

U6.

Data Visualization

- graphical representation of information & data using visual elements like charts, graphs & maps.
- helps users identify complex dataset quickly by identifying patterns, trends & outliers.

Importance

- makes large & complex data easier to understand.
- common tools could be used.
- used in BI, data analytics & decision making
- makes insights of the data clear & easy

Dis

- may cause misinterpretation
- complex visuals → overwhelm users.
- limitations of visualization tools.

Challenges in Big Data Visualization.

- ① Volume - visualizing immense datasets is difficult due to sheer size & complexity
- ② Variety:- combining variety of data source into a visual is challenging.
- ③ Velocity:- Real time data streaming requires fast processing & updating visuals.
- ④ Visual Noise:- high density of data points makes hard to distinguish individual obj's
- ⑤ Information loss:- data reduction has risk of info loss.
- ⑥ Large image perception:- Device constraints limits how data can be effectively displayed.
- ⑦ High rate of change:- Rapid data update can overwhelm the user
- ⑧ High performance requirement.

⇒ Solutions ⇒ ① Meeting the Need for speed ② understanding the data ③ Addressing the data Quality ④ only meaningful result ⑤ Dealing with outliers.

SPPU-TE-COMP-CONTENT - KSKA Git

Applications of Data Visualization

(1) Business sales Analysis. (bar chart)

- Track sales performance by region or product.

(2) Website traffic monitoring. (line chart)

(3) Customer feed back & Rating. (pie chart, histogram)

(4) Outlier Detection in finance & healthcare. (Box plot)

Types of data Visualization

(1) line chart → shows trends over time

(2) Bar chart → compare quantity across category

(3) Pie chart → % or proportional data.

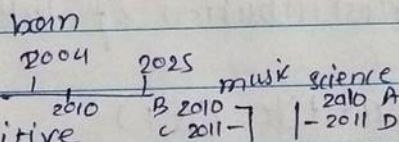
(4) Scatter plot → show relation b/w 2 variables.

(5) Histogram → distribution of numerical data.

(6) area chart

(7) flow chart

(8) timeline



(1) line chart.

- created using plot() func.

- show rel of x & y

x = [10, 20, 30, 40]

y = [20, 25, 30, 35]

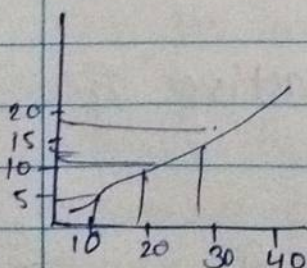
plt.plot(x, y)

plt.title("ABC")

plt.xlabel("X")

plt.ylabel("Y")

plt.show()



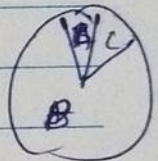
(2) pie chart.

cats = ['A', 'B', 'C']

data = [10, 20, 10]

plt.pie(data, labels=cats)

plt.show()

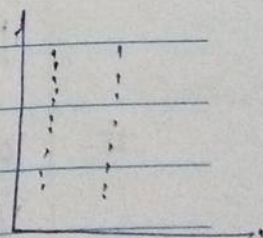


(3) x = data['day']

y = data['total']

plt.scatter(x, y)

plt.show()



(4) plt.hist(data, bins=10, color='skyblue')

data = np.random.randn(1000)

SPPU-TE-COMP-CONTENT - KSKA Git

⇒ # Tools used in data visualization.

(5). Density plot. - descuat

- smoothed version of histogram
- use to visualize distribution of data & identify peaks, spread & skewness.
- uses ~~kernel~~ kernel density estimation
- good for comparing multiple distributions.
- curve based.

eg. `data = np.random.normal(loc=0, scale=1, size=1000) // random data`

`sns.kdeplot(data, shade=True)`

`plt.show()`

- not fixed interval.

(6). Box plot.

- shows distribution of dataset. through 5 summary stats.

1. min.

2. first Quartile (Q1)

3. median (Q2)

4. Third Quartile (Q3)

5. max.

- shows outliers if present.

eg. `data = np.random.normal(0, 1, 1000)`

`sns.boxplot(data=data, color="skyblue")`

`plt.show()`

⇒ # Tools used in Data Visualization / Analytical tools in Big data visualization.

- analyzing massive datasets
- transform raw data into interactive & understandable visuals.

⇒ Tools are,

- ① Tableau
- ② Power BI
- ③ Qlik
- ④ Google data studio.

DS, JS
plotly
kibana.

pentaho
datameer.

