



# Introduction to Big Data - I



# Introduction to Big Data - I





# Introduction to Big Data - I

Due to the advent of **new technologies, devices, and communication means** like social networking sites, the amount of data produced by mankind is growing rapidly every year. The amount of data produced by us from the beginning of time till 2003 was 5 billion gigabytes. This rate is still growing enormously.

## What is Big Data?

**Big data** is a collection of large datasets that cannot be processed using traditional computing techniques. It is not a single technique or a tool, rather it has become a complete subject, which involves various **tools, techniques and frameworks**.



# Introduction to Big Data - I

## What Comes Under Big Data?

Big data involves the data produced by different **devices** and **applications**. Given below are some of the fields that come under the umbrella of **Big Data**.

- **Black Box Data** – It is a component of helicopter, airplanes, and jets, etc. It captures voices of the flight crew, recordings of microphones and earphones, and the performance information of the aircraft.
- **Social Media Data** – Social media such as Facebook and Twitter hold information and the views posted by millions of people across the globe.
- **Stock Exchange Data** – The stock exchange data holds information about the ‘buy’ and ‘sell’ decisions made on a share of different companies made by the customers.
- **Power Grid Data** – The power grid data holds information consumed by a particular node with respect to a base station.
- **Transport Data** – Transport data includes model, capacity, distance and availability of a vehicle.
- **Search Engine Data** – Search engines retrieve lots of data from different databases.

# Introduction to Big Data - I

Thus Big Data includes **huge volume**, **high velocity**, and extensible **variety** of data. The data in it will be of **three types**.

- **Structured data** – Relational data.
- **Semi Structured data** – XML data.
- **Unstructured data** – Word, PDF, Text, Media Logs.

## Benefits of Big Data

- Using the information kept in the social network like **Facebook**, the marketing agencies are learning about the response for their **campaigns**, **promotions**, and other advertising mediums.
- Using the information in the social media like preferences and product perception of their consumers, product companies and retail organizations are planning their production.
- Using the data regarding the previous medical history of patients, hospitals are providing better and quick service.



# Introduction to Big Data - I

## Big Data Technologies

Big data technologies are important in providing more accurate analysis, which may lead to more concrete decision-making resulting in greater **operational efficiencies**, **cost reductions**, and **reduced risks** for the **business**.

To harness the power of big data, you would require an infrastructure that can manage and process huge volumes of structured and unstructured data in realtime and can protect data privacy and security.

There are various technologies in the market from different vendors including **Amazon**, **IBM**, **Microsoft**, etc., to handle big data. While looking into the technologies that handle big data, we examine the following **two classes** of technology –

# Introduction to Big Data - I

## Operational Big Data

This include systems like **MongoDB** that provide operational capabilities for real-time, interactive workloads where data is primarily captured and stored.

**NoSQL** Big Data systems are designed to take advantage of new cloud computing architectures that have emerged over the past decade to allow massive computations to be run inexpensively and efficiently. This makes operational big data workloads much easier to **manage, cheaper,** and **faster** to implement.

Some NoSQL systems can provide insights into patterns and trends based on real-time data with minimal coding and without the need for data scientists and additional infrastructure.



# Introduction to Big Data - I

## Analytical Big Data

These includes systems like **Massively Parallel Processing** (MPP) database systems and **MapReduce** that provide analytical capabilities for retrospective and complex analysis that may touch most or all of the data.

MapReduce provides a new method of analyzing data that is complementary to the capabilities provided by SQL, and a system based on MapReduce that can be scaled up from single servers to thousands of high and low end machines.





# Introduction to Big Data - I

These two classes of technology are complementary and frequently deployed together.

	Operational	Analytical
Access Pattern	Writes and Reads	Reads
Queries	Selective	Unselective
Data Scope	Operational	Retrospective
End User	Customer	Data Scientist
Technology	NoSQL	MapReduce, MPP Database

# Introduction to Big Data - I

## Big Data Challenges

The major challenges associated with big data are as follows –

- Capturing data
- Curation
- Storage
- Searching
- Sharing
- Transfer
- Analysis
- Presentation

