



Module 1

Syllabus

Module- I

Introduction to Big Data Platform–History of Data Management-Structuring Bigdata-Elements of Big Data, Big data Analytics

Introduction to Data Mining: Data Mining Process, Data Mining Techniques - Clustering - Classification - Regression - Outlier detection - Association rules - Prediction, Big Data Vs Data Mining

Contents

- Introduction to Big Data Platform
 - History of Data Management
 - Structuring Big data
 - Elements of Big Data
 - Big data stack
 - Big data Analytics



DATA NEVER SLEEPS 2.0

How Much Data is Generated Every Minute?

Data is being created every minute of every day without us even noticing it. Given how much information is floating around these days, it's tempting to talk about big data only in terms of size. But data becomes the most valuable asset in digital activity by creating enough context and stories. So let's take a look at the things we do online every minute to make sure they're valuable. With every status we share, every article we read or every photo we upload, we're creating a digital trail that tells a story. Below, we explore how much data is generated in one minute.

PINTEREST
USERS PIN

3,472 images

VINE
USERS

SHARE
8,333
VIDEOS.

SKYPE
USERS
CONNECT FOR
23,300 HOURS.

YELP
USERS POST
26,380
REVIEWS.

APPLE USERS
DOWNLOAD
48,000
apps.

PANDORA
USERS LISTEN TO
61,141
HOURS OF
music.

YOUTUBE
USERS UPLOAD
72 HRS.
OF NEW
VIDEO.

EMAIL
USERS SEND
204,000,000
MESSAGES.

Google
RECEIVES OVER
4,000,000
SEARCH
QUERIES.

FACEBOOK
USERS SHARE
2,460,000
PIECES OF CONTENT.

TINDER
USERS SWIPE
416,667
TIMES.

WHATSAPP
USERS SHARE
347,222
PHOTOS.

EVERY
MINUTE
OF THE
DAY

TWITTER USERS
TWEET
277,000
TIMES.

INSTAGRAM
USERS
POST
216,000
NEW PHOTOS.

AMAZON
MAKES
\$83,000
IN ONLINE SALES.



THE GLOBAL INTERNET POPULATION GREW
14.3% FROM 2011-2013 AND NOW REPRESENTS

2.4 BILLION PEOPLE.

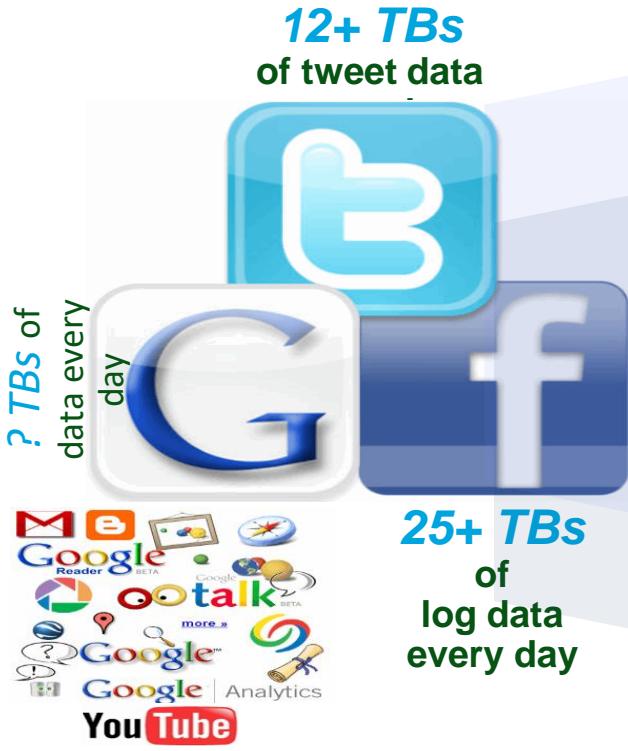
With each click, share and like, the world's data pool is expanding faster than we can comprehend. Businesses today are paying attention to scores of data sources to make crucial decisions about the future. The team at Domo can help your business make sense of this endless stream of data by providing executives with all their critical information in one intuitive platform. Domo delivers the insights you need to transform the way you run your business. Learn more at www.domo.com.

SOURCES:

INTELBLOCS, NYTIMES.COM, INTEL.COM, APPLE.COM, TIME.COM, DAILYMAIL.CO.UK, SKYPE.COM, STATISTICAIRAN.COM



Big Data – A growing torrent !



Introduction to Big Data

What is Big Data?

What makes data, “Big” Data?

Big Data

- Big data basically refers to huge volume of data that cannot be stored and processed using the traditional approach, within the given time frame.
- Organizations use data generated through various sources to run their businesses.

Contd...

- They analyze the data to understand and interpret market trends, study customer behavior, and take financial decisions.
- Big data consists of large datasets that cannot be managed efficiently by the common database management systems.
- These datasets range from terabytes to exabytes.

What is big data?

- Every second, there are around 822 tweets on Twitter.
- Every minute, nearly 510 comments are posted, 293,000 statuses are updated and 136,000 photos are uploaded on Facebook.
- Every hour, Walmart, a global discount departmental store chain, handles more than 1 million customer transactions.

Contd...

- Every day, consumers make around 11.5 million payments by using paypal.
- We live in a digital world where data is increasing rapidly because of the ever increasing use of the internet, sensors and so on...

Who's Generating Big Data

Social Media



Social media and networks
(all of us are generating data)



Scientific instruments
(collecting all sorts of data)



Mobile devices
(tracking all objects all the time)



Sensor technology and networks
(measuring all kinds of data)

The Model Has Changed...

- **The Model of Generating/Consuming Data has Changed**

Old Model: Few companies are generating data, all others are consuming data



New Model: all of us are generating data, and all of us are consuming data



- 90% of the data in the world today has been created in the last two years alone.
- This data comes from everywhere: sensors used to gather climate information, posts to social media sites, digital pictures and videos, purchase transaction records and cell phone GPS signals and so on... This data is big data.
- Data is everywhere, in every industry, in the form of numbers, images, videos and text.
- As data continues to grow, so does the need to organize it.

- The process of capturing or collecting Big data is known as ‘datafication’.
- Big data is ‘datafied’ so that it can be used productively.
- To extract meaningful value from Big data, we need optimal processing power, analytical capabilities and skills.

Definition

“***Big Data***” is data whose scale,
diversity, and complexity require new
architecture, techniques, algorithms, and
analytics to manage it and
extract value and hidden knowledge from it...

BIG DATA?

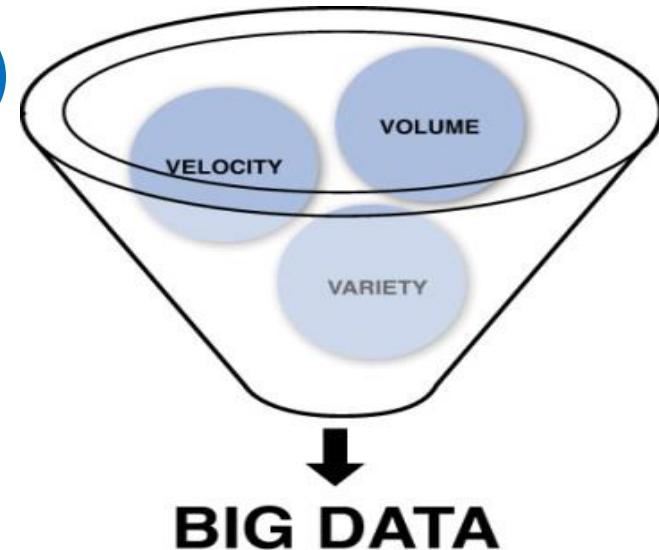
VOLUME
Large amounts of data.

VELOCITY
Needs to be analyzed quickly.

VARIETY
Different types of structured and unstructured data.

*What makes data,
"Big" Data?*

- In 2012, [Gartner](#) updated its definition as follows: "Big Data are **high-volume**, **high-velocity**, and/or **high-variety** information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization.
- **The 3 Vs (Volume, Velocity and Variety) makes the data as "Big" Data.**
- Big data is **structured, unstructured or semi-structured** in nature.



