## What is a data stream?

→ A stream is defined as a possibly unbounded sequence of data items or records. That may or may not be related to, or correlated with each other.

e.g instruments, many IOT appn areas, computer programs, websites or social media posts.

→ data each datas are timestamped and in some cases geo-tagged.

• Streaming data → sometimes referred to us event data → each data item is treated as an indv. event in a sync sequence. Synchronized sequence of events.

~~Data stream~~

### Streaming Data Systems

{
• Manage one record or small time window
• Near-real-time
• Independent computations
• Non-interactive.
}  Quiz

Dynamic steering ← part of streaming data management & processing.
→ self driving cabs.
~~data streamin~~

- Some streaming Data ~~stream systems~~ systems
  • amazon kinesis.
  • Apache storm
  • Flink
  • spark Streaming
  • Samze

## Why is cleaning Data different?

### Data -at-rest.
- mostly static data from one or more sources.
- collected prior to analysis.

### Data-in-motion
- analyzed as it is generated.
- stream processing ← analysis of stream data.

~~Data Processing ← analysis of~~

### Data Processing Algorithms:

static/Batch    size determines,
Processing →    time and space

Streaming
Processing → { unbounded size,
but ~~infinite~~ finite time
Quiz        & space.

### Streaming Data Management & Processing:
• compute one data element or a small window of data elements at a time.

ACID & BASE

is difficult to maintain in a BDMS.

↳ notares ACID
  ↳ BA: Basic Availability
  ↳ S: Soft state
  ↳ E: E-ventual Consistency

CAP Theory: A distributed computer system cannot simultaneously achieve.

  • Consistency
  • Availability
  • Partition Tolerance

The

——————×——————

Week. 4 notes

| Schema on Write | Schema on read |
|---|---|
| (1) data is ~~stored in~~ structured & enforced to adhere to a specific schema before being written to database/ datawarehouse | (i) data is stored without any predefined structure or schema. |
| | write vice versa |
| (2) ensures data integrity & consistency | (ii) schema is applied at the time of reading |
| | (iii) allows more flexibility in handling diverse & unstructured data. |
| (3) used in datawarehouses | (iv) used in data lakes |

Asterisk DB ← apache semi
└ used to deal with semi-structured data.
  └ fully-fledged, provide ACID properties like data integrity
  └ structured to deal with data that doesn't fit in rows &
    columns. like JSON

working
  └ organizes data into structures called dataverse & data types
                                                    ↑
                        act as namespace for data

                        define structure of
                        data.

e.g  tweets → nested parts like            } Asterix DB employs
             entities & user information     hierarchical structure
                                             in its schema.

└ handle geospatial data too, geographical info
└ uses a query language, XML similar to XLL
  called AQL.
       └ enables querying multiple languages like XQuery,
         Hadoop, Mapreduce, etc.

└ operates on clusters of machines.
└ divides data into partitions & executes queries

→ These computations can update metric monitor, & plots statistics

your query

○ Relatively fast & simple computation

○ No interaction with data sources.
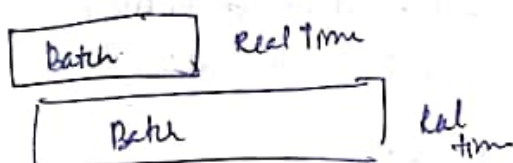
## Hybrid architecture

↳ lamda architecture - for processing streaming & back Jobs at the same time.

→ In these systems streaming data elements are pushed to a batch system and become available to access & process as batch data.

In such systems, stream storage layer is used to enable fast trees of streams & ensure data ordering & consistency.

## λ architecture

| Batch | Real Time |

| Batch | Real time |

streaming data changes over time.

size + frequency.

| size → unbounded |
| size & freq → Unpredictable |
| Processing → Fast & simple |

Quiz

Changes can be periodic or sporadic

_____

Periodic : evenings, weekends, etc.

Sporadic : major events

Quiz

other changes include dropping or missing data or even no data
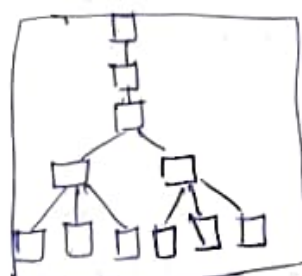
---

**Data Lake** ← big data storage & processing challenges.

→ part of data structure that many streams can flow into & get stored for processing in their original form

Quiz

### Data Warehouse vs Data Lake.

↓                          ↓
Hierarchical file          object storage
system

↓                          ↓
Structure format           raw format
↳ when data is in          ↳ data is
use, then stored in        stored as it
warehouse                  gets streamed.