①. Tableau

- powerful tool for interactive data visualization
- connects to big data platform like hadoop, spark, AWS & SQL
- has drag & drop interface
- also used for BI

②. Po- allows user to analyze, visualise & share data insights through interactive & shareable dashboard

- converts raw data to understandable format.

Key features

- ①. User friendly UI - No. coding required
- drag & drop functionality
- ②. Connects to multiple data source - SQL, Excel, google sheets.
- ③. Real time data Analysis. - supports live data connections & real time dashboard updates.
- ④. Interactive Dashboard. - visualization can be filtered, drilled down, & updated dynamically
- ⑤. Advanced chart types - scatterplot, histogram, geo plots.
- ⑥. Data blending & joining. - data from multiple sources.
- ⑦. Storytelling - creates a sequence of visualization to tell a data story

⑧. Tableau (Public / Desktop / server)

Public - free.

paid (full featured)

to share dashboard

Adv.

- easy to learn, highly interactive & dynamic.
- supports large datasets.

Limitations.

- high cost, limited custom visualization.
- Requires clean data.

Use case :- Analyze monthly sales, product wise performance.

SPPU-TE-COMP-CONTENT - KSKA Git

② Power BI

- Business analytics tool developed by Microsoft
- enables users to visualize & share insights across organizations
- helps create interactive dashboard, reports & data models.
- widely used in business & academic env. for data-driven decision making

Key features.

- ①. Multiple data source connectivity.
- ②. Interactive Visualization
- ③. Power Query editor - to clean, Transform & shape the data
- ④. DAX (Data analysis expressions)
 - a formula lang. used to make custom calc.
- ⑤. Dashboards & reports.
 - dashboard: shows real time KPI's
 - reports, multi paged visuals / charts.
- ⑥. Data modeling:-
 - Relationship b/w tables can be created & analysed
 - supports star & snowflake schema.
- ⑦. cloud & on premise Deployment.
- ⑧. AI ~~capabilities~~ compatibility

adv.

- user friendly
- Integrated with MS. tools.
- cloud based sharing & collaboration
- Real time data analysis

Limitation - slow with large dataset

DAX complex for beginners.

limited customization for advanced visualisations.

③ Qlik

- data analytics & BI platform
- allows users to create interactive visualizations, dashboards & data apps.
- they have 2 main types/products.
 1. Qlik View - script-based dashboard dev. tool
 2. Qlik Sense - self service & cloud friendly BI platform.
- Qlik is known for its in-Memory data processing & associative data model, - allows users to explore & analyze data without being restricted by predefined queries.

key features.

- ①. Associative Data model:- unlike SQL based tools, Qlik links data from multiple sources in non-linear associative way.
- ②. In memory processing. Data is loaded in memory → fast processing & calculations.
- ③. Data loading & ETL:- powerful scripting engine. allows data transform, cleaning & loading from diverse source
- ④. Rich visualization capabilities.
- ⑤. Self service BI → no heavy IT involvement
- ⑥. AI-powered insights.
- ⑦. Collaboration & sharing. via Qlik cloud & servers.

adv.

- fast & intuitive data discovery
- easy exploration.
- strong data compression & performance
- scalable (small, large) business.
- can integrate with R, python, Rest APIs

Limitations.

- less intuitive
- hard to learn.
- paid licence

SPPU-TE-COMP-CONTENT - KSKA Git

- supports data intensive distributed applications.

Hadoop Ecosystem.

- platform or framework to solve Big data problem.
 - Ecosystem refers to the various components in software library of hadoop (accessories & tools)
 - Hadoop is Java based framework for handling & analysing large set of data.
 - use of different parts of core hadoop set such as MapReduce, to handle vast amount of data. HDFS, a file handling system, YARN - hadoop resource manager.
- Tools of hadoop ecosystem

④. HDFS, Hive, Pig, YARN, MapReduce, Spark, HBase etc.

Ecosystem.

coordination	workflow & scheduling.	scripting (Pig)	ML (Mahout)	Query (Hive)	NO SQL HBase	Data integration
			(MapReduce)			
			Distributed processing			
			Distributed storage (HDFS)			

① HDFS: Hadoop distributed file system

- primary storage sys. of hadoop
- uses NameNode & DataNode architecture.
- it is a distributed file system.
- NameNode (Master) & DataNode (acts as slave)

②. Hive - ETL & Data warehouse tool

- used to query or analyze datasets.
- 3 main func.

①. Data summarization.

② Query

③ analysis of unstructured & semi-structured data.