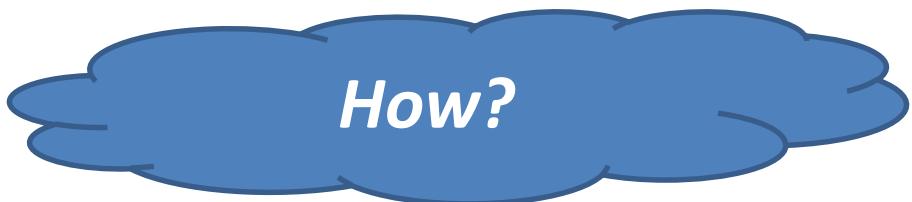
Features of Big data

Big Data

Is a new data challenge that requires leveraging existing systems differently

Is classified in terms of 4vs
Volume
Variety
Velocity
veracity

Is usually unstructured and qualitative in nature



How?

- It becomes difficult for computing systems to manage 'Big data' because of the immense speed and volume at which it is generated. Traditional data management, warehousing systems are not capable to analyze this type of data.
- Due to its complexity, big data is stored in distributed architecture file system.
- **Hadoop** by Apache is widely used for storing and managing big data; uses HDFS as the storage system.

Types and sources of data!

Big data is a pool of huge amounts of data of all types, shapes and formats collected from varied sources. Some common types of data and their sources:

| Type | Description | Source |
|--------------------|---|---|
| Social Data |  Information collected from various social networking sites and online portals | Facebook, Twitter and LinkedIn |
| Machine Data |  Information generated from RFID chips, bar code scanners, and sensors | RFID chip readings, Global Positioning System (GPS) results |
| Transactional Data | Information generated from online shopping sites, retailers, and B2B transactions | Retail websites like eBay and Amazon |

History of Data Management – Evolution of Big Data

- Challenges faced while handling data over the past few decades:

In the early 60s, technology witnessed problems with **velocity** or real-time data assimilation.

The need inspired the evolution of databases.



In the 90s, technology witnessed issues with **variety** (e-mails, documents, videos) leading to the emergence of non-SQL stores.



Today, technology is facing issues related to huge **volume**, leading to new storage and processing solutions.

Major milestones in the evolution of big data:

- Some of the major milestones in the evolution of big data:
- **1940s** : An American librarian speculated the potential shortfall of shelves and cataloging staff, realizing the **rapid increase in information and limited storage**.

- **1960s** : Automatic Data Compression was published in the communications of the ACM (Association of Computing Machinery).
 - The paper described '**Automatic Data Compression**' as a complete automatic and fast three-part compressor that can be used for any kind of information in order to reduce the slow external storage requirements and increase the rate of transmission from a computer system.

- **1970s** : In Japan, the Ministry of Posts and Telecommunications initiated a project to study information flow in order to track the volume of information circulating in the country.
- **1980s** : A research project was started by the Hungarian Central Statistics Office to account for the country's information industry. It measured the volume of information in bits.

- **1990s** : Digital storage systems became more economical than paper storage.
- **2000 onwards** : Various methods were introduced to streamline information.
- Techniques for controlling the volume, velocity and variety of data emerged...

- Structuring big data means, arranging the available data in a manner such that it becomes easy to study, analyze and derive conclusion from it.
- Structuring data helps in understanding user behaviors, requirements and preferences to make personalized recommendations for every individual.
- When a user regularly visits or purchases from online shopping sites, say eBay, each time he/she logs in, the system can present a recommended list of products that may interest the user on the basis of his/her purchases or searches.
- Thus presenting a specially customized recommendation set for every user. This is the power of big data analytics.

Types of data

- Data that comes from multiple sources, such as databases, Enterprise Resource Planning (ERP) systems, weblogs, chat history, and GPS maps, varies in format.
- Different formats of data need to be made consistent and clear to be used for analysis.
- Types of Data Sources:
 - Internal sources, such as organizational or enterprise data
 - External sources, such as social data.

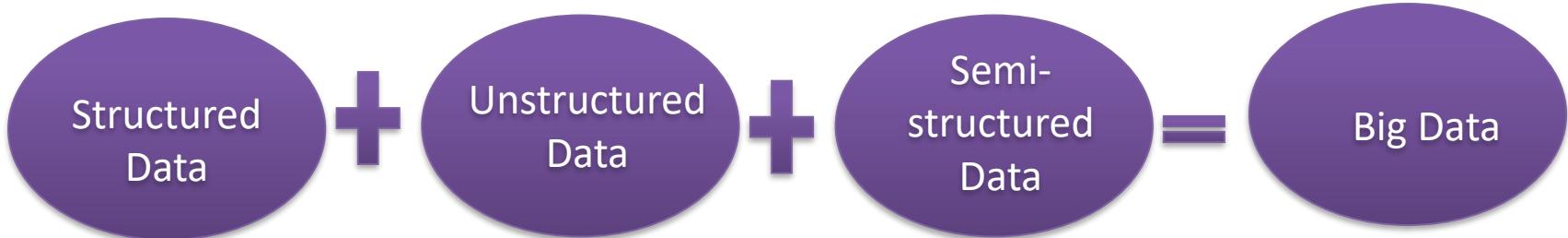
Internal & External Sources of data

| Data source | Definition | Examples of sources | Application |
|-------------|---|--|--|
| Internal | Provides structured or organized data that originates from within the enterprise and helps run business. | Customer Relationship Management (CRM), Enterprise Resource Planning (ERP) systems, customer details, products and sales details | This data is used to support daily business operations of an organization. |
| External | Provides unstructured or unorganized data that originates from the external environment of an organization. | Business partners, Internet, Government, Market research organizations | This data is often analyzed to understand the entities such as customers, competitors, market and environment. |

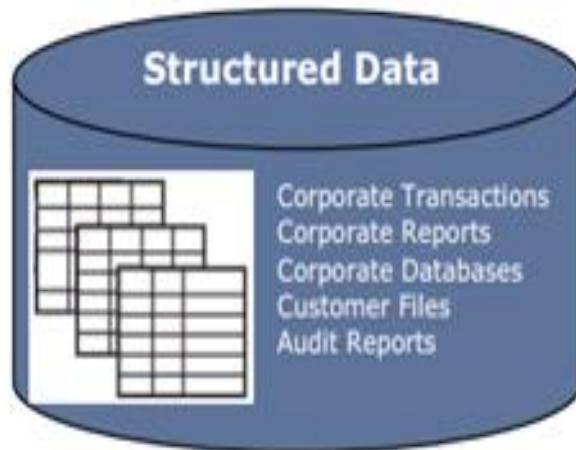
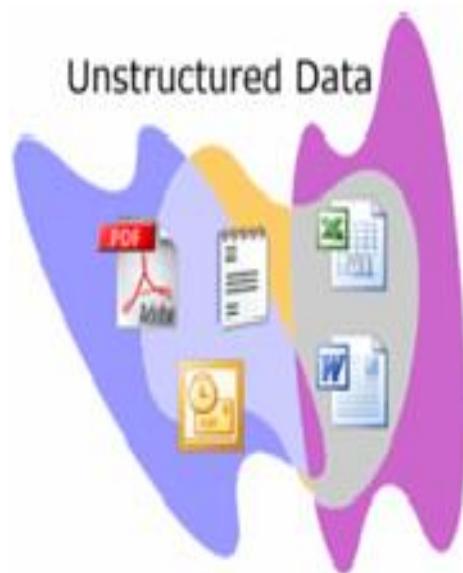
Types of data

- Based on the data received from the sources, the three types of data:
 1. Structured data
 2. Unstructured data
 3. Semi-structured data

In a real-world scenario, the unstructured data is larger in volume than the structured and semi-structured data, approximately 70% to 80% of data is in unstructured form.



Types of data



Semi-structured data

```
<Sensor>
<name>Sensor 193</name>
<attributes>
<Attribute>
<name>Alpha</name>
<x>101</x>
<y>20031</y>
</Attribute>
<Attribute>
<name>Beta</name>
<x>243</x>
<y>3037</y>
</Attribute>
</attributes>
<scale>1</scale>
</Sensor>
```

Structured Data

- Structured Data has a defined pattern. This pattern makes it easier for any program to sort, read and process the data.
- Structured data:
 - Is organized data in a predefined format
 - Is stored in tabular form
 - Is the data that resides in fixed fields within a record or file
 - Is formatted data that has entities and their attributes mapped
 - Is used to query and report.
- Processing structured data is much easier and faster than processing data without any specific repeating patterns.