→ Scheme-on-write    → Schema-on-read.

### Data Lake Object Storage

(i) Each data is stored in binary large object (BLOB) & is assigned a unique identifier.

(ii) Each data object is tagged with a # of of metadata tags.
                    ↑
            data is search using
            the tags.

### How a Data Lake Works

→ Load data from source
→ Store raw data
→ Add data model on read

Programming Models on Big Data
→ an abstraction or existing machinery or infrastructure.
→ set of run time libraries + programming libraries
→)

requirements: (i) Support big data operations. → split data values ~~values~~   volumes of data
                                              → Access data fast
                                                dis tr computation to node

(ii) Handle fault Tolerance → Replicate data partition
                            → recover files when needed.

(iii) Enable adding more racks

(iv) Optimize for specific datatypes

## Hadoop Goals
(i) Enable scalability
(ii) Handle (fault) tolerance
(iii) Ability to handle diff types of data
(iv) facilitate a chare env.
(v) Provides value ┌ community supported
                   └ wide range of Application

## Hadoop Ecosystem
(i) HDFS ┌→ scalability to large data sets
         └→ reliability to cope hardware failure

2 components
(a) Name nodes → keep track of filename, location in directory ,etc.
    → mapping of content in datanode
(b) Data nodes → stores file blocks
    → listens to namenodes for block evaluation,
      deletion & replication.
    → replication is done for fault tolerance & data locality.
YARN: flexible scheduling & resource management over HDFS.
Hive & Pig → A simplifies parallel computing
    − You only need to given MapReduce two function
      Map → applyes
      Reduce → summarize ()

Cloud Services ⟹ Service Models ⎡ Application
                              ⎢ Platform
                              ⎣ Infrastructure

Iaas = Get the Hardware only , Amazon EC2 cloud.
        ⮡ bare min^m rental service
Paas = - Platform as a service , google App engine
        - Get the computing env.  Microsoft Azure
Saas = Get full s/w on demand
        - service model
        - dropbox.

Decision depends on ⎡ skill
                    ⎢ demand
                    ⎢ capital
                    ⎣ security

Xaas = Anything as a services.
       ⎡ storage aaS
       ⎢ Marketting aaS
       ⎣ Communication aaS.

Hadoop layer diagram:
Low level iff → slow storage & scheduling, on the bottom
High level language & interactivity at the top.

SQL-like queries
facebook.

Yahoo
data flow scripting

Hive & Pig: augment data modeling of MapReduce with relational algebra & data flow modeling respectively.

Giraph: out large scale graph efficiently.

Storm, Spark, Flink: used for real time & in-memory processing of big data.

Hbase, cessandra, MongodB = No SQL for non-files

Zookeeper: management like synchronization, configuration & high-availability.

indexing {

Map Reduce: relies on YARN to schedule & execute parallel processing over distributed file system intents.
Map → Apply operation
Reduce → synchronize operation.
Map → Shuffle & Sort → Reduce

⭐ Where to use hardoop                    Where not to
→ Many platform over single data sou      → small dataset
→ High volume                             → Advenced volgo
→ High variety                            → Task parallelism
.. Random Access, Infrastructure replacement

Cloud Computing & Cloud service
(i) Build Resource            (ii) Clouds

## Steps in data Science Process

1. Acquire data
   - Identify/determine data set when to look
   - Retrieve data Suitable add
   - Query data → Relation db, R1

2. Prepare data
   - → Explore (Preliminary analysis, Understand nature of data)
   - → Pre-process data - (Clean, Integrate, Pkg).

3. Analyze Data
   - → Select analytical techniques
   - → Build model.

4. Communicate Results

5. Apply Results.

Acquiring data.
   - identify suitable data
   - acquire all available data.

