

UNIT-1

(1) (2)

- Big Data → 3V → high volume
[high velocity]
[high variety] → demands → cost effective
innovative forms
of information processing
for better decisions.

→ Black Box - Planes / Helicopters
→ Social Media Data
→ Stock Exchange Data → Search Engine Data

Types of Big Data

1. Structured Data [organized data
[predefined format
[defined repeating pattern
Data stored in fixed field in a record / file
Entities & their attributes are mapped.
Sources - Sensor data (GPS/RFID)
Weblog / Financial data
Stock market
Comp/Machine generated
↳ Point of sales data → barcode
cardswipe
Human generated - Samsung Health
in interaction with Comp
↳ IP data
↳ Click Stream - whenever click a link.
2. Unstructured Data [No logical, repeating pattern
[different formats - email / text / audio / video
Sources - Mobile / Social Media
Machine generated - Satellite image
Scientific data
Photograph / video
Radar / Sonar
Human generated - Social Media / Mobile Data / Texts, website content
3. Semi Structured Data [Schema less / self describing structure
Data stored inconsistently in rows and columns.

Types of Sources

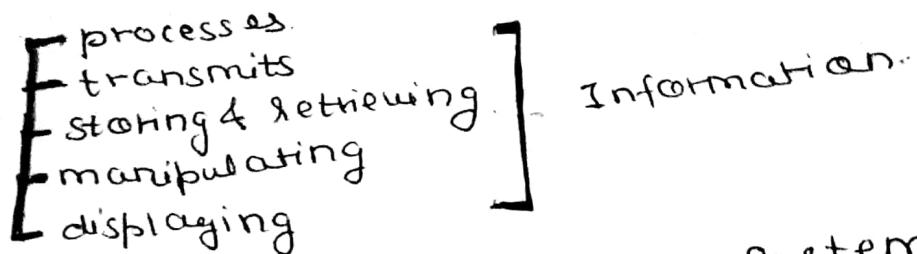
1. Internal Source [structured / organized data
[originates from within enterprise.
CRM / ERP
2. External Source [unstructured / unorganized data
[originates from external environment
of organization.
Ex: Business Partner
Internet / Govt'

• Information System

— system used for organizing & processing information
information technology + people's activities using technology to support operations / management
(combination of)

→ Interaction b/w people, algorithmic processes, data and technology.

→ Mediator b/w actions and technology.


processes
F transmits
F storing & retrieving
F manipulating
F displaying

Information

• Reasons to Use Information System

1. Operational Excellence

→ improves efficiency of operations
→ achieves higher profitability.
→ high level of efficiency & productivity

} can predict customer's choice

2. Improved Decision Making

3. Customer-Supplier Good Relations

3. New Product, Services

Elements of Big Data

1. Volume - amount of data
Nowadays → Petabytes
↳ Exabytes
By 2020
25 zettabytes
2. Velocity - speed of data processing.
↳ velocity at which it is created & integrated.
↳ rate at which data is generated & shared.
3. Variety - No. of types of data.
↳ data is generated from different types of sources [external & internal]
4. Veracity → uncertainty of data
↳ whether data obtained is correct or consistent.

• Uses of Big Data

1. Understanding & Targeting customers.
↳ understand customers & their behaviors.
↳ social media data / browse logs / text analytics.
2. Understanding & Optimizing Business Processes.
↳ optimize their stock based on predictions.
↳ predictions → social media data / web search trends.
3. Personal Qualification & Optimization.
↳ Benefit from data generated from wearing devices:
smart watches.

4. Improving Healthcare & Public Health

- { computing power of big data - decodes entire DNA strand within a minute.
- { can find new cures & predict disease patterns.

5. Sports Performance Improvement

- { video analytics - tracks performance of every player
- { sensor technology - gives feedback

6. Science & Research Improvement

7. Improving Security & Law Enforcement

- { detect & prevent cyber attacks
- { Big data tools to catch criminals
- { predict criminal activity
- { credit card companies - use big data → detect fraud

Big Data Management Cycle

Capture → Organize → Integrate → Analyze → Act

Big Data Analytics

→ process of examining data
↓
to uncover hidden patterns / unknown correlations
and other useful information

→ variety of sources / types.
→ volumes / complexities.

(3)

1. Traditional Big Data Analytics

→ enterprise will have computer to store & process data.
data stored in → RDBMS → Oracle DB / MS SQL Server.
→ sophisticated software written to interact with DB.
process data & present to users for analysis purpose.

Limitation
→ cannot deal with huge amount of data
→ can work well with only less volume of data.

Advanced Big Data Analytics

→ Map Reduce Algo - solves the problem of traditional analyt

↓
→ divides task into smaller parts.

→ assign those parts to many computers connected over network

→ collects results to form the final result dataset.

Distributing Computing

- multiple computing resources are connected in a network
- computing tasks are distributed across this network.

Sharing of tasks - increases speed & efficiency of system.

Source
Computer / Network
User

Latency Problem - ~~@@~~

- Latency - delay within a system based on delays in execution of a task.
- issue in every aspect of computing
 - communication
 - data management
 - system performance

Distributed computing can solve this problem.