Structured Data

Sources of structured data:

1. Relational databases
2. Flat files in the form of records (A flat file database is a database that stores data in a plain text file. Each line of the text file holds one record, with fields separated by delimiters, such as commas or tabs)
3. Multidimensional databases used in data warehouse technology.
4. Legacy Databases (It can refer to any of the older database technologies. Flat file, hierarchy, and network databases are usually referred as *legacy* databases.)

Structured Data

A sample of structured data in which the attribute data for every customer is stored:

| Customer ID | Name | Product ID | City | State |
|-------------|--------|------------|-----------|-----------|
| 22465 | Smitha | 141 | Bangalore | Karnataka |
| 45378 | George | 177 | Kochi | Kerala |
| 87231 | John | 445 | Chennai | Hyderabad |

Structured Data



| | | | | |
|-------|-------|-------|-------|-------|
| 0.103 | 0.176 | 0.387 | 0.300 | 0.379 |
| 0.333 | 0.384 | 0.564 | 0.587 | 0.857 |
| 0.421 | 0.309 | 0.654 | 0.729 | 0.228 |
| 0.266 | 0.750 | 1.056 | 0.936 | 0.911 |
| 0.225 | 0.326 | 0.643 | 0.337 | 0.721 |
| 0.187 | 0.586 | 0.529 | 0.340 | 0.829 |
| 0.153 | 0.485 | 0.560 | 0.428 | 0.628 |

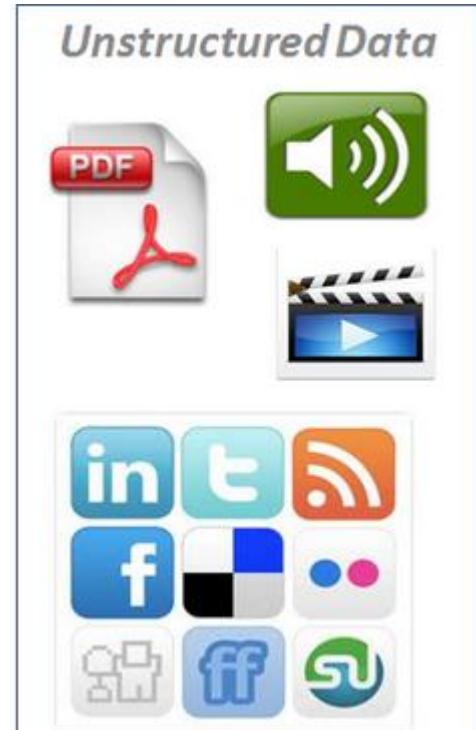
Unstructured Data

- Unstructured data is a set of data that might or might not have any logical or repeating patterns.
- Unstructured Data:
 - Consists of metadata
 - Comprises inconsistent data, such as data obtained from files, social media websites, satellites etc.
 - Consists of data in different formats such as e-mail, text, audio, video or images.

Unstructured Data

Sources of Unstructured Data:

- Internal and external to an organization: documents, logs, survey results, feedback and e-mails from both within and across the organization.
- Social media : data obtained from social networking platforms, including YouTube, Facebook, Twitter, LinkedIn and so on..
- Mobile data : data such as text messages and location information.



Unstructured Data

Challenges associated with unstructured data:

- Identifying the unstructured data that can be processed.
- Sorting, organizing and arranging unstructured data in different sets and formats.
- Combining and linking unstructured data in a more structured format to derive any logical conclusions.
- Costing in terms of storage and human resource needed to deal with exponential growth of unstructured data.

Semi-Structured Data

- Semi-structured data, also known as having a schema-less or self describing structure, refers to a form of structured data that contains tags or markup elements in order to separate elements and generate hierarchies of records and fields in the given data.
- Such type of data does not follow the proper structure of data models as in relational databases.
- In other words, data is stored inconsistently in rows and columns of a database.

Semi-Structured Data

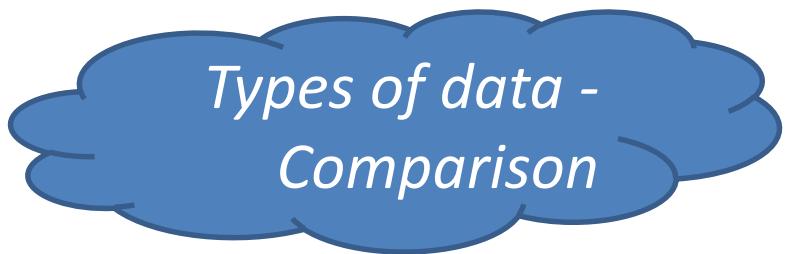
Sources:

- Files systems such as web data in the form of cookies
- Data exchange formats such as JavaScript Object Notation (JSON) data.

Example for semi-structured data

| SL.NO | NAME | E-MAIL |
|-------|--------------------------------------|--|
| 1 | Sam Jacobs | smj@xyz.com |
| 2 | First name: David Last name : Lee | davidlee@xyz.com |

- Example JSON
- { "name":"John", "age":30, "car":null }
- JSON objects are surrounded by curly brackets { }.
- JSON objects are written in key/value pairs.
- Keys must be strings, and values must be a valid JSON data type (string, number, array, boolean or null).
- Keys and values are separated by a colon.
- Each key/value pair is separated by a comma.



Types of data - Comparison

Big Data – Structured vs. Unstructured Data

| Features | Structured Data | Unstructured Data |
|--------------------------|------------------------------------|---|
| Representation | Discrete rows and columns | Less defined boundaries and easily addressable |
| Storage | Rational Databases or Spreadsheets | Unmanaged file structured |
| Metadata | Syntax | Semantics |
| Integration Tools | ETL or ELT | Batch processing or manual data entry that involves codes |
| Standard | SQL, ADO.NET, ODBC,... | OpenXML, JSON, SMTP, SMS, CSV,... |
| Databases | MSSQL, Oracle, Excel,... | Hadoop, HDInsight, MongoDB,... |
| Content | Typically Text | Text, Images, Audio, Video, Documents |

Elements of Big Data

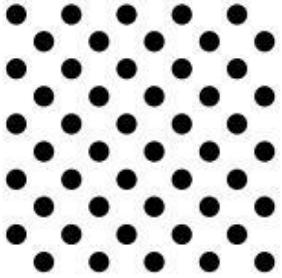
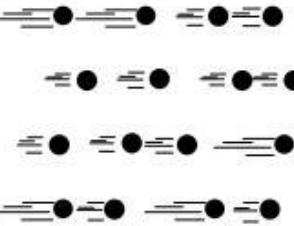
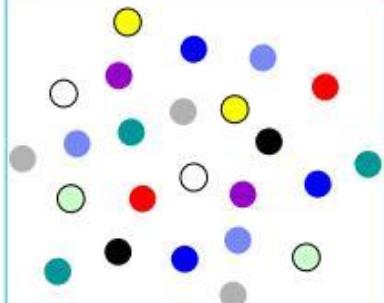
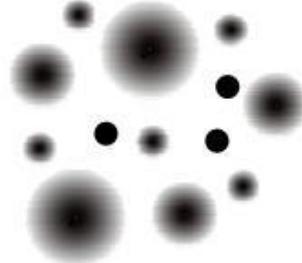
- The growth of data can be represented in terms of four Vs:
- **Volume**
- **Velocity**
- **Variety**
- **Veracity**

Elements of Big Data

4 V's

1. **Volume:** Amount of quantity of data an organization has to deal with.
2. **Velocity:** The speed which the data flows in an organization.
3. **Variety:** Varied forms of data, such as structured, semi-structured, or unstructured.
4. **Veracity:** Certainty of data

4V's

| Volume | Velocity | Variety | Veracity* |
|---|---|--|---|
|  |  |  |  |
| Data at Rest Terabytes to exabytes of existing data to process | Data in Motion Streaming data, milliseconds to seconds to respond | Data in Many Forms Structured, unstructured, text, multimedia | Data in Doubt Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations |

Elements of Big Data

5 V's

1. **Volume** -- data is getting higher/bigger than ever.
2. **Velocity** -- data is increasing e.g. Complex real time data.
3. **Variety** -- data is spiraling e.g. unstructured video & voice.
4. **Veracity** – certainty of data
5. **Value** -- how much useful is the data



Elements of Big Data

10 V's





Volume

- Volume is the amount of data generated by organizations or individuals.
- Today, the volume of data in most organizations is approaching Exabytes.
- Every minute, over 571 new websites are being created.
- It is estimated that by 2020, online business transactions will reach up to 450 billion per day.
- Internet alone generates a huge amount of data.
- Internet has around 14.3 trillion live web pages.
- Internet has around 672 exabytes of accessible data.
- Total data stored on the internet is over 1 yottabyte.