③. Map Reduce

- core component of processing (provides logic of processing)
- software framework helps in writing applications.

SPPU-TE-COMP-CONTENT - KSKA Git

that process large datasets.

- distributed & parallel flow of tasks.

(4) Apache Pig.

- high level scripting lang.
- to execute queries for large dataset.

(5) Apache Spark

- fast in memory data processing engine
- Supports, Java, py, scala, R & SQL

(6) Apache HBase

- hadoop eco-sys. component. for distributed database
- designed to store structured data
- scalable, distributed, NoSQL database.
- build on top of HDFS.

Hadoop architecture.

consist of (1) Name Node

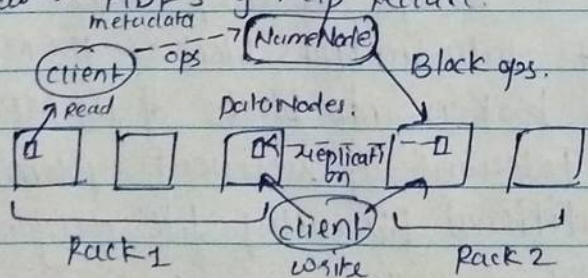
(2) Data Node

(3) Secondary Name Node

(4) HDFS client

(5) Block structure

- to process & store vast amount of data hadoop mostly uses. HDFS & Map Reduce.



HDFS is a block structured file system where each file is divided into blocks of pre-determined size

- blocks are stored across clusters
- HDFS follows Master slave architecture i.e. NameNode (Master) &

SPPU-TE-COMP-CONTENT - KSKA Git

- data Node (slave) → they process & store blocks.
- NameNodes - maintain meta data of dataNodes. & continues client access.
 - Info like Namespace info, block info
 - this info only brought to memory when some activity is requested else its stored in disk/persistent storage.
 - fsimage - logs of NameNode, (snapshot)
 - editlog - changes to read & write logs. (when namenode starts)
 - To restart a NameNode we need fsimage & to initialize & editlog. (to update the logs of NameNode - Time consuming)
 - If the NameNode shuts down the fsimage is used to start it at the same position as it was using the image it stores.
 - & the edit log updates the NameNode by updating its logs. it requires a lot of time.
 - to reduce that time Secondary NameNode updates the editlog to nameNode in a fixed interval.
 - Reducing the time for restart.
 - HDFS & MapReduce is must without that nothing can happen.
- ## # HDFS.

- Hadoop distributed file system
- used to scale a single hadoop cluster to 100's of nodes.
- min amount of data in HDFS block is 128MB
- these files are broken into blocks of 128MB size
- HDFS is fault tolerant & resilient system.
- support traditional hierarchy file organization
- directory based & follows namespace.
- follows Master Slave architecture using NameNodes & dataNodes.
- can be deployed on machines that operate on Java

SPPU-TE-COMP-CONTENT - KSKA Git

MapReduce

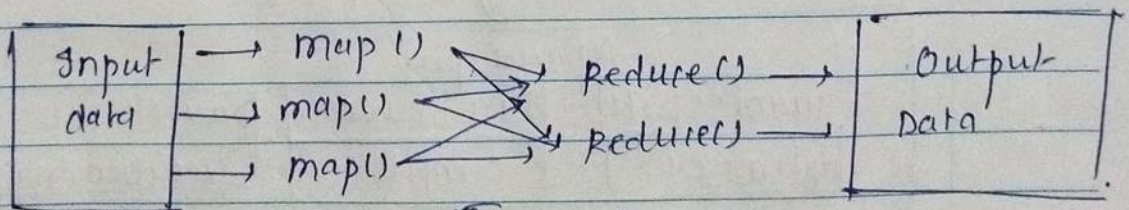
- is a programming model & software framework
- Intends to simplify the process of vast amount of data in parallel on large clusters of commodity hardware in a reliable fault-tolerant manner.

characteristics.

- works on large scale data
- write once read many times (data).
- allows parallelism without mutexes.
- map & reduce. main operations.

Map() & Reduce()

- performs actions like filtering, grouping & sorting
- reduce() func. aggregates the summarise the result
- result generated by map is (key, value) pair which acts as input to Reduce func.
- Every Map/Reduce program must specify a Mapper & a Reducer. mapper has map method that transform's input into (key, val), Reducer has reduce method. the transform (key, val) argument into any no. of output. (key, val) pairs



flow.

