# Agenda

- Course Overview
- Course Objectives and Outcomes
- Quick Review of Complete Handout
- Practical Aspects
- Different Types of Data and Storage for Data
- Big Data Characteristics, sources
- Big Data Systems Perspective – in-memory vs storage vs network
- Big Data challenges, applications/case studies
- Locality of Reference – principle examples

## BIG DATA ANALYTICS – emphasis / important ones

Introduction to Industry 4.0/5.0, BI & Analytics
Big Data Definition & Characteristics
Sources of Big Data
Challenges & Benefits of Big Data Systems
Case Studies / Applications

Hadoop Ecosystem & Map Reduce
Hadoop Intro
HDFS architecture
MAP REDUCE architecture
SPARK Usage

HIVE Database & Others
Hive Architecture and uses
Pig, Sqoop, Flume, OOzie
Other popular tools (open source & commercial)

Mongodb Architecture & Commands
Data Visualization;
Applications/projects – PowerBI/Tableau & Excel

**PRACTICAL / LAB EXAMPLES USING LINUX, SHELL, DB, SQL, JAVA, SQL, PYTHON, POWERBI.**
**MCQ – PUZZLES – CASE STUDIES/DISCUSSIONS**

**Focus Areas/Practicals**
✔ *VIRTUAL MACHINE*
✔ *UBUNTU LINUX COMMANDS*
✔ *EDITING FILES AND ENVIRONMENT*
✔ *SETTING UP HADOOP CLUSTER*
✔ *STARTING CLUSTER*
✔ *PRACTICING HADOOP AND HDFS COMMANDS*
✔ *SETTING UP AND RUNNING MAP REDUCE*
✔ *SETTING UP HIVE DB*
✔ *RUNNING QUERIES AGAINST HIVE DB*
✔ *EXPORT/IMPORT TO HIVE*
✔ *WORK WITH MONGODB*
✔ *QUICK DEMO… OF OTHER TOOLS…*

3

# A QUICK REVIEW OF DATA ANALYST & SCIENTIST

**Data Analysts and Scientists**
**Job Profile:** Data analysts help businesses develop well-informed strategies by creating charts and prepare visual presentations. Also, examining large data and identifying trends are the expected roles of a data analyst. Data scientists construct and develop new processes for data modelling and primarily use prototypes, algorithms, predictive models and custom analysis.
**Chief skills:** Python coding, Hadoop Platform, R programming, SQL Database/Coding, Apache spark, Machine learning and AI, Data Visualisation
**Big Data Specialists**
**Job Profile:** Utilise data analysis to evaluate the technical performance of an organisation. Also provides recommendations on system enhancements.
**Job skills:** Apache Hadoop, Apache Spark, NoSQL, Machine learning and Data Mining, Statistical and Quantitative Analysis, SQL, Data Visualization, General Purpose Programming language
**Digital Transformation Specialists**
**Job Profile:** Work in enhancing a company's technical performance. They analyse the company's infrastructure and the gaps in service.
**Job skills:** Technical aptitude, critical thinking abilities, excellent communication skills, adaptability, SQL, C++, HTML, CSS

## MATH & STATISTICS
- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

## PROGRAMMING & DATABASE
- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing package e.g. R
- ☆ Databases SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