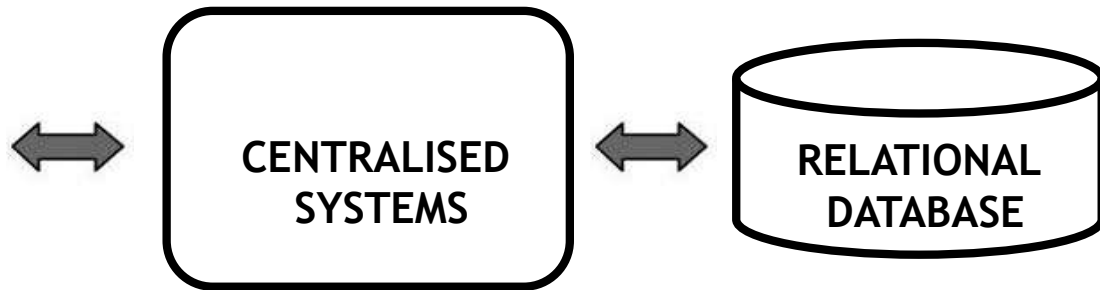
# Introduction to Big Data - I

## Traditional Approach

In this approach, an enterprise will have a computer to store and process big data. For storage purpose, the programmers will take the help of their choice of database vendors such as **Oracle**, **IBM**, etc. In this approach, the **user** interacts with the application, which in turn handles the part of **data storage** and **analysis**.



USER



## Limitation

This approach works fine with those applications that process **less voluminous data** that can be accommodated by standard database servers, or up to the limit of the processor that is processing the data. But when it comes to dealing with huge amounts of scalable data, it is a hectic task to process such data through a single database bottleneck.

# Introduction to Big Data - I

## Traditional Data Mining Life Cycle

In order to provide a framework to organize the work needed by an organization and deliver clear insights from Big Data, it's useful to think of it as a cycle with different stages. It is by no means linear, meaning all the stages are related with each other. This cycle has superficial similarities with the more traditional data mining cycle as described in **CRISP methodology**.

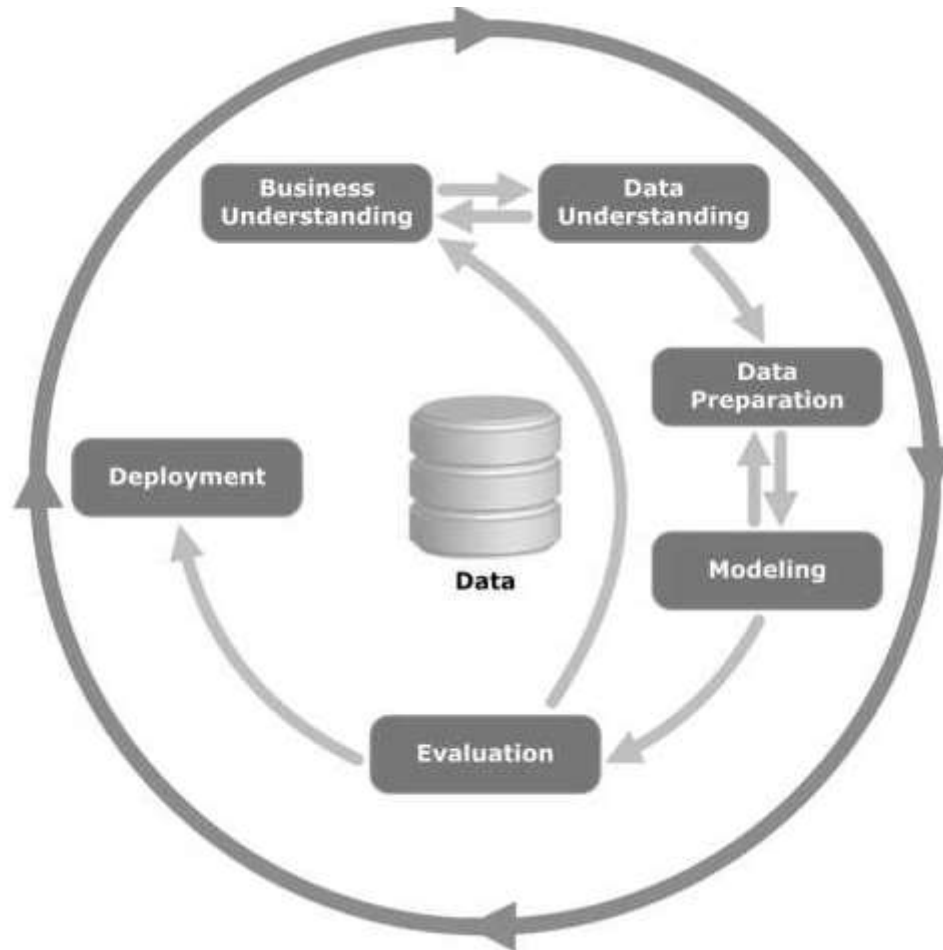
# Introduction to Big Data - I

## CRISP-DM Methodology

The **CRISP-DM methodology** that stands for **Cross Industry Standard Process for Data Mining**, is a cycle that describes commonly used approaches that data mining experts use to tackle problems in traditional BI data mining. It is still being used in traditional BI data mining teams.

