

Histogram :

→ A histogram is a type of bar chart that represents distribution of numerical data by grouping values into ranges called bins.

→ Width of bar is decided by the numeric proportions of grouped data.
Height is decided by the freq.

→ Eg: Consider the data of test scores:

14, 10, 12, 19, 18, 20, 10, 8, 5, 0, 2, 9

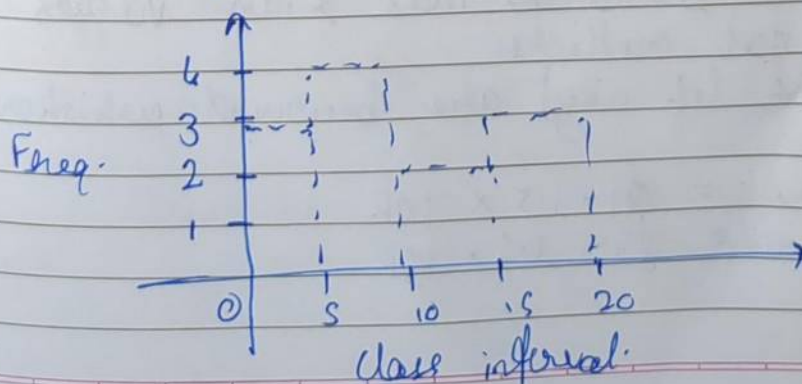
① Step ①: In ascending order

[0, 2, 5, 8, 9, 10, 10, 12, 14, 18, 19, 20]

② Step ②: Div. into bins, decide the range.

Class Interval	0-5	5-10	10-15	15-20
Freq.	3	4	2	3

③ ~~Plot~~ Step ③: Plot





→ Usage :

- ① Data distribution: Understand how data pts. spread over an interval.
- ② Detect skewness: Identify whether data is skewed or symmetric
- ③ Outliers: Bins w/ very low/high freq. or isolated may indicate outliers.

Box plot:

→ A box plot is a graphical repⁿ of data that show central tendency, spread and potential outliers.

→ Summary stats. viz. in box plot are:

- Minimum
- First quartile (Q_1) 25%
- Median / Second quartile (Q_2) 50%
- Third quartile (Q_3) 75%
- Maximum

→ Characteristics:

- Box shows Interquartile Range ($IQR = Q_3 - Q_1$)
- Vertical line in box shows Median (Q_2)
- Whiskers extend to max & min. values that are not outliers.
- Outliers (if any) are beyond whiskers.

→ Lower limit = $Q_1 - 1.5 \times IQR$
Upper limit = $Q_3 + 1.5 \times IQR$

→ Eg.:

SPPU-TE-COMP-CONTENT - KSKA Git

Consider data = 100, 120, 110, 150, 110, 140, 130, 170, 120, 220, 140, 110

① Ascending order :

[100, 110, 110, 110, 120, 120, 130, 140, 140, 150, 170, 220]

② Q1, Q2, Q3 :

$$Q2 / \text{median} = \begin{matrix} \text{odd} \\ \cancel{n}^{\text{th}} \end{matrix} \left(\frac{n}{2} \right)^{\text{th}} \cdot \text{OR} \begin{matrix} \text{even} \\ \left(\frac{n + (n+1)}{2} \right)^{\text{th}} \end{matrix}$$