Memory capacity hierarchy and conversion chart

| UNIT | ABBREVIATION | APPROXIMATE SIZE |
|------------|--------------|-----------------------------|
| bit | b | Binary digit, single 1 or 0 |
| nibble | — | 4 bits |
| byte/octet | B | 8 bits |
| kilobyte | KB | 1,024 bytes or 10^3 bytes |
| megabyte | MB | 1,024 KB or 10^6 bytes |
| gigabyte | GB | 1,024 MB or 10^9 bytes |
| terabyte | TB | 1,024 GB or 10^{12} bytes |
| petabyte | PB | 1,024 TB or 10^{15} bytes |
| exabyte | EB | 1,024 PB or 10^{18} bytes |
| zettabyte | ZB | 1,024 EB or 10^{21} bytes |
| yottabyte | YB | 1,024 ZB or 10^{24} bytes |

Volume

40 ZETTABYTES

[43 TRILLION GIGABYTES]

of data will be created by 2020, an increase of 300 times from 2005



6 BILLION PEOPLE
have cell phones



WORLD POPULATION: 7 BILLION



Volume SCALE OF DATA



It's estimated that

2.5 QUINTILLION BYTES

[2.3 TRILLION GIGABYTES]

of data are created each day



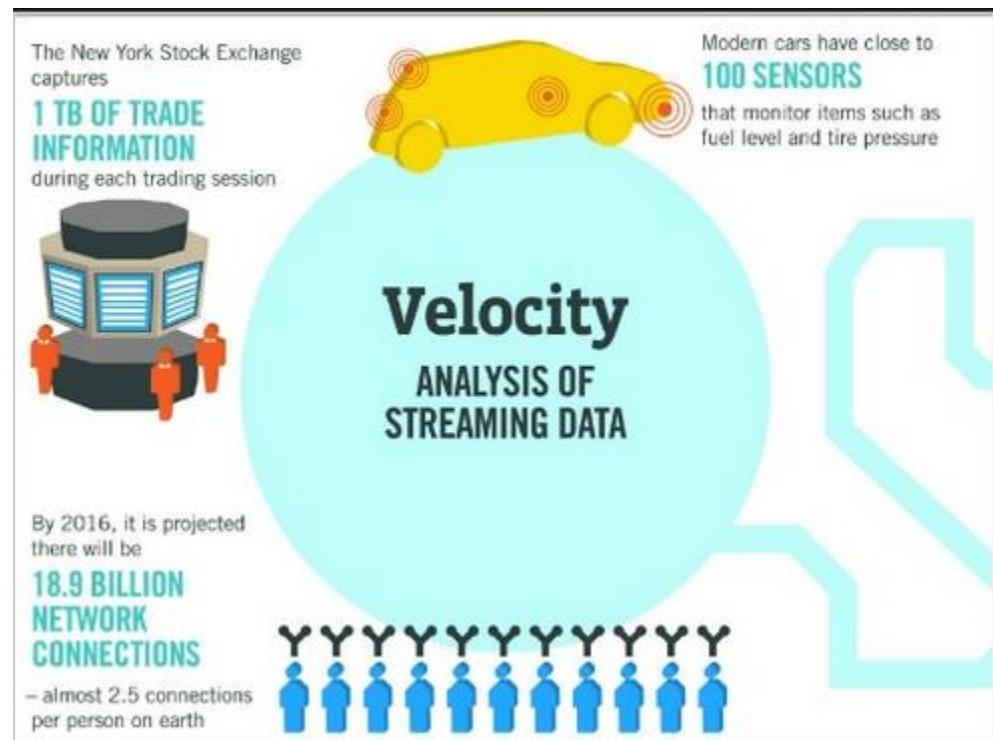
Most companies in the U.S. have at least

100 TERABYTES

[100,000 GIGABYTES]
of data stored

Velocity

- Velocity refers to the speed at which new data is generated and the speed at which data moves around.
- In other words, the rate at which data is generated, captured and shared in real time.





Variety

- Variety refers to the different types of data.
 - Structured data such as financial data Vs unstructured (80% of the world's data - text, images, video, voice, etc.)
- With big data technology we can now analyse and bring together data of different types such as messages, social media conversations, photos, sensor data, video or voice recordings.



Variety

As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES

[161 BILLION GIGABYTES]



**30 BILLION
PIECES OF CONTENT**

are shared on Facebook every month



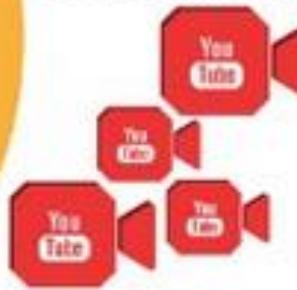
Variety

DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be

**420 MILLION
WEARABLE, WIRELESS
HEALTH MONITORS**

**4 BILLION+
HOURS OF VIDEO**
are watched on YouTube each month



400 MILLION TWEETS

are sent per day by about 200 million monthly active users.





Veracity

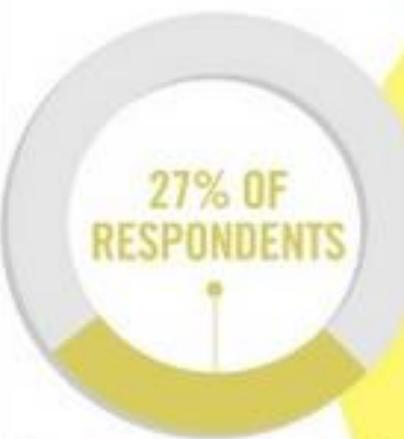
- Veracity refers to the uncertainty or trustworthiness of the data; ie, whether the obtained data is correct or consistent.
- In other words, it refers to the Data Accuracy.
- With many forms of big data quality and accuracy are less controllable (just think of Twitter posts with hash tags, abbreviations, typos and colloquial speech as well as the reliability and accuracy of content) but technology now allows us to work with this type of data.