Webservices
{ REST → Representational State Transfer
  SOAP, Websocket

NOSQL storage -
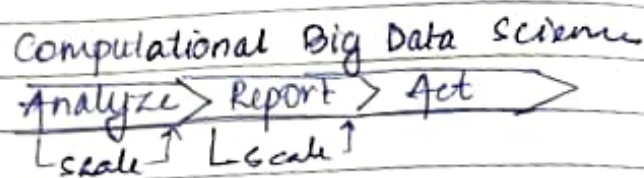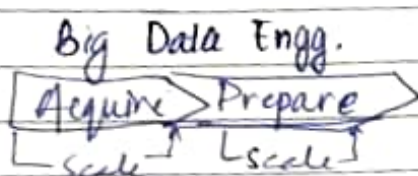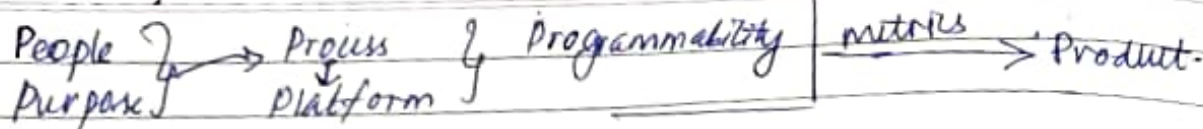| mongoDB, cassandra, neo4j, Apache, couch DB |

provide API allows user to access data → API — using webservice.
Webservice ← using → REST

eg → Acquiring Data from WILDFIRE

Historical Weather ← SQL
Curr. Weather ← Websocket
Real-time tweets ← REST
near fires

Trad$^n$ dB → SQL & query browser
Remote data → Webservices
Text files → Scripting Languages
No SQL tech/storage → Web Services, Programming Interface

5 P's
  Components of Data Science

People ⟩ → Process ⟩ Programmability | metrics → Product.
Purpose ⟩    Platform ⟩

Big Data Engg.              Computational Big Data Science

Acquire ⟩ Prepare ⟩         Analyze ⟩ Report ⟩ Act ⟩
└ Scale ┘ └ Scale ┘         └ Scale ┘ └ Scale ┘

Process, Build metrics for accountability
Cost
Timeline
Planning of deliverable
Expectations
Purpose

Asking the right question:

Assess situation → Risks, Benefits, Contingencies,
       regulations, resources, requirements
Define Goals → Objectives
            → Criteria

Steps to find right problem to tackle in data science

Define Problem
     ↓
Assess Situation      } formulate the Question
     ↓

Define Goals

→ Explore

Goal: To Understand Data

Why explore?    sales prices going up or down.    to detect errors

- explore dependencies of ~~variable~~ b/w diff variables of data

correlation    General trends    Outliers
↳ data pt distant from other data pt.

Describe your data
　Mean & Median ← represent the ~~distance~~ location of a set of values.
　Mode ← value that occurs most frequently
　range & standard deviation → measure of spread in your data

Visualize your data - Histograms
　　　　　　　　　- Line Graphs
　　　　　　　　　- Heatmaps
　　　　　　　　　- Scatter Plots
　　　　　　　　　- Box plots

Step 2B - Prep-processing data

　clean + transform
data quality issues — inconsistent values
　　　　　　　　　- Duplicate records
　　　　　　　　　- Missing values
　　　　　　　　　- Invalid data
　　　　　　　　　- Outliers

Addressing Data Quality Issues.

Remove data with missing values
Merge duplicate records

Generate best estimate for invalid value
Remove Outliers

Domain Knowledge

**Valence:** Measure of Connectivity.                          are related to each other

→ Data Connectivity → two data items are conn. when they
• Valence — fraction of data items that are connected out of total
# of possible connections.

Note: Valence increases over time
         Makes the data connections denser

Challenges: • More complex data exploration algo.
           • Modeling & prediction of valence changes.        ?
           • Group event detection.
           • Emergent behaviour analysis.

Strategy ⟸ Big data ⎡ Aim
                    ⎢ Policy
                    ⎢ Plan
                    ⎣ Action

1st step in determining a big data strategy ⟹ Business objectives ← Long term
                                                              ← Short term

Provide organizational buy-in
   • commitment , • sponsorship , • communication.
Build diverse teams: • Diverse expertise , • deliver as a team.
   Share data: • Remove barriers to data access
              • No data silos ← big heap or tall tower.
   ✳        • Data sharing mindset
Define big data policies: Privacy & Lifetime ( whom should be given access )
                         Curation & quality. data curation & quality
                         Interoperability & regulation.
Cultivate analytics-driven culture:
Taking all the decisions using analysis of data ⎰ Analysis
                                                 ⎱    + ⟹ Opportunities
                                                    Business

   Integrate analytics → Comm. goal → build team
                              ↑            initiate
                         Adopt for ⟳ share data
                         new direction

Velocity: Speed $\dfrac{\Delta x}{\Delta t}$

speed of creating, storing & str analyzing data.
- Real time processing ← gathering weather info for travel
                        ← sensors saving lives.

Batch Processing.