$$\text{Median} = \frac{120 + 130}{2} = 125$$

Q1 = First 6 vals ka median. If odd, exclude Q2 / Median

Q3 = Last 6 vals —, —

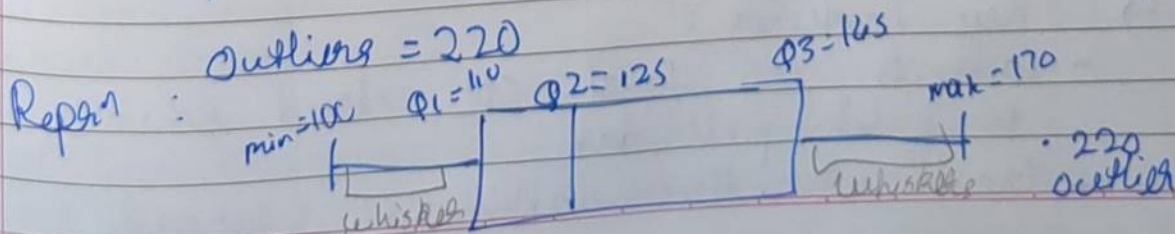
$$Q1 = \frac{110 + 110}{2} = 110 \quad Q3 = \frac{140 + 150}{2} = 145$$

$$IQR = Q3 - Q1 = 145 - 110 = 35$$

$$\text{lower limit} = Q1 - 1.5 \times IQR = 57.5$$

$$\text{upper limit} = Q3 + 1.5 \times IQR = 197.5$$

∴ min, max betⁿ [57.5, 197.5] are 100, 170 resp.





↳ Usage :

- Identify outliers
- Summarise large datasets
- Compare distributions of diff. groups.
- Skewness of data (if median not centered)
- Whiskers help identify spread of data.

Diff. of histogram from density plot :

- | | |
|---|---|
| ↳ A histogram repr ⁿ data freq. of bins whereas density plot repr ⁿ probability density funct ⁿ . | |
| ↳ Uses bars (hist) | ↳ Uses continuous curve (density) |
| ↳ No smoothing; depends on bin width | ↳ Smoothed repr ⁿ |
| ↳ Area of bars repr ⁿ count | ↳ Area under the curve = 1 (probab ^l) |
| ↳ outliers less visible | ↳ outliers easily identified. |

MS PowerBI :

- ↳ - It is a powerful BI & data viz. tool by MS.
 - Enables users to connect multiple data sources, transform raw data, create interactive dashboards & reports.
 - Easy to use & simple GUI
 - Cloud based solⁿ
 - Fast & efficient
- ↳ Key features :
- Connects w/ Excel, SQL, cloud services, APIs to fetch/pull data.
 - Create interactive dashboards & reports

SPPU-TE-COMP-CONTENT - KSKA Git

- AI integration for AI-driven insights and natural lang. queries.
- Cloud based so on any machine & easy to share.

→ Use cases:

- Business performance tracking
- Sales analysis
- Financial reporting
- Operational monitoring

Qlik

- - Data analytics & BI platform known for its associative data model.
- Allows users to explore data w/o being limited by predefined queries.

→ Key features:

- Fast in-memory processing.
- Connects w/ various data sources.
- Users can create their own dashboards & reports
- Enables free form data exploration.

→ Use cases:

- Market & customer analysis
- Healthcare analytics
- Supply chain optimization
- Risk mgmt.

QlikView → Script-based & developer-driven

Qlik Sense → more modern, user-friendly, focused on self-service analytics

Apache Pig: Acts as bridge betⁿ user & Hadoop

- Component of Apache Hadoop ecosystem.
- High level scripting language.
- Used for creating MapReduce programs used w/ Hadoop.
- Uses scripting lang. called Pig Latin that simplifies coding req. for processing large datasets.
- Supports batch & interactive modes.