## DOMAIN KNOWLEDGE & SOFT SKILLS
- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

## COMMUNICATION & VISUALIZATION
- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
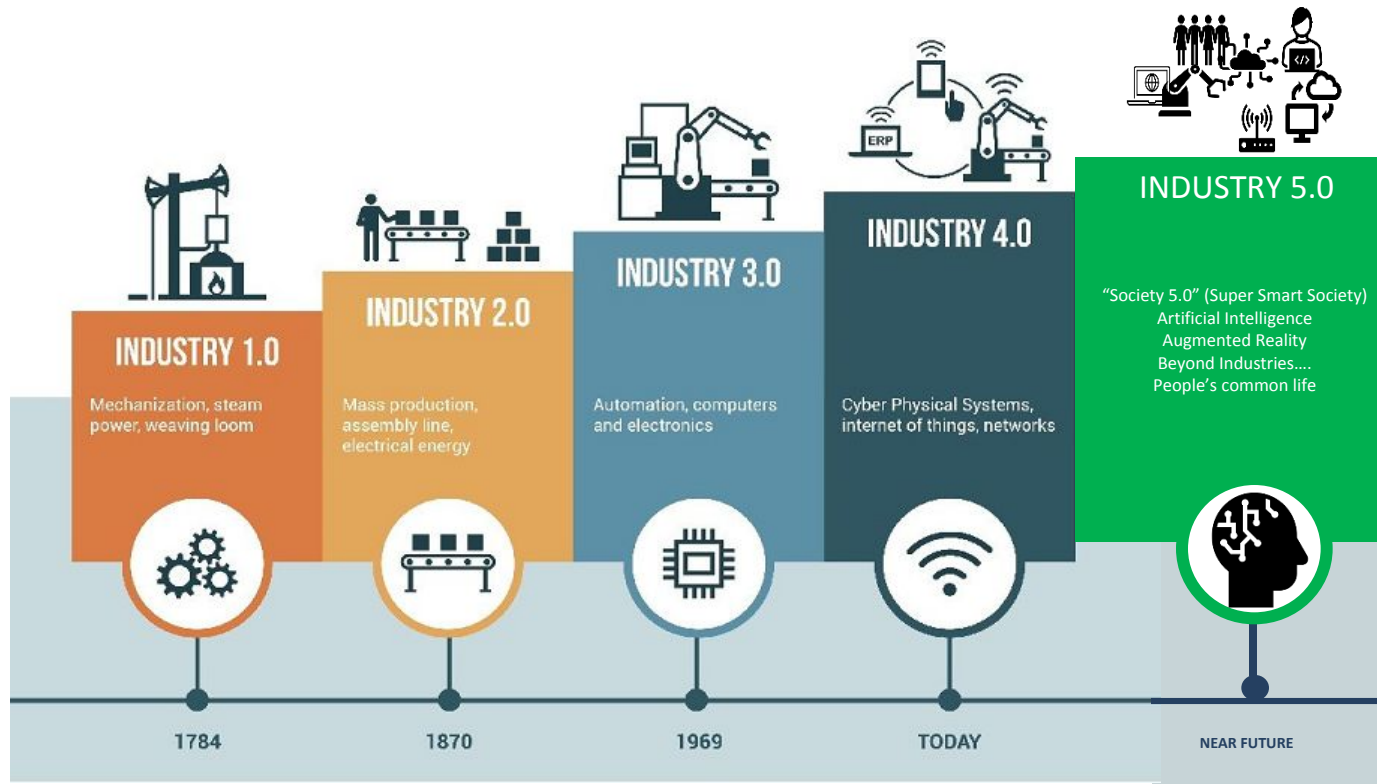- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

4

**Big Data Systems**

# A FEW million $ QUESTIONS



☐ **Data is woven into the everyday fabric of our lives. With the rise of mobile, social media, and smart technologies associated with the Internet of Things (IoT), we now transmit more data than ever before—and at a dizzying speed.** **Thanks to big data analytics!!**

☐ **Organizations can now use that information to rapidly improve the way they work, think, and provide value to their customers.**

