

Data Essentials

* Data, Information, Knowledge, And Wisdom

↳ Raw and unstructured facts

↳ contextualized, organized trend or pattern

↳ Ability to solve a problem

↳ Decision making and judgments

* Sources of data

Experimental data - clinical trials

Observational data - recorded through observations

compiled data - Manual data aggregated overtime and managed to database.

Simulated data - result of studying certain behaviors

Human Generated data - generated through mobilephones, web, audio, social media

Machine Generated data - Digital Devices, IoT

Business Generated data - stock exchange, Banks etc.

* Common Types of data

2 main types

* Quantitative data - measurable such as numbers and values, conclusive

* Qualitative data - Not measurable; descriptive. (unstructured data)

* Geospatial data - Earth data, GPS system

* Digital data - images, Audio, video, web based

* Data and from documentation and scripts - HTML, open text document.

* Data formats

Structured data - Data having a defined format, model, structured

Easy to search and analyze EX: CSV, XLS

Semi-structured data - Data having flexible schema

schema maybe descriptive, partial or evolving EX: XML, JSON, Gmail

Unstructured data - No structure, format or model EX: Pdf, Jpeg, mp4

* Data Terminologies

Primary data - collected by a person or group by themselves.

Secondary data - collected by a third party

Data Analytics - Exploring Data

Data science - uses scientific methods tools to extract knowledge and insights from the data.

RANKA

DATE / /

PAGE

* Database - Repository of data

Data Analytics - process/technology used to explore data

Data Aggregation - procedures to gather data

Metadata - details about data

Time Series - analysis studying temporal patterns.

* Anonymized Data - Individuals identities are kept secret

Ex: Feedback form

Augmented Reality - Data analytics utilizing ML and NLP

Data literacy - reading and communicating data as information

File format - XLS, CSV

Data science -

Deep learning - Neural network based ML Algorithms

AI - The ability of machines to demonstrate human like behavior

Data Storage and Backup

Full Backup

Backup for all data

requires high space

slowest speed

Lots of duplicate data

Differential Backup

Backup of revised of last backup

medium to high storage

Fast speed

Intermediate amount of duplicate data

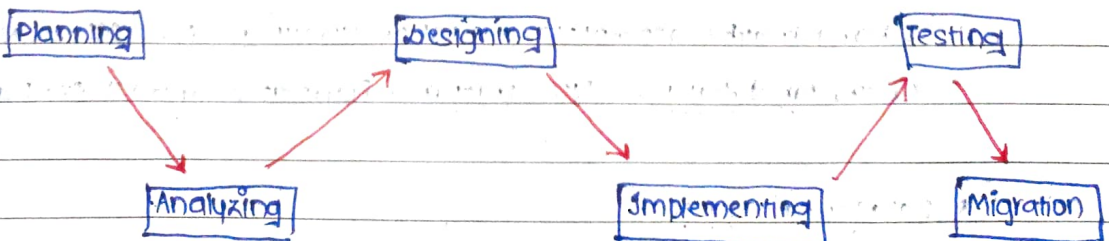
Incremental Backup

only revisions or changes to previous small storage

fastest speed

little or no duplicates

Data Migration : Transfer of data from one storage system to another



Planning

Identifying Stakeholders

Identify data

Risk, Mitigation and Backup strategy

Analyzing

Explore the data

determine impact, migration

cut-off point

Create a data dictionary

Design

Develop "source to target" mapping

Determine ETL Transformations

Implementation

configure tools

Write migration scripts

Testing

Deploy all required tools

create test plan with accurate data coverage

design a migration Validation engine

Migration

Execute migration steps

Execute
 - proper planning
 - go live including migration

ETL

Extract

retrieves and verifies
data

Transform

processes and
organizes extracted data

Load

moved transformed
data to repository

Data Visualization - representing data in the form of charts, graphs and maps.