Big Data companies - need low latency because of BD requires high speed & volume.

Sources of Big Structured Data

computer / Machine Generated

- Sensor data
- Web log data
- Point of Sale
- Financial data

Human Generated

- Input data
- Click-stream data
- Game related data

Sources of Big Unstructured Data

Machine Generated

- Satellite Images.
- Scientific Data
 - ↳ seismic
 - ↳ atmospheric.
- Radar / Sonar
- Photographic & Video

Human Generated

- Text internal to company - logs / surveys / results / emails.
- Social media data
- Mobile Data
- Website content.

Role of Relational Database in Big Data

Data persistence — how db retains itself when modified.

RDBMS — primitive technique for data persistence.

→ data is stored in table.

→ have a schema — structural representation of what is in the db

→ defines table, fields in table & relationships b/w two.

→ Data stored in rows & columns.

→ SQL can be used to retrieve data.

• Role of CMS IN Big Data Management

CMS - Content Management System

can manage complete lifecycle of content

web content/documents, content & other forms media

comprises - strategies / methods / tools.

capture data

manage & store, preserve

& deliver content & documents

Include technologies

collaboration

web

content

document management

web content management

collaboration

document

content

Big Data Stack

UNIT-2

① Redundant Physical Infra

fundamental to scalability & operation of BD architecture.

Supports unanticipated / unpredictable volume of data

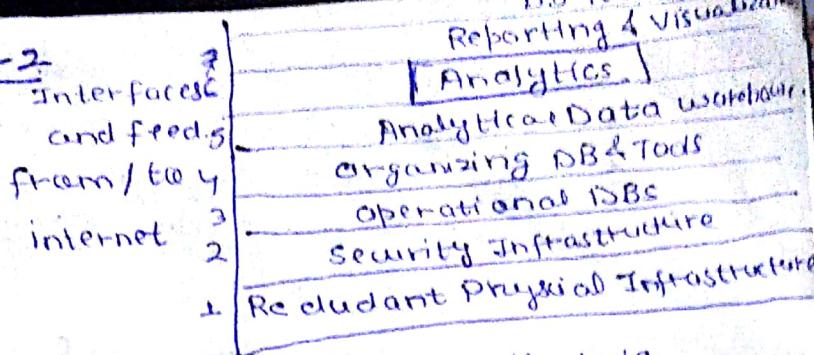
② Security Infrastructure

More important in BD analysis, more important is to source data

③ Operational Data sources

(Structured / Unstructured / Semistructured) Incorporate all the data sources to give complete picture of business

highly structured data managed in relational database.



REDUNDANT PHYSICAL INFRASTRUCTURE

Network Design

Hardware → Storage & servers
well-managed environment

① Physical Redundant Big Data Networks

Network redundant → enough capacity → accommodate anticipated volume & velocity of data
should be. & have
Designers plan - expected increase. inbound & outbound data
physical implementation - elastic.

② Manage Big Data Hardware: Storage & Servers

Hardware - sufficient speed & capacity

→ to handle all expected big data capabilities

If slow servers → bottleneck in network

fast servers → overcome variable N/w performance.

3) Big Data Infrastructure Operations Management

better performance and flexibility is only possible
in a well managed environment.

Datacenter managers → able to anticipate & prevent failures
to the integrity of data.

- Identify the type of Data needed for Big Data.
 - Exploratory stage for big data - identifying patterns in the data.
 - Codifying stage for big data - pattern identification to implementation
 - data integration & incorporation stage - after analysis.
 - └ integrate results to real-time business actions
- identify the type - explore the data → codifying → data integration
(Identify pattern) (pattern identification)

3) SECURITY INFRASTRUCTURE

- 1. Data Access - Only authorized user who have legitimate business needs should examine or interact with data.
 - └ core data storage platforms - rigorous security schemes.
- 2. Application Access - APIs - offers protection from unauthorized usage or access.
- 3. Data Encryption: - most challenging aspect of security in BD.
 - └ Identify data elements requiring more security to provide more. & fast computational capability.
 - └ encrypt only necessary items.
- 4. Threat Detection - mobile devices / social networks exponentially increases both amount of data and security threats.
 - └ Multi perimeter approach to security.

③ Operational Databases

(2)

core of big data environments → database engines having a collection of data
↳ needs to be fast/scalable /Rock Solid

A - Atomicity → all or nothing - If any part of transaction fails - entire transaction fails.
C - consistency → Transactions only with valid data will be performed.
I - Isolation → Multiple, simultaneous transactions will not interfere.
D - Durability → After the data from transaction is written to db
↳ it stays forever.

- ↳ keeps track of critical data required for real time, day to day operation.
- ↳ continuously updated based on transactions happening.
- ↳ able to scale up - to support thousands of users.
- ↳ accurate representation of business - must blend structured & unstructured data.

④ Organizing Data Services and Tools

↳ capture, validate and assemble various big data elements into relevant collections.

Organizing Data Services → Map Reduce Engines.



Designed to optimize the organization of big data.

↳ ecosystem of tools & technologies

↳ used to gather & assemble data in preparation for further processing.