☐ **With the assistance of tools and applications, big data can help you gain insights, optimize operations, and predict future outcomes.**
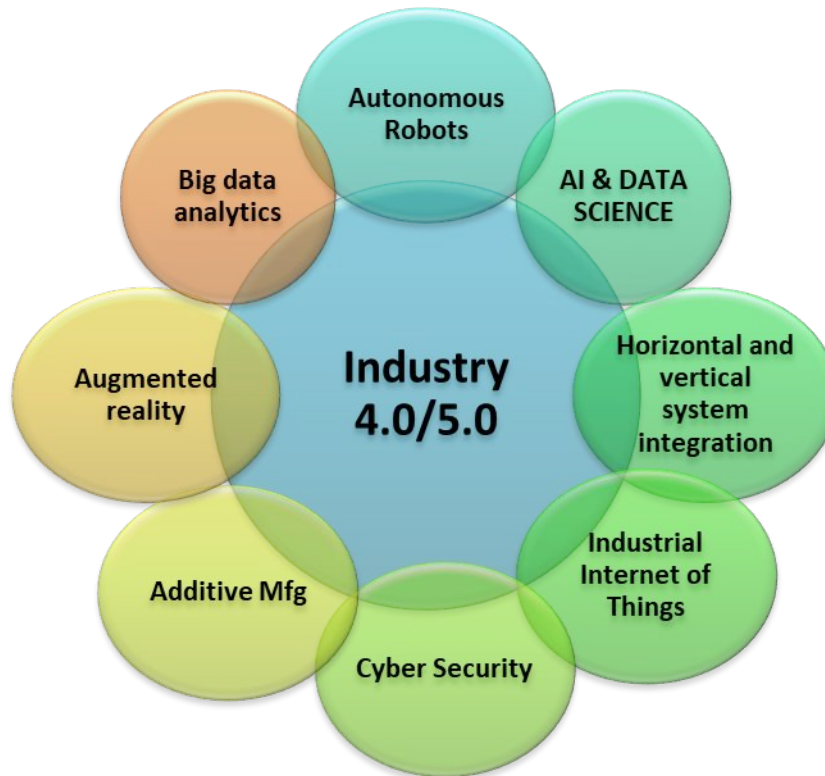
| | |
|---|---|
| Information | Internet Of Everything/Things |
| Science | Basics/Formulas/neurons.. |
| Engineering | Structural Approaches/principles/methods |
| Technology | Enablers… Wifi.. 5S.. etc |
| Intelligence | AI, ML, DL, ANN |
| Automation | Robots, Chatbots, Programs, Daemons, etc |
| Machine Intelligence | Making machine to think.. and communicate with other machines as well as human |
| B-B C-B..!? Management | Data Analytics, Visualization, Statistics |

**Societal Values & Cost Savings & New Way of Lazy Life??**
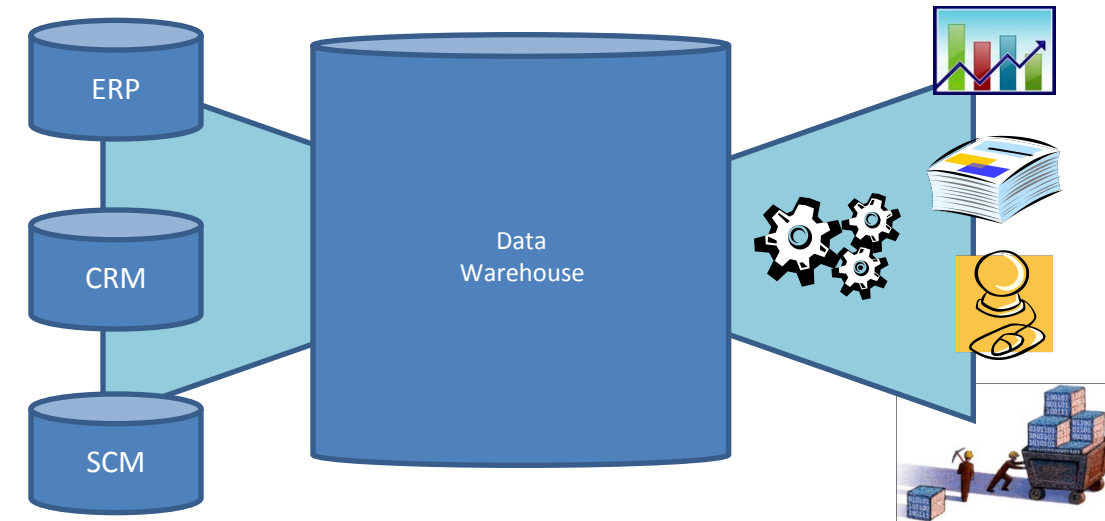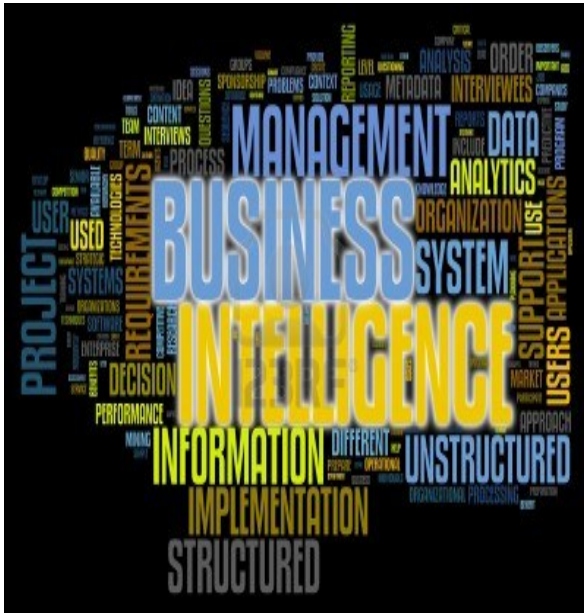
# EVOLUTION OF INDUSTRIES



