

UNIT-II

Types of digital data

Digital data

- repr other forms of data using specific language, then it can be interpreted by various technologies.
- Eg: binary system.

It stores audio, video, doc & text info in the form of binary characters (011)

Digital data storage → hard drive

- Mostly offline storage [back up]
- Use cloud storage. (online)
- Server, contains all uploaded data.
- data stored in the form of code.

Different ways to store data

① RAM → Temporary data storage.
→ allows to read data fast.

② hard disk → Permanent data storage.
→ It can be external / Internal.

③ Cloud → we can store our data in cloud. we can access data from cloud anywhere and anytime.

④ Digital data center

- Save space, money & time.
- 24/7 security.
- quickly access the data.

⑤ Digital asset management

- Servers.
- Stores data in diff location.
- well suited for files & docs.
- Easily share the data.
- Organizes files efficiently.
- Digital data storage should not help storage space.
- A chance to quickly retrieve files.

Data

explanations

- set of facts, such as descriptions,
Observations & numbers.

- Three types.

(1) Structured data

(2) un " " " "

(3) semi " " " "

Structured data

- Tabular data.

- repr by rows & columns.

- called relational database.

↳ formed set of columns.

- SQL used for structured data.

Eg
→

Student-id	name	age
1	Jim	20
2	Sam	19
3	Ram	18

sub_id	subject	Teach
001	Maths	Joh
002	English	Vich
003	Science	Bal

stu-id	sub-id	grade
1	001	90
1	002	100
1	003	92
1	001	95
2	002	80
2	;	:

Semi- Structured ^g data

- It does not consist of structured data
 - but still some structure on it.
 - Contains both structured & unstructured information.
 - Common in object-oriented database
 - Does not follow any tabular data structure

eg

Email	XML
CSV	JSON → Javascript object notation.

Unstructured data

- Does not have any pre-defined data model.
 - data does not organized in pre-defined manner.
 - e.g. video, audio, binary data, Image, audio, sensor data.
 - Collection of various data types stored in native formats.

eg

Unstructured data

The university has 5600 students. John's ID is 1, he is 18 years old and holds B.Sc degree. David's id is 2, he is 32 years old and holds a degree Ph.D. Robert's id is 3, he is 22 years old, he holds BE degree.

Semi-structured data

<university>

<stud-ID="1">

<NAME>John</NAME>

<AGE>18</AGE>

<DEGREE>B.Sc</DEGREE>

</stud>

<stud-ID="2">

<NAME>David</NAME>

<AGE>32</AGE>

<DEGREE>Ph.D</DEGREE>

<stud-ID="3">

</stud>

<stud-ID="3">

<NAME>Robert</NAME>

:

Structured data

ID	Name	age	Degree
1	John	18	B.Sc
2	David	32	Ph.D
3	Robert	22	B.E

Characteristics of Structured

- Relational → structured
- Non-relational = unstructured

Data quality

Def

The degree to which data meets a company's expectations of accuracy, validity, completeness and consistency.

→ follows some formats.

→ DOB → follow specific format

↓
entered valid.

→ If don't, not valid.

Completeness → contains all relevant info

Consistency → same info stored in multiple sources, matches equally.

Data quality characteristics / dimensions

- Accuracy.
- completeness
- consistency
- Integrity → Accuracy + Completeness + Consistency.
- ~~Responsible~~ Reasonability → Whether data meets expected range, type & value
- Timeliness
- Deduplication. → speed of dissemination
- Validity.
- Accessibility. → spreading
eliminate redundancy
Make our data easy to get.

Data quality challenges

① Privacy and protection laws

- protection of private data which is accessible by people. (Individual)
 - ↑ → provides some guidelines for General data protection regulation (GDPR)
 - & California Consumer Privacy Act (CCPA)
 - organizations must store individual info securely and should not miss even a fraction of the collected data.

② AI & ML

- Most of the companies uses AI & ML applying to their BI.
- Difficult to keep up the data secure.
- Chance to ^{happen} mistakes & inaccuracies.
- Monitoring complicated tasks are more challenging.

③ Data governance practices

→ DBMS

- It ensures all data is consistent, validity & trustworthy.
- Should provide right data governance.
- If not, Confusions may occur.

6 ways to determine data quality

(same)

characteristics.



Data classification

(1) Data sensitivity levels:

↳ High sensitivity data

Medium

Low

High → ~~Unauthorized transactions.~~

→ Include personal info.