Veracity

1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions

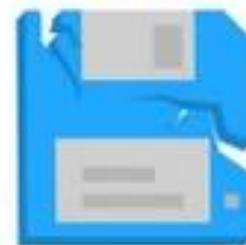


27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate

Veracity UNCERTAINTY OF DATA

Poor data quality costs the US economy around
\$3.1 TRILLION A YEAR

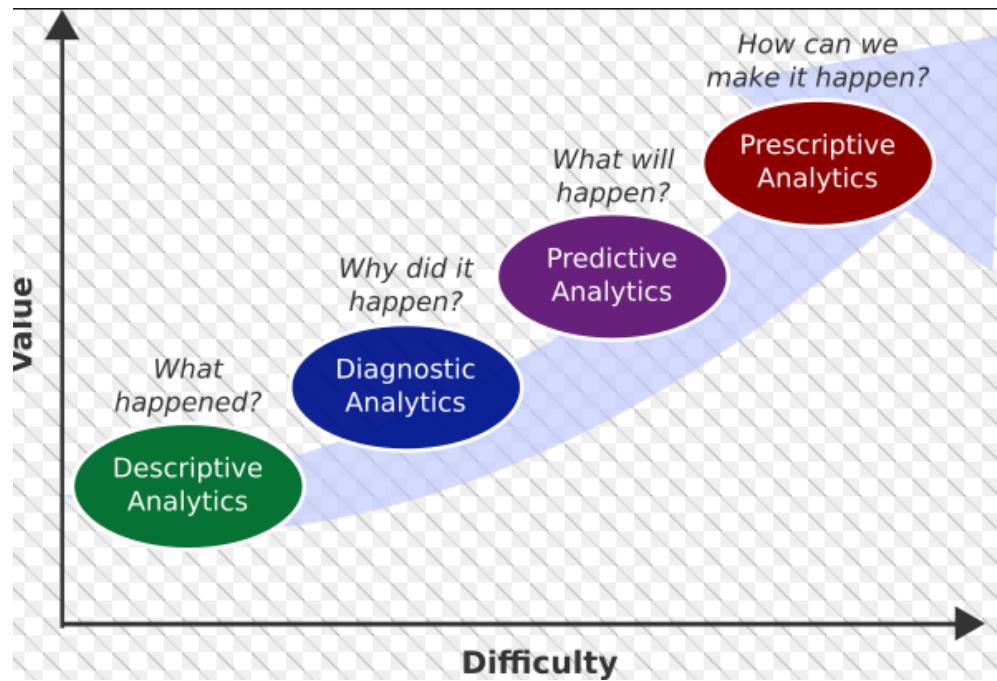


Big Data Analytics

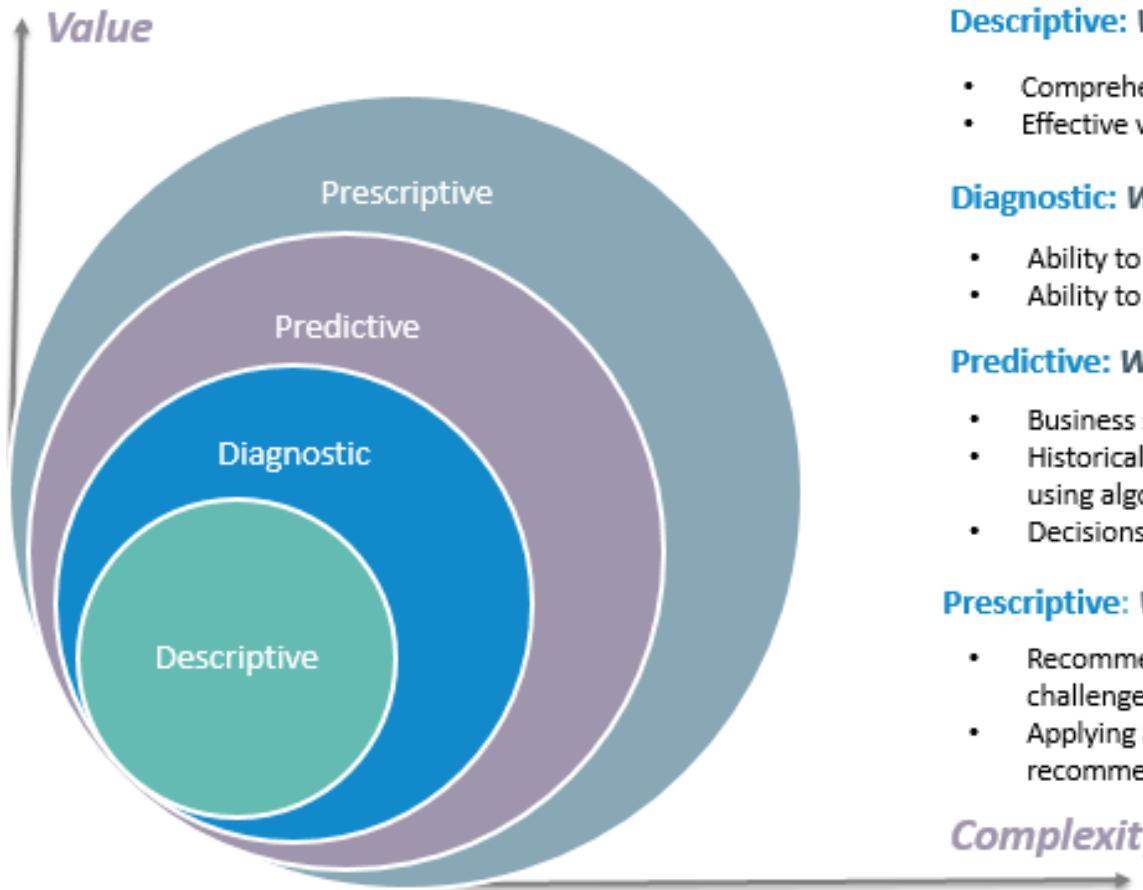
- Big data analytics helps to conduct business in many ways, such as it improves decision making, business process management etc...
- Business analytics uses the data and different techniques like information technology, features of statistics, quantitative methods and different models to provide results.

Big Data Analytics

- Main types of business analytics are
 1. Descriptive Analytics
 2. Diagnostic Analytics
 3. Predictive Analytics
 4. Prescriptive Analytics



4 types of Data Analytics



What is the data telling you?

Descriptive: *What's happening in my business?*

- Comprehensive, accurate and live data
- Effective visualisation

Diagnostic: *Why is it happening?*

- Ability to drill down to the root-cause
- Ability to isolate all confounding information

Predictive: *What's likely to happen?*

- Business strategies have remained fairly consistent over time
- Historical patterns being used to predict specific outcomes using algorithms
- Decisions are automated using algorithms and technology

Prescriptive: *What do I need to do?*

- Recommended actions and strategies based on champion / challenger testing strategy outcomes
- Applying advanced analytical techniques to make specific recommendations

Complexity

Descriptive Analytics

- It is the most widespread form of analytics and it serves as a base for advanced analytics.
- It answers the question **“What happened in the business”**.
- It analyzes a database to provide information on the trends of past or current business events that can help managers, planners, leaders etc. to develop a road map for future actions.
- It performs an in-depth analysis of data to reveal details such as frequency of events, operation costs, and the underlying reason for failures.
- It helps in identifying the root cause of the problem.

Diagnostic analytics

- Diagnostic analytics are used for discovery or to determine **why something happened.**
- While Descriptive Analytics might be concerned with describing how large or significant a particular outcome, Diagnostic Analytics will be more focused on determining what factors and events contributed to the outcome.
- Much of Diagnostic Analytics lives in the area of probabilities, likelihoods, and the distribution of outcomes.

Predictive Analytics

- Predictive analytics is about understanding the predicting the future and answers the question “**what will happen**” by using statistical models and different forecast techniques.
- It predicts the near future probabilities and trends and helps in what-if analysis.
- (**What-If Analysis** is the process of changing the values in cells to see how those changes will affect the outcome).
- In predictive analytics, we use different statistics, data mining techniques and machine learning to analyze the future.

Prescriptive Analytics

- Prescriptive analysis answers “**what should we do**” on the basis of complex data obtained from descriptive and predictive analyses.
- By using the optimization technique, prescriptive analytics determines the finest substitute to minimize or maximize some equitable finance, marketing and many other areas.

Advantages of Big data Analytics

- The right analysis of the available data can improve major business processes in various ways.
- For example, in a manufacturing unit, data analytics can improve the functioning of the following processes.
- **Procurement:** To find out which suppliers are more efficient and cost-effective in delivering products on time.
- **Product Development:** To draw insights or innovative product and service formats and designs for enhancing the development process and coming up with demanded products.

Advantages of Big data Analytics

- **Manufacturing:** To identify machinery and process variations that may be indicators of quality problems.
- **Distribution:** To enhance supply chain activities.
- **Marketing:** To identify which marketing campaigns will be the most effective in driving and engaging customers and understanding customer behaviors and channel behaviors.
- **Human resources:** To find out the characteristics and behaviors of successful and effective employees, as well as other employee insights for managing talent better.

Advantages of Big data Analytics

- **Price management:** To optimize prices based on the analysis of external factors.
- **Merchandising:** To improve merchandise breakdown on the basis of current buying patterns and increase inventory levels and product interest insights on the basis of the analysis of various customer behaviors.
- **Sales:** To optimize assignment of sales resources and accounts, product mix, and other operations.
- **Storage operations:** To adjust inventory levels on the basis of predicted buying patterns, study of demographics, weather, key events and other factors.

Advantages of Big data Analytics



Transportation: Big data has greatly improved transportation services. The data containing traffic information is analyzed to identify traffic jam areas.

Distributed sensors are installed in handheld devices, on the roads and on vehicles to provide real-time information.

This data is analyzed and passed to the traffic control authority.

Advantages of Big data Analytics



Education:

- e-learning helps to analyze the students' ability and thus impart education effectively in accordance with each student's needs.
- The analysis can be done by studying the responses to questions, recording the time consumed in attempting those questions, and analyzing other behavioral signals of the students.

Advantages of Big data Analytics



Travel

- Maintains complete details of all customer records that are then analyzed to determine certain behavior patterns in customers.



- For example, In the airline industry, big data is analyzed for identifying personal preferences or spotting which passengers like to have window seats.
- Big data track customers who regularly fly between specific routes to understand their behavior, what kind of items are preferred.