**INDUSTRY 1.0**
Mechanization, steam power, weaving loom

**INDUSTRY 2.0**
Mass production, assembly line, electrical energy

**INDUSTRY 3.0**
Automation, computers and electronics

**INDUSTRY 4.0**
Cyber Physical Systems, internet of things, networks

**INDUSTRY 5.0**
"Society 5.0" (Super Smart Society)
Artificial Intelligence
Augmented Reality
Beyond Industries....
People's common life

1784    1870    1969    TODAY    NEAR FUTURE

6

# BUILDING BLOCKS OF INDUSTRY 4.0 & 5.0



All the building blocks create a huge opportunities to the universities/engineering institutes to prepare or update the curriculum and teaching methodologies to train the engineers to be industry 4.0 read & 5.0 ready.

7

# What Is Business Intelligence & How it works?



ERP

CRM

SCM

Data Warehouse

"Getting data in"

"Getting data out"

# The Scope of Business Intelligence



**Smaller organizations:**
**Excel spreadsheets**

**Larger organizations:**
**Data mining, Predictive,**
**Prescriptive, analytics,**
**dashboards**

# WHO IS DATA SCIENTIST ??

- <mark>There are two basic types of Predictive Analytics / Data Science problems:</mark>

- 1. Internal Predictive Analytics / Data Science problems, such as bad data, reckless analytics, or using inappropriate techniques.

- Internal problems are not business problems; they are internal to the Predictive Analytics / Data Science community.

- Therefore, the fix consists in training data scientists to do better work and follow best practices.

- 2. Applied business problems are real-world problems for which solutions are sought, such as fraud detection or identifying if a factor is a cause or a consequence.

- These may involve internal or external (third-party) data.

- <mark>These are the characteristics of the modern trends in Predictive Analytics / Data Science which one you  should be aware of:</mark>

  - **In-memory analytics**
  - MapReduce and Hadoop
  - NoSQL, NewSQL, and graph databases
  - Python and R
  - Data integration: blending unstructured and structured data (such as data storage and security, privacy issues when collecting data, and data compliance)
  - Visualization
  - Analytics as a Service, abbreviated as AaaS
  - Text categorization/tagging and taxonomies to facilitate extraction of insights from raw text and to put some structure on unstructured data

- <mark>What Knowledge does a Data Scientist need:</mark>

- Thus, data scientists also need to be good communicators to understand, and many times guess, what problems their client, boss, or executive management is trying to solve.

- Translating high-level English into simple, efficient, scalable, replicable, robust, flexible, platform-independent solutions is critical.

- Predictive Analytics / Data Science = Some (computer science) + Some (statistical science)
- + Some (business management) + Some (software engineering) + Domain
- expertise + New (statistical science), where

- Some () means the entire field is not part of Predictive Analytics / Data Science.
- New () means new stuff from the field in question is needed

13

- <mark>Horizontal Versus Vertical Data Scientist</mark>

- <u>Vertical data scientists</u> have deep technical knowledge in some narrow field.