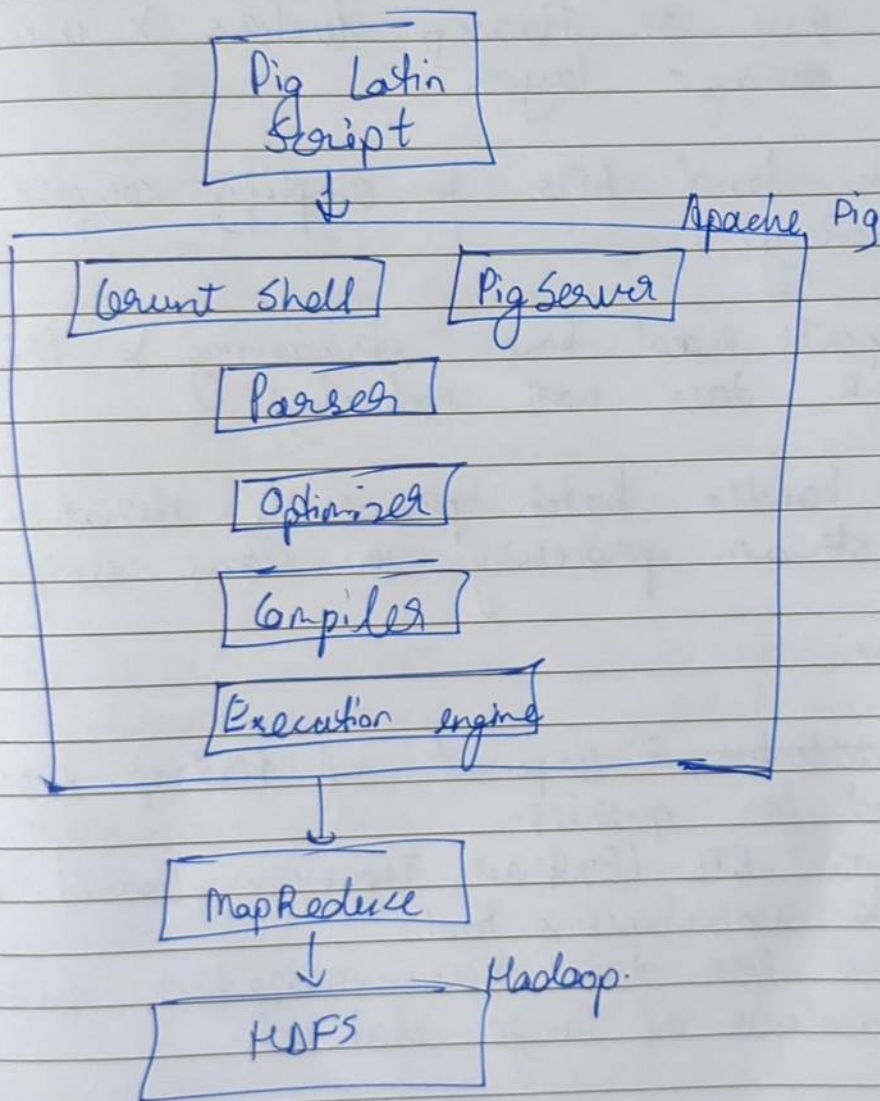
Architecture:

- ① Pig Latin Scripts - Used for writing data analytics logic.
- ② Grunt Shell - Interactive CLI for users to write & execute Pig Latin Scripts.
- ③ Pig server - Acts as interface betⁿ user & Pig runtime env., managing execution of Pig scripts.
- ④ Parser - Takes Pig Latin script & converts into logical plan. I/P \rightarrow O/P
- ⑤ Optimizer - Improves logical plan by applying optimization rules to make data processing more efficient. O/P \rightarrow optimize logical plan (L/P)
- ⑥ Compiler - Takes I/P \rightarrow translates into series of MapReduce Jobs.
- ⑦ Execution engine - Runs compiled MapReduce

SPPU-TE-COMP-CONTENT - KSKA Git

jobs on Hadoop cluster, processing data stored in HDFS.

② MapReduce executes data processing tasks. HDFS stores I/P & O/P data.



Role of Pig in Apache Hadoop:

- Pig scripts converted to MapReduce Jobs auto by Pig engine.
- Improves productivity of developer by abstracting complex mapR. logic.
- Useful for iterative data processing & prototyping.

Spark complementing Hadoop:

- Can run workloads upto 100x faster than traditional MapR job due to in-mem. data processing engine.
- Can run on Hadoop clusters & use HDFS as storage layer.
- High level APIs to simplify complex data processing.
- Supports real time streaming & ML which MapR does not natively.
- Can handle batch processing, iterative algos & stream processing in one unified engine.