Advantages of Big data Analytics

Healthcare:

- The pharmacy and medical device companies use big data to improve their research and development practices.
- Big data helps doctors to analyze the requirement and medical history of every patient and provide individualistic services to them, depending on their medical condition.

Advantages of Big data Analytics

Telecom

- The mobile revolution and the Internet usage on mobile phones have led to a tremendous increase in the amount of data generated in the telecom sector.
- Big data analytics allows telecom industries to utilize this data for extracting meaningful information that could be used to understand their performance, improving customer service, and generating more business opportunities.

Advantages of Big data Analytics



Consumer Goods Industry:

- Consumer goods companies generate huge volume of data in varied formats from different sources, such as transactions, billing details, feedback forms etc.
- For example the data generated from the Point-of-Sale (POS) systems provides significant real-time information about customer's preferences, current market trends, the increase and decrease in demand of different products at different regions etc...

Introduction to Data Mining

- ▶ **Data mining** is one of the most useful techniques that help entrepreneurs, researchers, and individuals *to extract valuable information from huge sets of data.*
- ▶ Data mining is also called ***Knowledge Discovery in Database (KDD)***.
- ▶ The **knowledge discovery process** includes Data cleaning, Data integration, Data selection, Data transformation, Data mining, Pattern evaluation, and Knowledge presentation.

What is Data Mining?

- The *process of extracting information* to identify patterns, trends, and useful data that would allow the business to take the data-driven decision from huge sets of data is called **Data Mining**.

- ▶ Data mining is the act of automatically searching for large stores of information to find trends and patterns that go beyond simple analysis procedures.
- ▶ Data mining utilizes complex mathematical algorithms for data segments and evaluates the probability of future events.
- ▶ Data Mining is a process used by organizations to extract specific data from huge databases to solve business problems.
- ▶ It primarily turns **raw data into useful information.**

Benefits of data mining

- ▶ The business benefits of data mining come from the increased ability to uncover hidden patterns, trends, correlations and anomalies in data sets.
- ▶ *That information can be used to improve business decision-making and strategic planning through a combination of conventional data analysis and predictive analytics.*

Specific data mining benefits include the following:

- ▶ **More effective marketing and sales :** Data mining helps marketers better understand customer behavior and preferences, which enables them to create targeted marketing and advertising campaigns.

- ▶ Similarly, sales teams can use data mining results to improve lead conversion rates and sell additional products and services to existing customers.
- ▶ **Better customer service** :Companies can identify potential customer service issues more promptly and give contact center agents up-to-date information to use in calls and online chats with customers.
- ▶ **Improved supply chain management** : Organizations can spot market trends and forecast product demand more accurately, enabling them to better manage inventories of goods and supplies.
- ▶ Supply chain managers can also use information from data mining to optimize warehousing, distribution and other logistics operations.

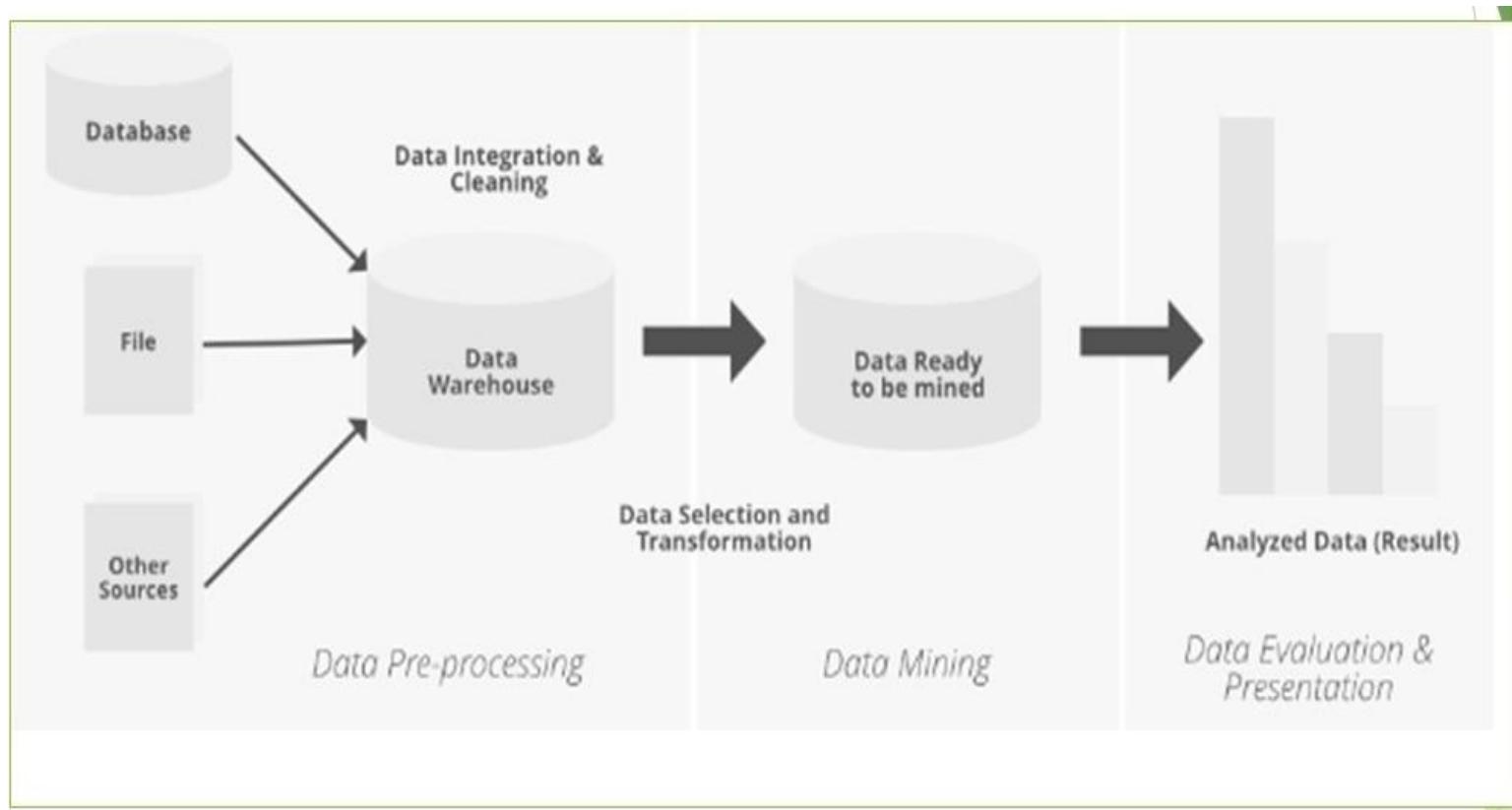
- ▶ **Increased production uptime** : Mining operational data from sensors on manufacturing machines and other industrial equipment supports predictive maintenance applications to identify potential problems before they occur, helping to avoid unscheduled downtime.
- ▶ **Stronger risk management** : Risk managers and business executives can better assess financial, legal, cybersecurity and other risks to a company and develop plans for managing them.
- ▶ **Lower costs** : Data mining helps drive cost savings through operational efficiencies in business processes and reduced redundancy and waste in corporate spending.

Data Mining Process

KNOWLEDGE DISCOVERY (KDD) PROCESS

Here is the list of steps involved in the knowledge discovery process:

- **Data Cleaning** – In this step, the noise and inconsistent data is removed.
- **Data Integration** – In this step, multiple data sources are combined.
- **Data Selection** – In this step, data relevant to the analysis task are retrieved from the database.
- **Data Transformation** – In this step, data is transformed or consolidated into forms appropriate for mining by performing summary or aggregation operations.
- **Data Mining** – In this step, intelligent methods are applied in order to extract data patterns.
- **Pattern Evaluation** – In this step, data patterns are evaluated.
- **Knowledge Presentation** – In this step, knowledge is represented.