- For instance, they might be any of the following:

- Computer scientists familiar with computational complexity of all sorting algorithms

- Statisticians who know everything about eigenvalues, singular value decomposition and its numerical stability, and asymptotic convergence of maximum pseudo-likelihood estimators

- Software engineers with years of experience writing Python code (including graphic libraries) applied to API development and web crawling technology

14

- <mark>Horizontal Versus Vertical Data Scientist</mark>

- Database specialists with strong data modeling, data warehousing, graph database, Hadoop, and NoSQL expertise

- Predictive modelers with expertise in Bayesian networks, SAS, and SVM

- The key here is that by "vertical data scientist" we mean those with a more narrow range of technical skills,
- such as expertise in all sorts of Lasso-related regressions but with limited knowledge of time series, much less of any computer science.

- <mark>Horizontal Versus Vertical Data Scientist</mark>

- <u>Horizontal data scientists</u> are a blend of business analysts, statisticians, computer scientists, and domain experts.

- They combine vision with technical knowledge.

- They might not be experts in eigenvalues, generalized linear models, and
- other semi-obsolete statistical techniques,
- but they know about more modern data-driven techniques applicable to unstructured, streaming, and big data.

- They can design robust, efficient, simple, replicable, and scalable code and algorithms.

16

# Data Science

- Statistical and Operations research techniques
- Machine Learning
- Deep Learning

Data Science refers to an emerging area of work concerned with the **Collection, Preparation, Analysis, Visualization, Management, and Preservation** of **large collections of information**.

Data Science involves using methods to analyze **massive amounts of data** and **extract the knowledge** it contains.

Data Science includes **Data Analysis** as an important component of the skill set required for many jobs in this area, but is not the only necessary skill.

**Data Science is fundamentally an interdisciplinary subject**

17

Data Science involves using methods to analyze **massive amounts of data (Big Data)** and extract the knowledge it contains.

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

18

# Nearly everyone across the organization engages with software



**Yet, fewer than 25% of workers have access to analytical insights**

19

# BI Questions

> *Developing a business intelligence strategy is an important first step in implementing a BI solution.*

- ✔ *Who are the key stakeholders? Who will be using this system?*
- ✔ *What departments need business intelligence and what will be measured?*
- ✔ *What support do content authors and information consumers need?*
- ✔ *Focus on Questions That Are Aligned With Your Business Strategy*
- ✔ *Ask BI Questions That Give You Actionable Insights*
- ✔ *Questions that Identify Opportunities for ROI*
- ✔ *Questions That Identify Opportunities For New Sources of Revenue*
- ✔ *Questions Identifying Cost-cutting Opportunities*
- ✔ *Questions That Begin With "Why"*
- ✔ *How Can You Understand Your Customers Better?*

*As technology companies like Amazon, Meta, and Google continue to grow and integrate with our lives, they are leveraging big data technologies to monitor sales, improve supply chain efficiency and customer satisfaction, and predict future business outcomes.*

*Currently, there is so much big data that International Data Corporation (IDC) predicts the "Global Datasphere" will grow from 33 Zettabytes (ZB) in 2018 to 175 ZB in 2025. That's equal to a trillion gigabytes.*

What if you could empower everyone with analytics anywhere decisions are made?

# Today, BI extends to everyone



**3rd wave**
End user BI

Everyone

**2nd wave**
Self-service BI

Analyst to end user

1st wave
Technical BI

IT to end user

# Two Data Factors

```
                          ┌─── On-Premises
           Data Location ─┼─── Cloud
                          └─── Public

                          ┌─── Sharing
           Data Access  ──┤
                          └─── Row-Level Security
```

# The Good

Experiments, observations, and numerical simulations in many areas of science and business are currently generating terabytes of data, and in some cases are on the verge of generating petabytes and beyond. Analyses of the information contained in these data sets have already led to major breakthroughs in fields ranging from genomics to astronomy and high-energy physics and to the development of new information-based industries.
- *Frontiers in Massive Data Analysis, National Research Council of the National Academies*