Take a look at the following illustration. It shows the major stages of the cycle as described by the CRISP-DM methodology and how they are interrelated.

# Introduction to Big Data - I



# Introduction to Big Data - I

Let us now learn a little more on each of the stages involved in the CRISP-DM life cycle.

- **Business Understanding** – This initial phase focuses on understanding the project objectives and requirements from a business perspective, and then converting this knowledge into a data mining problem definition. A preliminary plan is designed to achieve the objectives. A **decision model**, especially one built using the **Decision Model** and **Notation standard** can be used.
- **Data Understanding** – The data understanding phase starts with an initial data collection and proceeds with activities in order to get familiar with the data, to identify data quality problems, to discover first insights into the data, or to detect interesting subsets to form hypotheses for hidden information.
- **Data Preparation** – The data preparation phase covers all activities to construct the final dataset (data that will be fed into the modeling tool(s)) from the initial raw data. Data preparation tasks are likely to be performed multiple times, and not in any prescribed order. Tasks include table, record, and attribute selection as well as transformation and cleaning of data for modeling tools.

# Introduction to Big Data - I

- **Modeling** – In this phase, various modeling techniques are selected and applied and their parameters are calibrated to optimal values. Typically, there are several techniques for the same data mining problem type. Some techniques have specific requirements on the form of data. Therefore, it is often required to step back to the data preparation phase.
- **Evaluation** – At this stage in the project, you have built a model (or models) that appears to have high quality, from a data analysis perspective. Before proceeding to final deployment of the model, it is important to evaluate the model thoroughly and review the steps executed to construct the model, to be certain it properly achieves the business objectives.

A key objective is to determine if there is some important business issue that has not been sufficiently considered. At the end of this phase, a decision on the use of the data mining results should be reached.

- **Deployment** – Creation of the model is generally not the end of the project. Even if the purpose of the model is to increase knowledge of the data, the knowledge gained will need to be organized and presented in a way that is useful to the customer.



# Introduction to Big Data - I

## SEMMA Methodology

**SEMMA** is another methodology developed by SAS for data mining modeling. It stands for **S**ample, **E**xplore, **M**odify, **M**odel, and **A**ssess. Here is a brief description of its stages –

- **Sample** – The process starts with data sampling, e.g., selecting the dataset for modeling. The dataset should be large enough to contain sufficient information to retrieve, yet small enough to be used efficiently. This phase also deals with data partitioning.
- **Explore** – This phase covers the understanding of the data by discovering anticipated and unanticipated relationships between the variables, and also abnormalities, with the help of data visualization.
- **Modify** – The Modify phase contains methods to select, create and transform variables in preparation for data modeling.
- **Model** – In the Model phase, the focus is on applying various modeling (data mining) techniques on the prepared variables in order to create models that possibly provide the desired outcome.
- **Assess** – The evaluation of the modeling results shows the reliability and usefulness of the created models.

# Introduction to Big Data - I

The main difference between **CRISM-DM** and **SEMMA** is that SEMMA focuses on the **modeling aspect**, whereas CRISP-DM gives more importance to **stages of the cycle prior to modeling** such as understanding the business problem to be solved, understanding and preprocessing the data to be used as input, for example, machine learning algorithms.

# Introduction to Big Data - I

## Big Data Life Cycle

In today's big data context, the previous approaches are either incomplete or suboptimal. For example, the SEMMA methodology disregards completely data collection and preprocessing of different data sources. These stages normally constitute most of the work in a successful big data project.

A big data analytics cycle can be described by the following stage -

- **Business Problem Definition**
- Research
- **Human Resources Assessment**
- Data Acquisition
- **Data Munging**
- Data Storage
- **Exploratory Data Analysis**
- Data Preparation for Modeling and Assessment
- **Modeling**
- Implementation

# Introduction to Big Data - I

Now, we will discuss on each of these stages of big data life cycle.