1. Data Cleaning

2. Data Integration

1. Data Cleaning:

- **Removing noisy data**
- **Fill the missing values**
- **Remove inconsistent data**

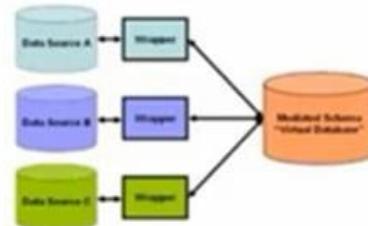
- Remove Noise and Inconsistent Data



2. Data Integration:

- **The process of combining data from multiple heterogeneous data sources is known as Data integration**

- Where multiple data sources may be combined



3. Data Selection

4. Data Transformation

- Where data relevant to the analysis task are retrieved from the data base



4. Data Transformation:

- Data need to be transferred into appropriate form for various applications**

3. Data Selection:

- The process of selecting task relevant data.**

- Where data are transformed and consolidated into forms appropriate for mining by performing summary or aggregation operation



5. Data Mining

6. Pattern Evaluation

5. Data Mining :

- **The functionalities such as characterization, classification, regression, association, clustering, outlier analysis applied to mine data.**

An essential Process where intelligent methods are applied to extract data patterns



6. Pattern evaluation:

- All the discovered data may not be interested or we may have discovered huge amount data as interested. Thus interesting measures such as support, confidence, lift etc used for finding only interested patterns.

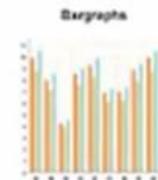
- To identify the truly interesting patterns representing knowledge based on interestingness measures



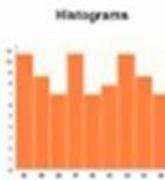
7. Knowledge representation

7. knowledge presentation:

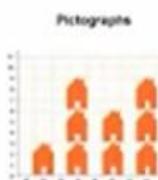
- Discovered knowledge presented visualization of discovered knowledge in the form of pie charts, bar charts, curves, plots, tables, cross tabs, trees etc.



Bar graphs



Histograms

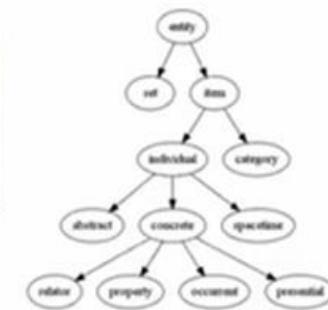


Pictographs



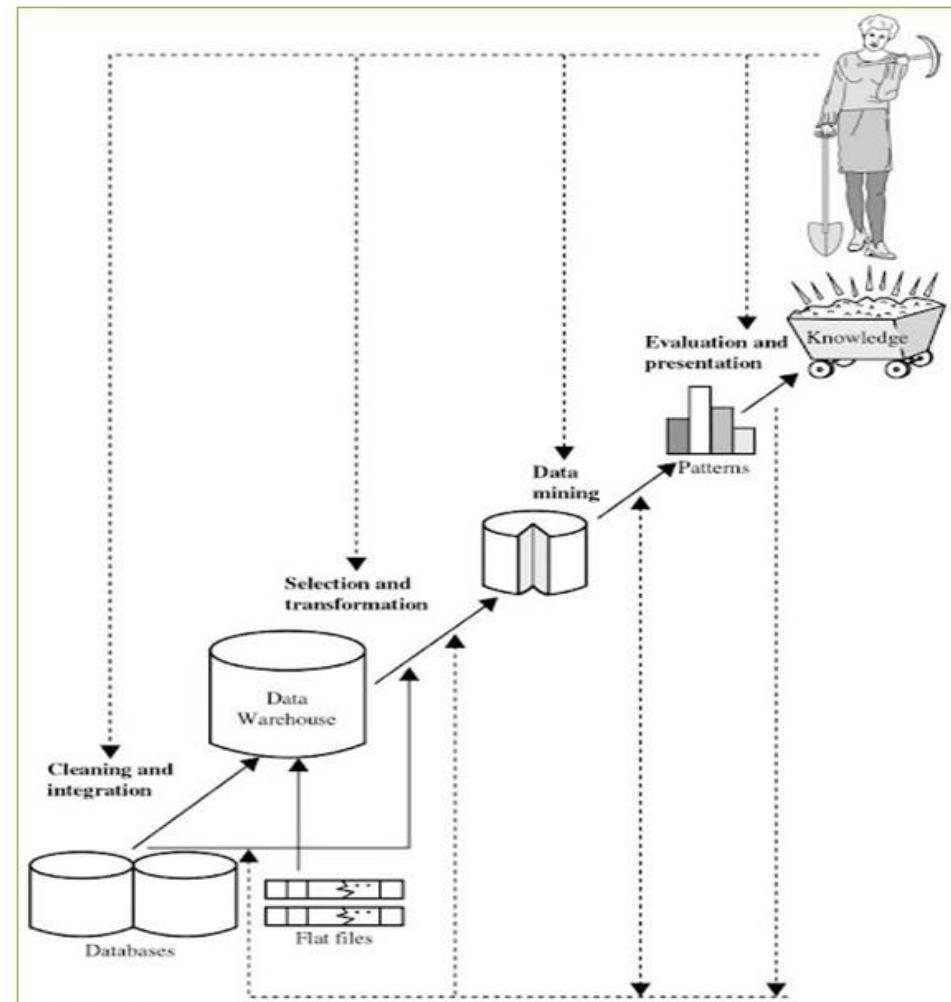
Pie Charts

| | | Age | | | | | | |
|--------------------|-------|----------|-------|-------|-------|-------|-----|--|
| | | Under 10 | 10-14 | 15-19 | 20-24 | 25-29 | 30+ | |
| Genre | Total | 18 | 47 | 23 | 13% | 20 | 8 | |
| | | 29% | 22% | 13% | 13% | 10% | 4% | |
| Frequency of visit | | | | | | | | |
| Daily | | 12 | 8 | 3 | 4 | 2 | 1 | |
| Twice a week | | 15 | 11 | 8 | 5 | 3 | 2 | |
| Weekly | | 14 | 10 | 8 | 5 | 3 | 2 | |
| Monthly | | 17 | 13 | 11 | 9 | 5 | 3 | |



- IF-THEN
- Independent part, combined with other pieces, to produce better result
- Model of human behavior
- Examples
 - IF condition, THEN conclusion
 - Conclusion, IF condition
 - If condition, THEN conclusion1 (OR) ELSE conclusion2

KDD process

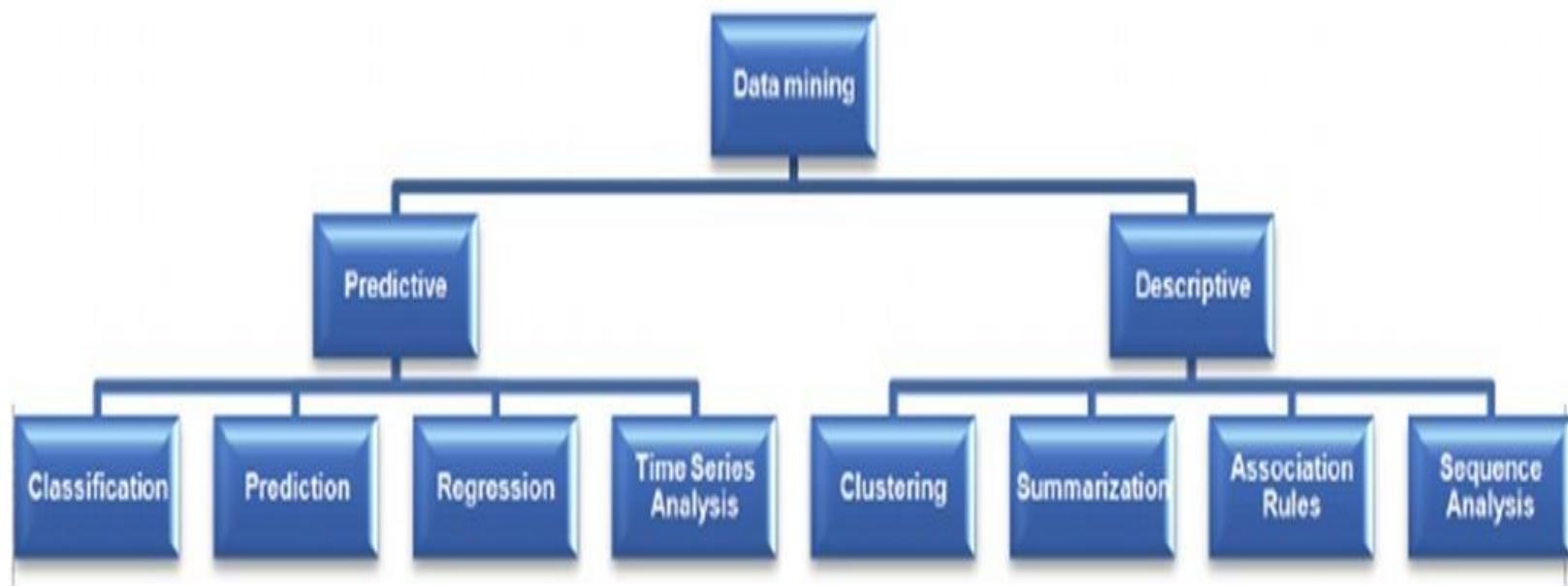


Data Mining Techniques

- To extract information from data, a wide variety of data mining techniques are employed.