# The Bad

Given a large mass of data, we can by judicious selection construct perfectly plausible unassailable theories—all of which, some of which, or none of which may be right.   - *Paul Arnold Srere*
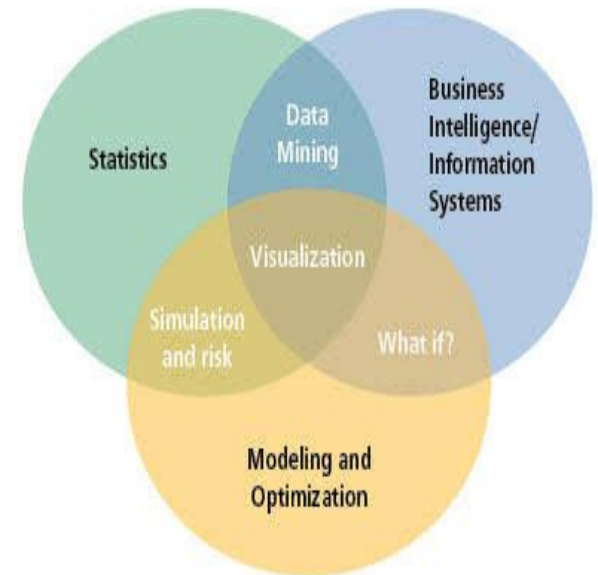
# DATA ANALYTICS vs BUSINESS ANALYTICS

- Data analytics is a broad umbrella for finding insights in data
- Data analytics can refer to any form of analysis of data—whether in a spreadsheet, database, or app—where the intent is to uncover trends, identify anomalies, or measure performance.
- Additional mathematics or IT skills can help data analysts do everything from managing a database of subscribers to calculating yields for a potential investment.
- Data analytics (DA) is the technical process of mining data, cleaning data, transforming data, and building the systems to manage data. Data analytics takes large quantities of data to find trends and solve problems. Data analytics is not just confined to business applications—it's used across disciplines, from the government to science.

- Business analytics focuses on identifying operational insights.
- Business analytics focuses on the overall function and day-to-day operation of the business.
- A business analyst would deal less with the technical aspects of analysis and more with the practical applications of data insights.
- Some job responsibilities might include creating a streamlined workflow or choosing the best vendors.
- Business analytics (BA) refers to the process of taking your company's raw data and turning it into useful information, including identifying trends, predicting outcomes, and more.



Statistics — Data Mining — Business Intelligence/Information Systems — Visualization — Simulation and risk — What if? — Modeling and Optimization

25

# DIFFERENCE BETWEEN BUSINESS ANALYTICS AND DATA ANALYSIS



| Business Analysis | Characteristics | Data Analysis |
|---|---|---|
| Refers to a person who uses data to employ concrete and practical decisions in a business. | Definition | Refers to a person tasked with collecting, processing,& analyzing already available data & producing vital insights that can improve business efficacy and solve existing problems. |
| Works on the frontlines of the data pipeline and helps businesses improve products, processes, and services. | Roles | The role of a data analyst involves data mining, data cleaning, using statistical techniques, managing data, and fixing bugs by designing databases and programs. |
| Identified weaknesses in organizational procedures and drives organizational growth through the analysis and measurement of growth. | Importance | Gathers information and creates actionable strategies for new and existing business opportunities. |
| Have good knowledge of programming and statistical tools. | Background | Have a strong background in statistics, math, and computer science. |

# DIFFERENCE BETWEEN DATA ANALYTICS AND BIG DATA ANALYTICS

| | DATA ANALYTICS | BIG DATA |
|---|---|---|
| **NATURE** | Like a book where you can find a solution to your problems. | considered as a Big Library where all the answers to all the questions are there but difficult to find the answers to your questions |
| **STRUCTURE OF DATA** | mostly structured data. It analyses the structured data to answer complex business queries, find solutions to business challenges, etc. | unstructured and raw data. The main aim of big data is to convert the raw data into meaningful data. |
| **TOOLS** | simple tools for statistical modelling and predictive modelling because the data to analyze is already structured and not complicated. | sophisticated technological tools such as automation tools or parallel computing tools to manage the Big Data |
| **TYPE OF INDUSTRY** | Used by IT Industries, Travel Industries, and Healthcare Industries. Data Analytics helps these industries to create new developments which are done by using historical data and analyzing past trends & patterns. | Used by industries such as banking industries, retail industries and many more. Big Data helps these industries in many ways to take some strategic business decisions. |