→ " Medical / health info

e.g. Credit card / debit card / Military id no

→ Confidential.

Medium sensitivity data

- Internal business info & financial info.
- emails, documents

Low sensitivity data

- Public info
- e.g. Public website content

② Based on performance.

① Content-based classification

- Classify doc and files → based on review

② Content based classification

- based on meta data

③ User based classification

- based on manual judgement

Data states

- There are 4 states

* Rest

* Process

* Transit

* Confidential: Takes less time to classify data

Data format

2 types
structured → DB
unstructured → email, video
↳ take more time to classify the data.

Data discovery

- collect data from various sources.

- used to understand trends & patterns of data.

Data stored in

- DB
 - cloud → googledoc
 - Big data platforms.
 - files → PDF, email
- Before data classification, we must perform accurate data dictionary.

Create own classification policy

- Create our own classification policies.

- We should consider:

- Which person/org/ → owns the info.
- Who is responsible for accuracy of info.
- Where the info stored.

structured

unstructured

1) displayed in rows/columns

cannot display data in table format.

2) Nos, dates, strings

Images, video, audio, emails

3) hold 20% enterprise data

80% enterprise data

4) Requires less storage

Requires more storage

5) Easy to manage

More difficult to manage.

Data warehouse 3

- RDBMS
- designed for query and analysis data.
- It derived data from various sources.
- It cannot perform transaction.
- Mainly used for making decisions.

DW attributes

- DB, investigating data from various apps.
- Support relatively small no. of clients.
- Contains current & historical data.
- Tables.

Characteristics of DW

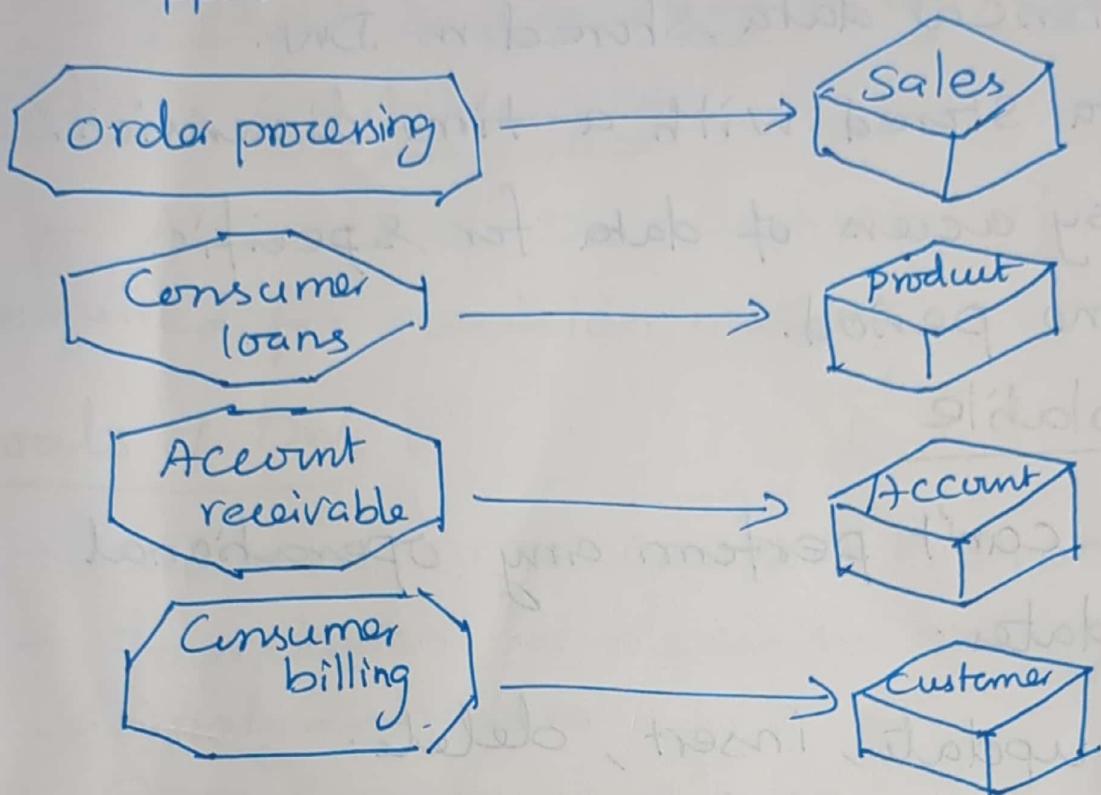
- (1) Subject-oriented.
- (2) Integrated
- (3) Time-Variant
- (4) Non-Volatile.

