive models of default and repayment to guide their underwriting, pricing, and marketing actions.

IV. BIG DATA CHALLENGES

Heterogeneity, scale, timeliness, complexity, and privacy problems with Big Data impede progress at all phases of the pipeline that can create value from data. The problems start right away during data acquisition, when the data tsunami requires us to make decisions, currently in an ad hoc manner, about what data to maintain and what to reject, and how to store what we keep unfailingly with their metadata. A great deal of data today is not natively in structured format; for example, tweets and blogs are weakly ordered pieces of text, while images and video are structured for storage and display. But not for semantic content and look for. Transforming such content into a structured format for later analysis is a main test. The value of data explodes when it can be associated with other data. Thus data integration is a major creator of value. The majority of data is directly generated in digital format today; we have the opportunity and the challenge both to influence the creation to facilitate later linkage and to automatically link before created data. Data analysis, organization, recovery, and modeling are the foundational challenges. Data analysis is a clear

bottleneck in a lot of applications, both due to the small loss of capability of the original algorithms and due to the complexity of the data that needs to be analyzed. Lastly, the presentation of the results and its clarification by non-technical domain experts is vital to extracting actionable knowledge.

Volume of data

The volume of data, especially machine-generated data, is exploding, how fast that data is growing every year, with new sources of data that are emerging. For instance, in the year 2000, 800,000 petabytes (PB) of data were stored in the world. According to IBM it is anticipated to reach 35 zettabytes (ZB) by 2020. Social media plays a key role. Twitter generates 7+ terabytes (TB) of data every day. Facebook, 10 TB. Mobile devices also play an important role.

Big data skills are in short supply

There's already a shortage of data scientists in the market. This includes a shortage of people who know how to labor well with large volumes of data and big data sets. Companies need the right merge of people to help make sense of the data streams that are coming into their organizations. This includes skills for applying prophetic analytics to big data, a skill set that even most data scientists are short of.

V. CONCLUSIONS

The availability of Big Data, low-cost commodity hardware, and analytic software has had a unique moment in the history of data analysis. The union of these trends means that we have the capabilities required to analyze amazing data sets quickly and cost-effectively for the first time in history. All these capabilities are neither theoretical nor trivial. They represent a real leap forward and a clear chance to realize enormous gains in terms of efficiency, productivity, income, and profitability.

Requirements for dealing out that may seem unbelievable today will soon be routine when big data systems are available. We learn how to exploit them. Not very many years ago, systems the scale of Facebook and Google would have seemed like science fiction. At that time 100 transactions per second for airline and banking systems was a stretch. Several new requirements will all combine data from many sources, not all of which will become a

owned. For instance, some will make use of 'open data' from government. Lots of opening for innovators!

REFERENCES

- [1] "Data, data everywhere". *The Economist*. 25 February 2010. Retrieved 9 December 2012.
- [2] "Data, data everywhere". *The Economist*. 25 February 2010. Retrieved 9 December 2012.
- [3] "Community cleverness required". *Nature* **455** (7209): 1. 4 September 2008. doi:10.1038/455001a.
- [4] "Sandia sees data management challenge spiral". *HPC Projects*. 4 August 2009.
- [5] META Group. "3D Data Management: Controlling Data Volume, Velocity, and Variety." February 2001.
- [6] Hellerstein, Joe (9 November 2008). "Parallel Programming in the Age of Big Data". *Gigaom Blog*.
- [7] "Welcome to Apache™ Hadoop®!". *hadoop.apache.org*. Retrieved 2015-09-20.