Colled Data → Clean Data → Feed in chunks → Wait → Acts } incomplete

Real-Time Processing

Instantly capture Streaming → Feed real time to machines → Prouss Real Time → Act } fast

Rate needed for data-driven actions
↓

Rate of gen' & processing of data.

Streaming data =         +        Streaming data >
"what's going on                   gets generated at a varied
right now"                         rates

| Real-time processing |

Agile & adaptable business decisions

Veraity: Quality - validity, volatility
- Accuracy of data
- Reliability of the data source
- Context within analyses.

# Week-2

Characteristics of Big Data
- **Volume** — vast amt of data that is generated every second. etc.
- **Variety** — diff forms & can come in. That data —
- **Velocity** — pace at which data moves from one pt to the next
- **Veracity** — refers to biases, noise & abnormality in data
- **Value** — connucheaness of data

4 v's

**Volume:**

Challenges:
Storage $\xrightarrow{\text{distr}}$ Processing

volume↑ | performance↓ | cost↑

Volume = size $\rightarrow$ Challenges
- storage
- Access
- Processing

we use qualitative vs quantitative measur
age can be a # or we can prep as juvenile, infants

**Variety?**

Axes of Data Variety

Semantic Variety
how to interpret & operate on data.

diff in the representation of the data → Structured Variety — formats & models

satellite images in wikipedia very diff from tweets sent by ppl.

EkG signal to newspaper article

medium in which data gets delivered ← { Media Variety — medium in which data get delivered.

eg - audio of a speech & transcripts.

Availability variations -
real-time ? ← traffic amms
Intermittent ? ← satellite data

Scalability Issue: Impact of data variety:
- harder to ingest, • difficult to create common storage.
- difficult compare & match data across variety, • diff. to integrate
- management & policy challenges.

**File System** → long term info storage system
→ It can access recent of process data
→ store large amt of info
→ enable access of multiple process

## Distributed File System (DFS):

→ replicated the data b/w the racks & also computer across the geographical location.

→ DFS makes the system more built tolerant,

## High concurrency vs Low consistency

Data Partitioning

Data replication
↓

Data scalability

Fault tolerance

High concurrency

## Scalable computing over l/n.

Commodity Clusters= are affordable parallel computer with avg # of computing nodes conn. to each other via fast n/w.
→ class specialized.
→ reduce computing cost.

Distributed Computing → computing in one or more of these cluster across a local area n/w on Internet.
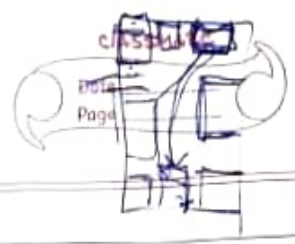enables data parallelism.

## Common failure in commodity cluster:

- failure of entire rack.
- failure of connection b/w rack & n/w.
- failure of connection b/w two nodes in the a rack.

Failure ↛ Complete Restart

Redundant
data storage

Data parallel
Job restart

# Getting Data in Shape
also called Data Munging, Data Wrangling, Data Preprocessing.

when dataset has large # of dimensions

**Data Munging :**
- Dimensionality Reduction (3D to 2D)
- Data Manipulation
- Transformation ← reduce noise & variability.
- Scaling — scale values b/w zero & one
- Feature Selection — remove, combine, add feature
  - redurdant or irrevalant features
  - feature

raw data, manipulated to be in correct format for analysis

comparing heights & weights, weight in mnns, your ise. scaling will equalize contributions

## Categories of Analysis Techniques:

Predicted vs correct values

(i) Classification ← Predict category

(ii) Regression → Predict numeric values e.g sales, stock price

(iii) Clustering ← Organize similar items into groups.

(iv) Association Analysis ← Find rules to capture associations b/w items.

investigate & validate

(v) Graph Analytics : Use graph structures to find connections b/w entities.

## (Modeling)
Select Techniques → Build Model → Validate Models

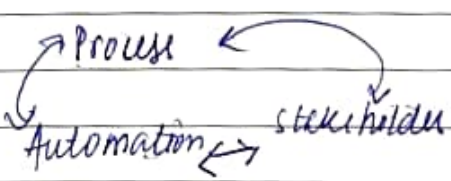| Communication | present using visualization tools.
R, Python, D3, Leaflet, Tableau, Google Charts, Timeline

create visualizations in your public profile
cross platform compatibility
allows to create timeline.

**Action :**
Process ←
Automation ↔ stakeholder

**Assess Impact :** Monitor
Measure
Evaluation