Subject-oriented

- The data is organized in specific subject.
- So, we can access and track the data easily.

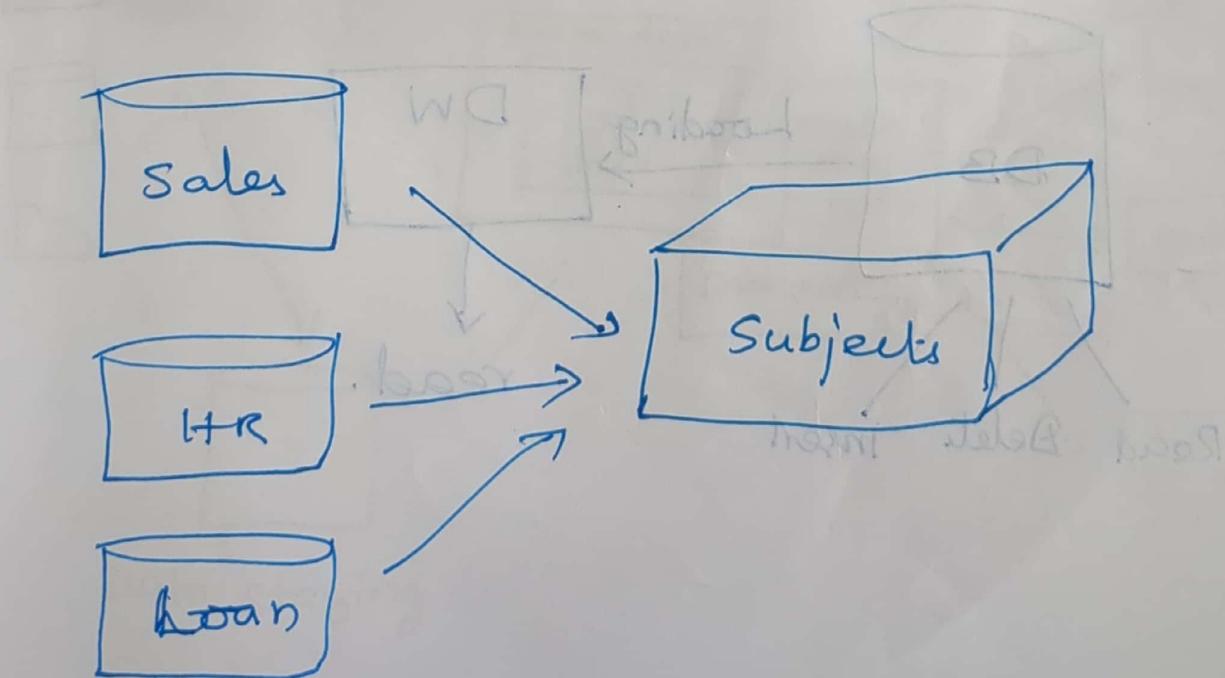
Operational Applns

Dw subjects



Integrated

- DW integrates various data sources like RDBMS, flat files & online transaction records
- Done data cleaning.