- **Business Problem Definition**

This is a point common in traditional BI and big data analytics life cycle. Normally it is a non-trivial stage of a big data project to define the problem and evaluate correctly how much potential gain it may have for an organization. It seems obvious to mention this, but it has to be evaluated what are the expected gains and costs of the project.

- **Research**

Analyze what other companies have done in the same situation. This involves looking for solutions that are reasonable for your company, even though it involves adapting other solutions to the resources and requirements that your company has. In this stage, a methodology for the future stages should be defined.

- **Human Resources Assessment**

Once the problem is defined, it's reasonable to continue analyzing if the current staff is able to complete the project successfully. Traditional BI teams might not be capable to deliver an optimal solution to all the stages, so it should be considered before starting the project if there is a need to outsource a part of the project or hire more people.

# Introduction to Big Data - I

- Data Acquisition

This section is **key** in a big data **life cycle**; it defines which type of profiles would be needed to deliver the resultant data product. **Data gathering** is a non-trivial step of the process; it normally involves gathering unstructured data from different sources. To give an example, it could involve writing a **crawler** to retrieve reviews from a **website**. This involves dealing with text, perhaps in different languages normally requiring a significant amount of time to be completed.



# Introduction to Big Data - I

- **Data Munging**

Once the data is retrieved, for example, from the web, it needs to be stored in an easy-to-use format. To continue with the reviews examples, let's assume the data is retrieved from different sites where each has a different display of the data.

Suppose one data source gives reviews in terms of rating in stars, therefore it is possible to read this as a mapping for the response variable  $y \in \{1, 2, 3, 4, 5\}$ . Another data source gives reviews using two arrows system, one for **up voting** and the other for **down voting**. This would imply a response variable of the form  $y \in \{\text{positive}, \text{negative}\}$ .

In order to combine both the data sources, a decision has to be made in order to make these two response representations equivalent. This can involve converting the first data source response representation to the second form, considering one star as negative and five stars as positive. This process often requires a large time allocation to be delivered with good quality.



# Introduction to Big Data - I

- **Data Storage**

Once the data is processed, it sometimes needs to be stored in a database. Big data technologies offer plenty of alternatives regarding this point. The most common alternative is using the **Hadoop File System** for storage that provides users a limited version of SQL, known as **HIVE Query Language**. This allows most analytics task to be done in similar ways as would be done in traditional BI data warehouses, from the user perspective. Other storage options to be considered are **MongoDB, Redis, and SPARK**.

This stage of the cycle is related to the human resources knowledge in terms of their abilities to implement different architectures. Modified versions of traditional data warehouses are still being used in large scale applications. For example, **teradata** and **IBM** offer SQL databases that can handle **terabytes of data**; open source solutions such as **postgreSQL** and **MySQL** are still being used for large scale applications.

# Introduction to Big Data - I

- **Data Storage (Cont'd)**

Even though there are differences in how the different storages work in the background, from the client side, most solutions provide a SQL API. Hence having a good understanding of SQL is still a key skill to have for big data analytics.

This stage a priori seems to be the most important topic, in practice, this is not true. It is not even an essential stage. It is possible to implement a big data solution that would be working with real-time data, so in this case, we only need to gather data to develop the model and then implement it in real time. So there would not be a need to formally store the data at all.



# Introduction to Big Data - I

- **Exploratory Data Analysis**

Once the data has been cleaned and stored in a way that insights can be retrieved from it, the data exploration phase is mandatory. The objective of this stage is to understand the data, this is normally done with statistical techniques and also plotting the data. This is a good stage to evaluate whether the problem definition makes sense or is feasible.

- **Data Preparation for Modeling and Assessment**

This stage involves reshaping the cleaned data retrieved previously and using statistical preprocessing for missing values imputation, outlier detection, normalization, feature extraction and feature selection.

- **Modelling**

The prior stage should have produced several datasets for training and testing, for example, a predictive model. This stage involves trying different models and looking forward to solving the business problem at hand. In practice, it is normally desired that the model would give some insight into the business. Finally, the best model or combination of models is selected evaluating its performance on a left-out dataset.

# Introduction to Big Data - I

- **Implementation**

In this stage, the data product developed is implemented in the data pipeline of the company. This involves setting up a validation scheme while the data product is working, in order to track its **performance**. For example, in the case of implementing a predictive model, this stage would involve applying the model to new data and once the response is available, evaluate the model.

# Introduction to Big Data - I

**THANK YOU**