27

# Evolution of BI

# Characteristics of Data for Good Decision Making

## Better Quality Data — Characteristics

### Accuracy

**Accurate enough for intended purposes:**
- Balance with use, cost, effort, timeliness
- Capture close to point of activity
- Make accuracy compromises clear

### Validity

**Compliance with requirements:**
- Application of definitions
- Consistency over time
- Consistency with others

### Reliability

**Collection processes consistent:**
- ...over time
- ... for multiple collection points
- ... between collection systems

### Timeliness

**To Influence Decisions:**
- Capture quickly after the event
- Available quickly enough
- Available frequently enough

### Relevance

**Data relevant to intended purpose:**
- Periodic review of requirements
- Quality assurance and feedback process
- Use carefully for other purposes

### Completeness

**Monitor quality to match data needs:**
- Missing data
- Invalid data
- Incomplete data

29

# CRISP – DM
## Cross-Industry Standard Process for Data Mining



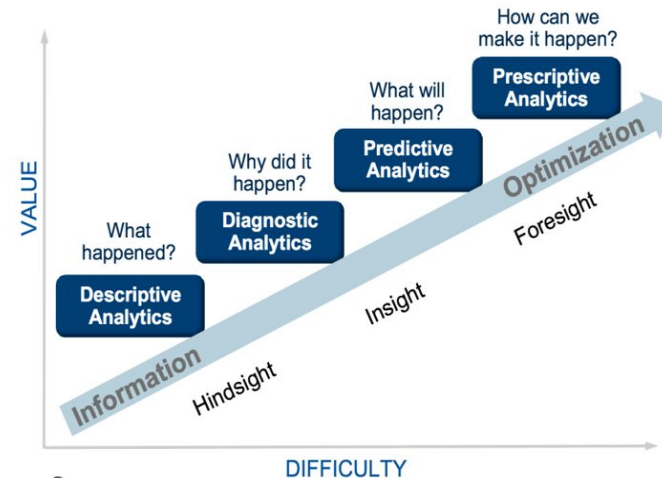| Business Understanding | Data Understanding | Data Preparation | Modeling | Evaluation | Deployment |
|---|---|---|---|---|---|
| **Determine Business Objectives** Background Business Objectives Business Success Criteria | **Collect Initial Data** Initial Data Collection Report | *Data Set* Data Set Description | **Select Modeling Technique** Modeling Technique Modeling Assumptions | **Evaluate Results** Assessment of Data Mining Results w.r.t. Business Success Criteria Approved Models | **Plan Deployment** Deployment Plan |
| **Situation Assessment** Inventory of Resources Requirements, Assumptions, and Constraints Risks and Contingencies Terminology Costs and Benefits | **Describe Data** Data Description Report | **Select Data** Rationale for Inclusion / Exclusion | **Generate Test Design** Test Design | **Review Process** Review of Process | **Plan Monitoring and Maintenance** Monitoring and Maintenance Plan |
| **Determine Data Mining Goal** Data Mining Goals Data Mining Success Criteria | **Explore Data** Data Exploration Report | **Clean Data** Data Cleaning Report | **Build Model** Parameter Settings Models Model Description | **Determine Next Steps** List of Possible Actions Decision | **Produce Final Report** Final Report Final Presentation |
| **Produce Project Plan** Project Plan Initial Asessment of Tools and Techniques | **Verify Data Quality** Data Quality Report | **Construct Data** Derived Attributes Generated Records | **Assess Model** Model Assessment Revised Parameter Settings | | **Review Project** Experience Documentation |
| | | **Integrate Data** Merged Data | | | |
| | | **Format Data** Reformatted Data | | | |

**Business Intelligence**
- Improves the decision making process
- Makes it easy to Access and Share information
- Helps you know your business
- Enables real-time analysis with Quick Navigation
- Reduces the risk of bottlenecks
- Helps identify waste in the system