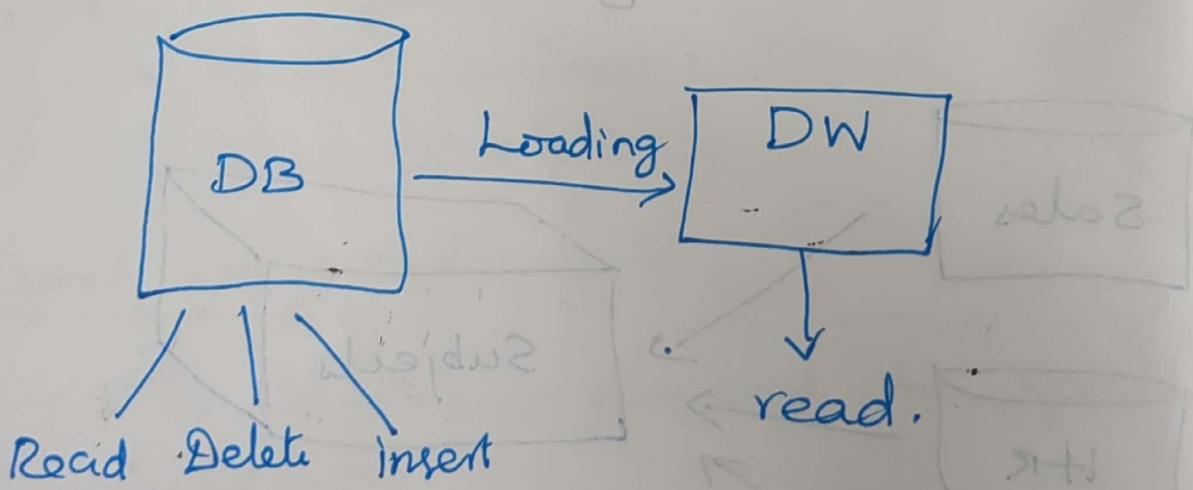


Time Variant

- Historical data stored in DW.
- Data stored with a time dimension.
- Easy access of data for specific time period.

Non-volatile

- We can't perform any operational updates.
- 2 fns < loading of data
Accessing of data.
- Once the data entered in to the Datawarehouse, data should not change.



History of DW

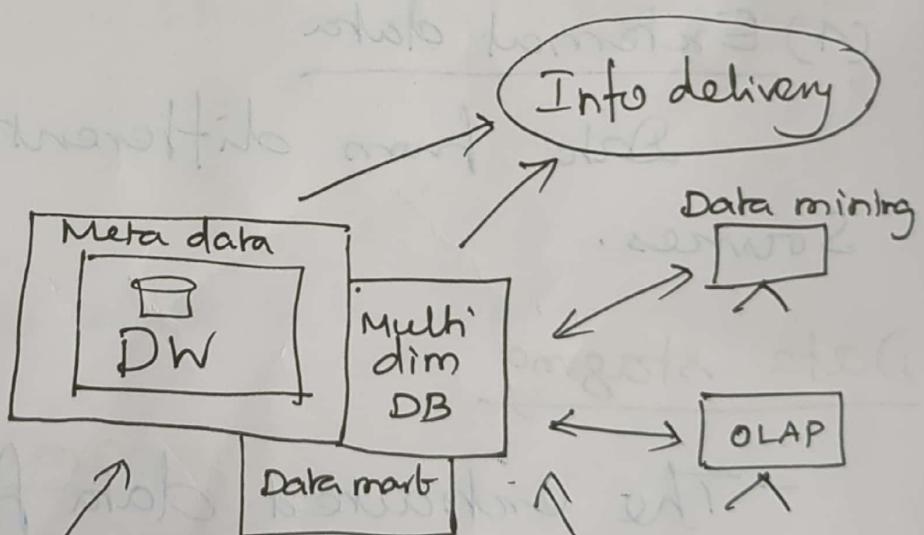
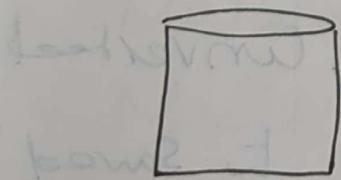
- Established in 1980's by Barry Devlin & Paul Murphy
- Absence of DW, vast amt of space required for decision making.

Goals of DW

- Maintains organisation's historical info
- foundation for decision making.
- Analyzing data.
- Data consistency & quality.
- High response time.

Components of DW (Building Blocks)

Sources
data external



Data staging

Source data component

4 categories

(1) Production data

- data comes from the different OS of the enterprise.

(2) Internal data

- Client keeps their private info like spreadsheet, reports, profile, dept info.

(3) Archived data

- Periodically take the old data and store in archived files.

(4) External data

Data from different external sources.

Data staging

- The extracted data from various sources need to be changed, converted and ready in a format \rightarrow to saved in DW.

2 processes

(1) Data Extraction

- appropriate techniques

(2) Data transformation

- Cleaning the data

- Error data

- unwanted "

→ remove

- Standardization

- Combining pieces of

data from various sources.

- Transformation

- Storing → merging of

data.

Data loading

Loading the data into DW.