These include:

1. Classification
2. Clustering
3. Association rule
4. Regression
5. Anomaly(outlier) detection
6. Prediction



1. Clustering

- ▶ **Clustering** is a *technique for grouping related entries in a database into clusters based on their similarities.*
- ▶ Whereas classification assigns variables into known categories, **the clustering technique first singles out these clusters in the dataset and then groups variables based on their characteristics.**
- ▶ **For example, you can cluster customers into groups according to the sales data – those who regularly buy pet food or specific drinks and who are stable in their preferences and customer behavior.**
- ▶ Once you establish these clusters, you can easily target them with customized advertisements.

2. Classification

- ▶ Classification is used to divide data into **predetermined groups or classes**.
- ▶ *This data mining technique determines* the class to which a record belongs based on the values of several attributes.
- ▶ The goal is to sort data into predefined classes.

Clustering has a wide range of applications:

- ▶ medical diagnostics
- ▶ computational biology
- ▶ text mining
- ▶ web analytics
- ▶ **Another example of clustering**, there are two clusters named as mammal and reptile. A mammal cluster includes human, leopards, elephant, etc.
- ▶ On the other hand, reptile cluster includes snakes, lizard, komodo dragon etc.
- ▶ **The tools mainly used in cluster analysis** are k-mean, k-medoids , density based, hierarchical and several other methods.

- ▶ For example, in a banking application, the customer who applies for a loan may be classified as a safe and risky according to his/her age and salary.
- ▶ This type of activity is also called **supervised learning**.
- ▶ The constructed model can be used to classify new data.
- ▶ The **learning step** can be accomplished by using already defined training set of data.
- ▶ Each record in the training data is associated with an attribute referred to as a class label, that signifies which class the record belongs to.
- ▶ Some **common classification algorithms** are *decision tree*, *neural networks*, *logistic regression*, etc.

What is this?

And what is this?

What do we understand from this example?

- A 5-year old boy can adapt and react to different objects.
- We want to teach similar adaptation to a *Machine* to recognize and provide actions.

• apple
• orange

Discriminatory power of humans/machines to recognize objects, is classification.

Bag of Fruits, can we arrange?

We want to arrange similar fruits together

We found 3 bags

This is called Clustering

Key Differences Between Classification and Clustering

1. Classification is the process of classifying the data with the help of class labels. On the other hand, Clustering is similar to classification but there are no predefined class labels.
2. Classification is geared with supervised learning. As against, clustering is also known as unsupervised learning.
3. Training sample is provided in classification method while in case of clustering training data is not provided.

- ▶ Classification and clustering are the methods used in data mining for analyzing the data sets and divide them on the basis of some particular classification rules or the association between objects.
- ▶ Classification categorizes the data with the help of provided training data.
- ▶ On the other hand, clustering uses different similarity measures to categorize the data.

Examples

Netflix

A well-known application of clustering algorithms are Netflix recommendation systems. Although the company is quite discreet with its algorithms, it is confirmed that there are about 2,000 clusters or communities that have common audiovisual tastes. Cluster 290 is the one that includes people who like the series "Lost", "Black Mirror" and "Groundhog Day". Netflix uses these clusters to refine its knowledge of the tastes of viewers and thus **make better decisions in the creation of new original series.**

Fraud Detection

Classification is commonly used in the financial sector to guarantee data security. In the era of online transactions where the use of cash has decreased markedly, it is necessary to determine whether movements made through cards are safe. Entities can classify transactions as correct or fraudulent using historical data on customer behavior to **detect fraud very accurately.**

Regression

- Regression analysis is the data mining process, that is used to identify and analyze the relationship between variables because of the presence of the other factor.
- It is used to define the probability of the specific variable.
- Regression, primarily a form of planning and modeling. **For example**, we might use it to project certain costs, depending on other factors such as availability, consumer demand, and competition.
- Primarily it gives the exact relationship between two or more variables in the given data set.

Some more examples are:

Example 3:

Based on the expected number of customers and the previous days' data given below, predict the number of burgers that will be sold by a KFC outlet .



| No. of Customers | Burgers Sold |
|------------------|--------------|
| 798 | 444 |
| 450 | 324 |
| 5067 | 1054 |

Example 4:

Calculate the life expectancy for a group of people with the average length of schooling based on the data given below:



| Average Length of Schooling (in years) | Life Expectancy |
|--|-----------------|
| 14 | 63 |
| 20 | 80.4 |
| 10 | 57.1 |

- ▶ One example is forecasting customers' age based on their purchase history.
- ▶ We can also predict costs based on such variables as consumer demand – for example, a surge of prices on the secondary market due to the increased demand for cars in the US.

4. Anomaly detection(Outlier Analysis)

- ▶ **Anomaly detection** is a *data mining technique used to identify outliers (values that deviate from the norm).*
- ▶ **For example,** in e-commerce datasets, *it can detect unusual sales during a given week at a store location.*
- ▶ Among other things, it can be used to discover credit or debit fraud and identify intrusion or interruption in the network.

5. Association rule

- ▶ **Association rule** learning **discovers if-then patterns** between two or more variables.
- ▶ The simplest example is the association between buying bread and butter.
- ▶ *People that buy bread usually get butter with it, and vice versa.*
- ▶ *That is why you will find these two products close to one another in a grocery store.*
- ▶ However, the link may be not that direct.

- ▶ While the concepts behind association rules can be traced back earlier, association rule mining was defined in the 1990s, when computer scientists Rakesh Agrawal, Tomasz Imieliński and Arun Swami developed an algorithm-based way to find relationships between items using point-of-sale (POS) systems.
- ▶ Applying the algorithms to supermarkets, the scientists were able to discover links between different items purchased, called *association rules*, and ultimately use that information to predict the likelihood of different products being purchased together.
- ▶ *For retailers, association rule mining offered a way to better understand customer purchase behaviors.*
- ▶ Because of its retail origins, **association rule mining** is often referred to as *market basket analysis*.

6. Prediction

- ▶ Another process of data analysis is prediction.
- ▶ It is used to find a numerical output.

- Prediction uses a combination of other data mining techniques such as clustering, classification, etc.
- It analyzes past events or instances in the right sequence to **predict** a future event.

- ▶ Predicting the value of a house depending on the facts such as the number of rooms, the total area, etc., is an **example for prediction**.
- ▶ *For example*, suppose the marketing manager needs to predict how much a particular customer will spend at his company during a sale.

Big Data Vs Data Mining

| Data Mining | Big Data |
|--|---|
| <p>It is one of the method in the pipeline of Big Data.</p> | <p>Big Data is a technique to collect, maintain and process the huge information. It explains the data relationship.</p> |
| <p>Data mining is a part of Knowledge Discovery of the Data. It is close view of the data.</p> | <p>It is about extracting the vital and valuable information from huge amount of the data. It is a technique of tracking and discovering of trends of complex data sets. It is a large or overall view of the data.</p> |
| <p>The goal is same as Big Data as it is one of the tool of Big Data.</p> | <p>The goal is to make data more vital and usable i.e. by extracting only important information from the huge data within existing traditional aspects.</p> |

It is manual as well as automated in nature.

It is only automated as computing huge data is difficult.

It only focuses on only one form of data. i.e. structured.

It focuses and works with all form of data i.e. structured, unstructured or semi-structured.

It is used to create certain business insights. Data mining is a manager of the mine.

It is mainly used for business purposes and customer satisfaction. Big Data is a mine.

It is a sub set of Big Data. i.e. one of the tools.

It is a super set of Data Mining.

It is a tool to dig up the vital information from the large data. Data can be large as well as small.

It is more involved with the processes of handling voluminous data. Data can only be large.