**Descriptive**
Describes what happened.

**Diagnostic**
Describes why something happened.

**Predictive**
Describes what's going to happen.

**Prescriptive**
What we should do about it.



VALUE — DIFFICULTY

- Descriptive Analytics — What happened? — Information / Hindsight
- Diagnostic Analytics — Why did it happen? — Insight
- Predictive Analytics — What will happen? — Optimization
- Prescriptive Analytics — How can we make it happen? — Foresight

# BIG DATA SYSTEMS PERSPECTIVE…

Big Data Systems

# Systems Perspective: <span style="color:red">Processing Data</span>

**In-Memory Processing**

**Characteristics**:

– Data processed directly in RAM, avoiding disk I/O.

– Extremely fast (low latency).

**Advantages**:

– Ideal for real-time analytics and low-latency applications (e.g., Spark, Apache Flink).

– Supports iterative algorithms and machine learning.

**Challenges**:

– Limited by available RAM.

– More expensive compared to disk-based solutions.

# Systems Perspective: Processing from Secondary Storage

**Characteristics**:
- Data processed from hard drives or SSDs.
- Disk I/O introduces latency.

**Advantages**:
- Can handle massive datasets that don't fit in memory.
- Suitable for batch processing (e.g., Hadoop MapReduce).

**Challenges**:
- Slower compared to in-memory processing.
- Requires optimized data locality and access patterns.

# Systems Perspective: Processing over the network

**Characteristics**:
- Distributed processing across multiple nodes in a network.
- Data often stored in HDFS, S3, or similar distributed systems.

**Advantages**:
- Scalability: Can process petabytes of data by leveraging many nodes.
- Redundancy: Fault-tolerant systems with data replication.

**Challenges**:
- Network latency can become a bottleneck.
- Requires efficient task scheduling and data shuffling (e.g., Apache Hadoop, Spark).

# LOCALITY REFERENCES…

Big Data Systems

# Locality of Reference: Principle & Examples

- - **Definition**: Locality of reference is a principle in computing that describes how programs tend to access a relatively small portion of their address space at any given time.

- - **Types of Locality**:

-   - **Temporal Locality**: Recently accessed data is likely to be accessed again soon.

-   - **Spatial Locality**: Data near a recently accessed location is likely to be accessed soon.

- - **Examples**:

-   - **Code**: Sequential instruction execution (loops).

-   - **Data**: Consecutive array accesses in loops.

-   - **Memory Allocation**: Reuse of stack/heap data.

# Impact of Locality on Performance

- - **Reduced Latency**:
- - Data in the CPU cache is faster to access than RAM or disk.
- - Better locality = fewer cache misses = reduced latency.
- - **Optimized Resource Usage**:
- - CPU pipelines stay efficient.
- - Reduced memory bandwidth contention.
- - **Examples**:
- - Loop unrolling and blocking in matrix multiplication.
- - Optimized database query plans.

# Algorithms & Data Structures Leveraging Locality

- - **Sorting Algorithms**:
- - Merge Sort benefits from spatial locality during merging phases.
- - **Search Structures**:
- - B-Trees/B+ Trees: Designed for efficient disk access.
- - **Dynamic Programming**:
- - Uses temporal locality by storing reusable subproblem solutions.

# Data Organization for Better Locality

- - **On-Disk Data Layout**:

- - Contiguous allocation for files (e.g., ext4, NTFS).

- - Index structures like clustered B-Trees for databases.

- - **In-Memory Data Structures**:

- - Arrays vs. Linked Lists: Arrays have better spatial locality.

- - Cache-Aware Algorithms: Tailored for specific cache sizes.

# Mitigating Latency Through Locality Optimization

- - **Software Optimizations**:
- - Code restructuring to improve data locality.
- - Using cache-friendly algorithms (e.g., blocking in matrix operations).
- - **Hardware Optimizations**:
- - Multi-level caches (L1, L2, L3).
- - Prefetching mechanisms.
- - **Real-World Applications**:
- - High-performance computing.
- - Database systems optimized for query efficiency.