Data storage Components

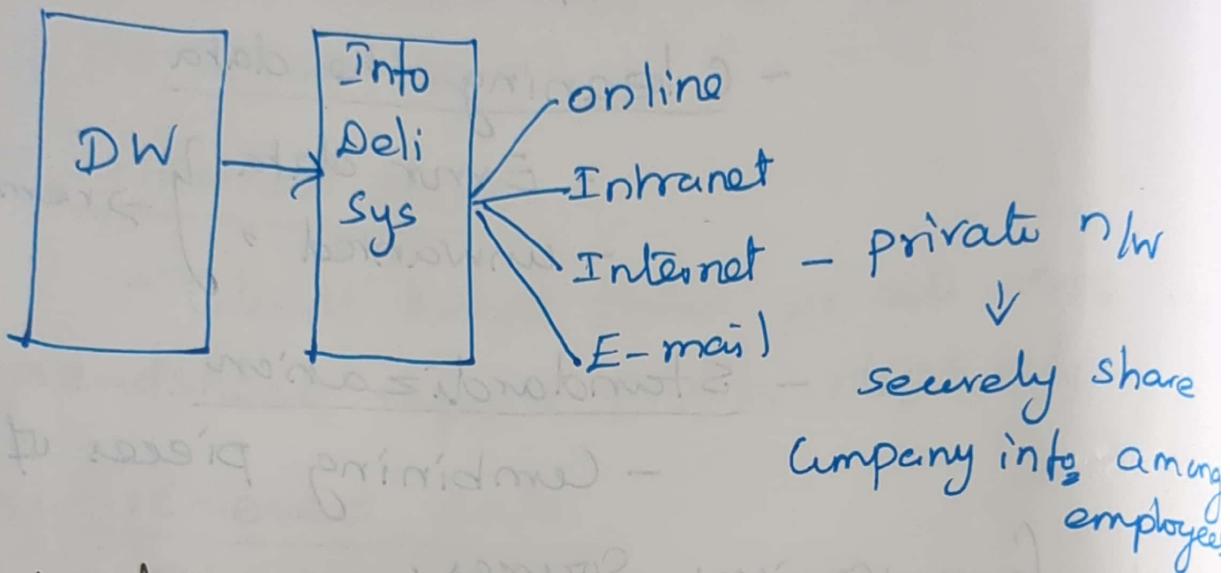
- Metadata

- Multidimensional DB.

- Data mart.

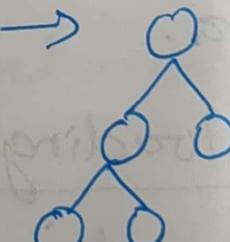
Information Delivery Component

- process of transferring DW files to one or more destinations.



Meta data

- Data about data.
- * logical structure of data →
 - records
 - addresses.



Data mart

- Contains particular topic subject.
- Several data marts available in organization.

Management

- Coordinate the services & functions in DW.
- Control data transformation.
- ensure the data correctly saved in repositories.

Diff b/w DB & DW

Database	Data warehouse
① Used for OL+P	Used for OLAP.
② Tables & joins are complicated. Response time - $\frac{\text{Max}}{\text{Slow}}$	Response time \rightarrow Min $\frac{\text{Seconds}}{\text{Min}}$
③ Data is <u>dynamic</u> \downarrow changes made on data.	Data is largely <u>static</u> \downarrow fixed data
④ Read operation, Write, update, Modify	Write Read operation
⑤ Performance is low	High performance.

Data lake

- When we want to storing big data.
2 options are available.

(1) Data lake

(2) DW.

(1) Types of data

- Structured data

- un " =

- semi " =

Purpose

- cost effective storage of large

amt of data from various sources.

- Allow any structure of data.

- Doesn't need to fit a specific schema

Users

- Data analysts

- Business "

- Data Engineers → Maintains data

- Data Scientist

Task performed by Data lake

- Data engineers use data lakes to store incoming data.
- Perform big data analytics.
↳ Using Apache Spark & Hadoop

Diff b/w data lake & DW

	Data lake	DW
Type of data	unstructured & structured data	Historical data.
Purpose	Cost effective big data storage	Analytics for business purpose.
Users	Data scientist's & engineers	Data analyst & business analyst.
Task	- Storing large amt of data. - analyze big data. - deep learning, real-time analytics	Aggregating & summarizing data.
Size	Stores all data &	Only stores relevant data.

Business reporting

- Shows business related info like facts, figures & analysis.
- Used to make decisions & plan for the future.

Reporting process involves

- (1) Compiling relevant data.
- (2) Reviewing info with specific areas.
- (3) Analysis.
- (4) Drawing conclusions.
- (5) Recommendations

Types

(1) Analytical reports

- Used for taking decisions on a business pbm.

(2) Informational reports

- Gives the info abt no. of employees, in which dept they belong, role of employee
- Show data in diff ways graph, chart, table.

(3) Research report

- Plan to launch new product in the market
- Plan to expand office area.

Benefits

(1) Adhere to regulation & Competence

- Obey the laws, regulation, guidelines and specifications of the business organisation.