(Optional)

Input → split → map → combine → shuffle → Reduce → o/p
 ↓ ↓
 split data to process
 small pieces each split
 according to map() func.

SPPU-TE-COMP-CONTENT - KSKA Git

How MapReduce fits in Hadoop ecosys.

- (1) Data Processing Engine.
- (2) Working with HDFS.
- (3) distributed processing
- (4) Data locality optimization.
- (5) fault tolerance
- (6) Job scheduling & resource management
- (7) Scalability (horizontally) (among hardware nodes)

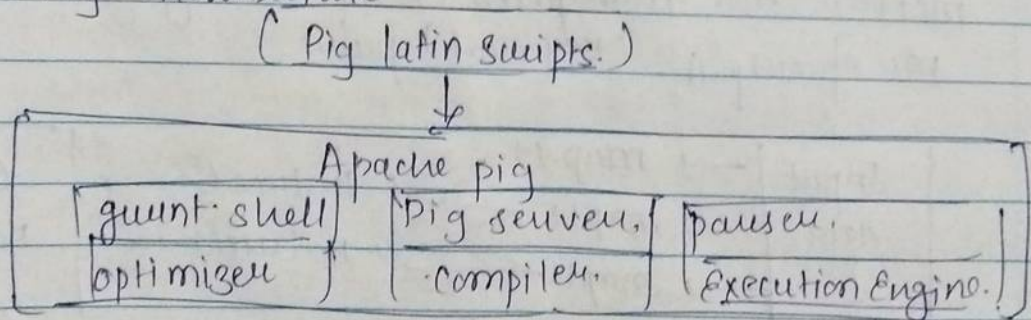
Apache Pig.

- high level platform for creating mapReduce programs. used with hadoop.
- uses scripting language called pig latin.
- simplifies coding & required for processing large dataset.

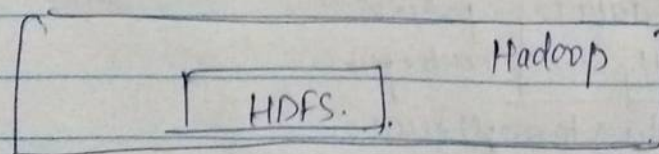
Features

- handles structured & ^{semi} unstructured data.
- provides data flow language
- supports user defined func
- used for ETL (Extract, Transform, load)

Pig Architecture



(Map Reduce)



SPPU-TE-COMP-CONTENT - KSKA Git

- ① Pig latin script. → user writes script to define data operations
- ② Grunt shell → interactive command line interface (CLI) where user write pig latin script.
- ③ Pig driver → acts as interface b/w user & Pig runtime env. managing the execution of pig scripts
- ④ Parser → parses the script & checks syntax producing a logical plan
- ⑤ Optimizer → optimizes the logical plan by applying optimization rules to make data processing more efficient
- ⑥ compiler :- converts the optimized plan to a series of mapreduce jobs.
- ⑦ Execution Engine :- runs compiled MapReduce jobs on hadoop cluster, processing the data stored in HDFS.
- ⑧ MapReduce :- Executes the tasks generated by pig
- ⑨ HDFS - stores input & output.

- Pig acts as interface b/w user & hadoop
- Reduces complexity of writing raw MapReduce code.
- Supports batch & interactive mode.

HIVE (Apache)

- data warehouse infrastructure build on top of hadoop.
- provides an SQL like interface (HiveQL) to query & manage large dataset stored in HDFS.
- for users who are comfortable with SQL & not with Java

Key features.

- ①. HiveQL (Query lang) :- similar to SQL, easy data analysis
- ②. Schema on read. - allows flexibility
- ③. Batch processing
- ④. Build on top of hadoop
- ⑤ Supports user defined functions.

SPPU-TE-COMP-CONTENT - KSKA Git

Hive architecture

- (1). UI (CLI/Web UI) - interface for queries
- (2). Driver - manage lifecycle of HiveQL statement
- includes parser, planner, optimizer, executor.
- (3). Compiler - converts query into logical execution plan
- (4). Metastore - stores metadata like table schema.
- usually stored in RDBMS
- (5). Execution Engine - interact with hadoop to run job & return result
- (6). HDFS :-
stores data.

adv.

- SQL like lang
- easy to use.
- work on masive dataset
- Scalable with HDFS & MapReduce
- supports partitioning & Bucketing.

limitations.

- No real time processing
- high latency.
- limited transactional support.
- No row level operations.

use cases

1. data summarization. & reporting
2. log analysis in batch mode
3. ETL