↳ Tools provide

↳ integration
↳ translation
↳ normalization.

⑤ Analytical Data warehouses

Data warehouses - contain normalized data gathered from variety of sources & assembled to optimize decision making.

↳ Simplify creation of reports & utilization of disparate data items.

⑥ Big Data Analytics

Analytics Tools & techniques - helpful in making sense of big data.

↳ Able to work with large amount of potentially real time data.

Reporting & Dashboards - user friendly representation of info from various sources.

Visualization - Next step in evolution of reporting

↳ O/P - tends to be highly interactive & dynamic in nature.

↳ Difference b/w reports & visualized O/P → animation.

Analytics & Advanced Analytics

↳ Reaching into data warehouse & process the data for human consumption.

VIRTUALIZATION

- Process of using comp resources to imitate other resources.

Rather than assigning a set of dedicated physical resources to each task, a pool of virtual resources can be allocated as per need.

Increases IT resource utilization efficiency scalability Improves Latency.

Application can use a resource without any concern that where it resides how it is implemented which platform it uses. How much of it is available.

Separates resources & services from underlying physical delivery environment.

Enables to create many virtual machines within a single physical machine.

Advantages

1. Virtualization of physical resources (server/storage/N/W).
 - Improves the resource-utilization
 - control over usage & performance of IT resources.
2. Provides a level of automation & standardization.
3. Provides foundation for cloud computing.

- Characteristics - to support scalability & operation efficiency in big data environment
1. Partitioning - many VMs can be supported in single physical system by partitioning available resources.
 2. Isolation - Each VM is isolated from its host & other VMs. If one VM crashes others are not affected.
 3. Encapsulation - A VM can be represented & stored as single file. Easy identification based on services provided.

Network Virtualization - Software defined networking.

Efficient way - treat networking as a pool of connection resources.

Instead of relying on physical N/W for managing traffic b/w connections.

↓
can create multiple virtual networks.

are utilizing same physical implementation.

- creates a logical software-based view of the h/w and S/w networking devices (switches, routers etc).
- Physical N/W devices - simply forward the packet
- Virtual N/W - provides an intelligent abstraction that makes it easy to deploy & manage N/W services.

Cost saving / Efficiency / Security / Flexibility.

Processor & Memory Virtualization

Processor Virtualization [helps to optimize processor
maximize performance

Memory Virtualization - decouples memory from servers.

→ In Big data analytics - repeated queries of large dataset & creation of advanced analytic algorithms are there, to discover hidden patterns.

Advanced analytics - require lots of processing power (CPU) and memory (RAM). Without sufficient CPU & memory, these computation can take long time.

Processor & memory virtualization can help speed the processing & get analysis results sooner.

- Server Virtualization → increases efficiency / comp.
- One physical server is partitioned into many virtual servers
 - H/w & resources → RAM / CPU / hardware can be virtualized
 - Each virtual machine → has its own OS & runs its applications
 - Software representation of physical machine

Hypervisor → Technology that manages traffic b/w VM & physical machine

Server virtualization uses Hypervisor → provides efficiency in the use of physical resources

Ensures that platform can scale up as needed to handle large volume of data as we don't know the amount of data before beginning of analysis.

FOUNDATION THAT ENABLES MANY CLOUD SERVICES USE AS.

DATA SOURCE IN BIG DATA ANALYSIS.

Application Virtualization

- efficient way to manage applications in context with customer demand
- Application is encapsulated
 - in a way that removes dependencies from physical comp. system
- improves overall manageability & portability of application

Data & Storage Virtualization

Data Virtualization - can create a platform for dynamic linked data services.

allows data to be easily searched & linked through a unified reference source.

- It provides an abstract service that delivers data in a consistent form regardless of underlying physical database
- Exposes ~~data~~ cached data to all application to improve performance

Storage Virtualization

Some data may be unstructured
not easily stored using traditional methods.

- combines physical storage resources so that are more effectively shared.
- reduces cost of storage
- makes it easier to manage data for analysis.

SV makes it easier to store large & unstructured data.

DV + SV → makes it easier, less costly

to retrieve & analyze large volumes of data

Managing Virtualization with Hypervisor

Hypervisor → technology responsible for ensuring that resource sharing takes place in an orderly & repeatable way.

- ↳ allows multiple O.S to share a single host.
- ↳ creates and runs virtual machines. — practical way of getting things virtualized quickly & efficiently.
- It seems like every O.S has physical resources all to itself.
- In big data, one has to support multiple O.S — hypervisor becomes load delivery mechanism for it.
- It lets you share same application on multiple systems without having to physically copy that application on each system.

Each virtual machine running on a physical machine → Guest Machine

Hypervisor → schedules the access to guest O.S to have everything including CPU, memory, disk I/O etc.

↳ It does all the heavy lifting

↓
Guest O.S ~~does~~ not have any idea that it's running in virtual ~~host~~ partition.

Type 1 Hypervisor → runs directly on hardware platform.

↳ achieves higher efficiency — running directly on platform

Type 2 Hypervisor → runs on host operating system

↳ used - when need exists to support broad range of I/O devices.