(2) Transparency

- Provide clear picture for everyone involved in the organization.

(3) Improving efficiency

- Business Managers & owners need to share their experience to their employees in all the way of taking decision & improving efficiency.

(4) Troubleshooting

- Identify pbm in time.
- Escalating pbms to higher levels.

(5) Customer focus

- Improve the organization on every level and provide better customer experience.

Visual Analytics

Data visualization

- Showing data in visual format.
- Makes easier to understand the concepts.
- In the form of
Chart, graph, list, maps,
dashboards, score card.

Visual analytics

It does the following jobs

- * - Identify & reveal patterns and trends.
- It prepares data for data visualization
 - * Examine data
 - * Align algorithms with requirements
 - * Identify the patterns
 - * Gain meaningful insights from complex data sets.

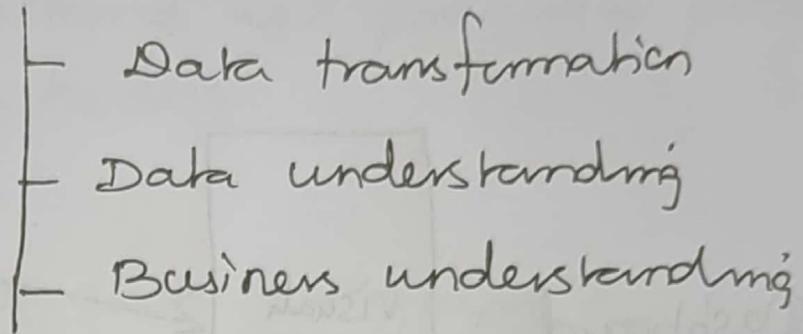
Role of visualization in Analytics

- Data visualization is either Static or Interactive

Static → Single view of dataset.

Interactive → Various views of same dataset.

Interactive visualization



Types of Analytics

3 types.

(1) Descriptive → simplest type.

- process of using current and historical data to identify trends.

ex → / Extract

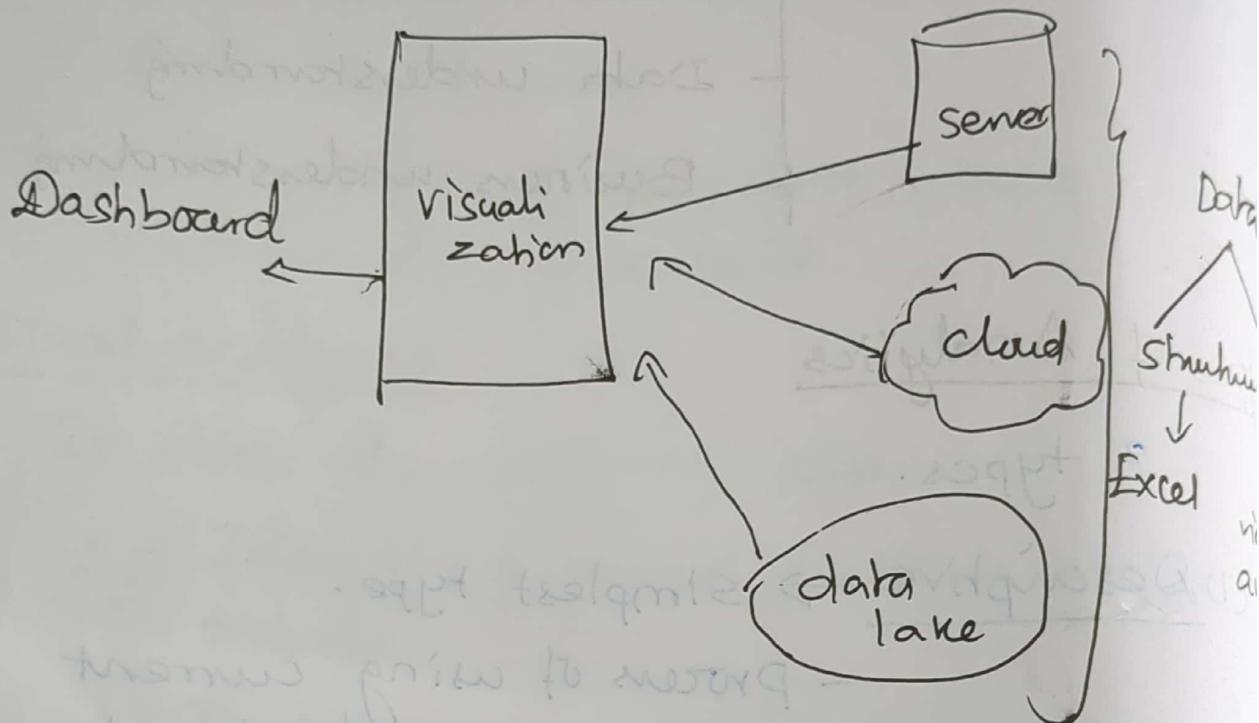
(2) Prescriptive analytics

- yields recommendations from customers for next step.
- ML algos are used here.