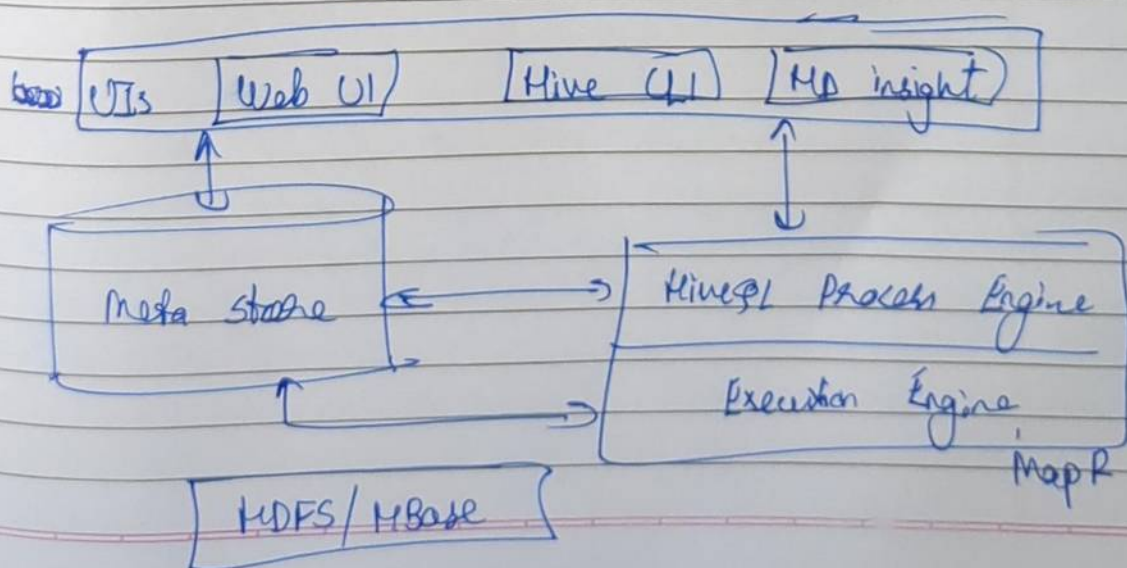
Hive

- ~~Architecture~~ - Component of Hadoop ecosystem
Used for queries.
- It is ETL (Extract, Transform, Load) & data warehousing tool.
- Used for data summarization, querying & analysis of large datasets.
- Architecture:

① UI - Hive has a CLI for writing & running HiveQL queries & a web UI to interact w/ Hive.

SPPU-TE-COMP-CONTENT - KSKA Git

- ② Meta store - stores metadata about Hive tables (Schema, data types, table location in HDFS).
Managed using traditional RDBMS like MySQL.
- ③ HiveQL Process Engine - Accepts & processes HiveQL (SQL-like) queries. Parses & compiles & optimizes the queries. Converts HiveQL queries to execution plans.
- ④ Execution engine - Works w/ process engine & MapReduce. Executes query plans by launching corresponding jobs & handles job monitoring & collects final result.
- ⑤ MapReduce - Default execution framework for Hive. HiveQL queries are internally converted to MapR programs. Enables parallel processing.
- ⑥ HDFS / HBase - HDFS is default storage system for Hive. HBase is used when real-time read/write access is needed. Stores data in tables in distributed format.



→ Characteristics:

- Has SQL like lang. (HiveQL) thus easy for users familiar w/ SQL.
- Schema is applied when reading data.
- Batch processing for large scale offline analysis.
- Metadata stored in RDBMS-based metastore.

→ Features:

- Scalable
 - Compatible w/ Hadoop
 - Supports summarization, querying & reporting
 - Supports partitioning & bucketing to improve performance of data access.
-