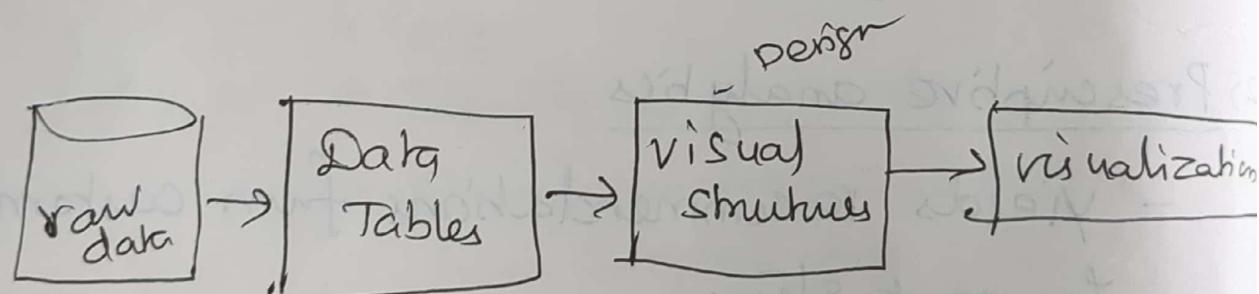
(3) Predictive analytics

- Complex
- analyze
 - * forthcoming scenarios
 - * New trends
 - * Make them efficient & effective

Process



Different types of charts & graphs



Types

* Chart

bar chart
pie chart

* Table

* Graphs

* Map

* Dashboard

Uses of chart & graph

- Easy to understand abt particular process.
- Impress Stake holders.
- Build trust

How to choose right chart or graph?

- ① Identify the goal for presenting the data.
- ② Gather data.
- ③ Based on select graph / chart

Comparing values

* Column chart

* Bar

* Pie

* Line

* Scatter plot

easily, show low & high values in data set.

Composition values

* Pie , * Area

* Waterfall ,

+ Mekko

Distributing data

Scatter, pie, bar, Line

Analyzing trends → Line, column

Relationship b/w data sets → Line, Bubble, Scatter plot

① Bar graph

3

- Helps to compare data b/w diff groups.

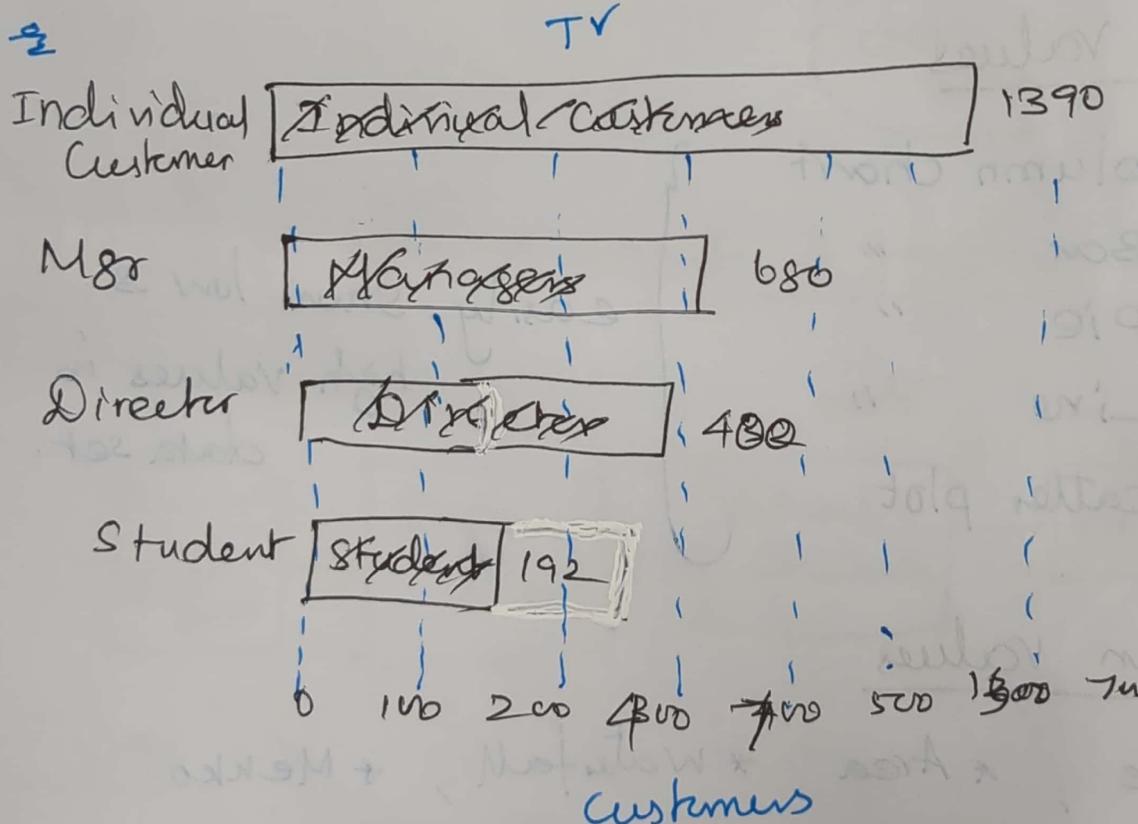
* Product Comparisons

* Price

* Category

- Attributes

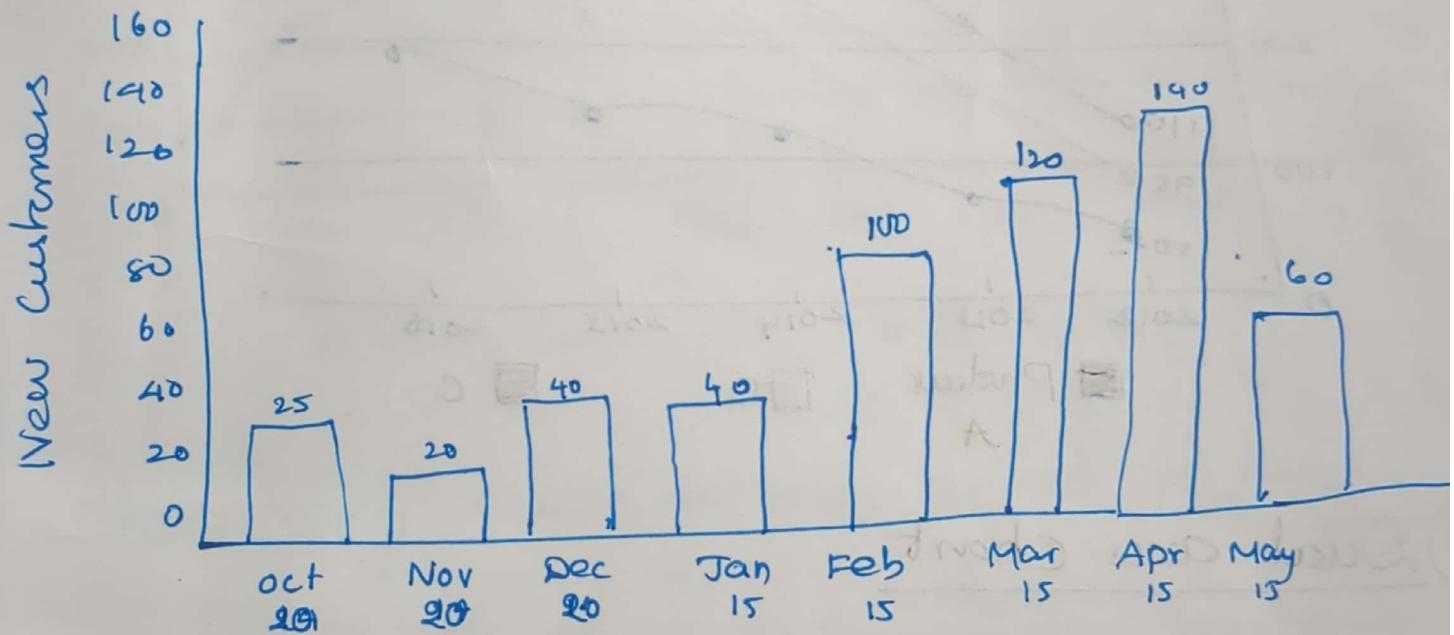
Color, Label, axis



② Column Chart

- Comparison among diff items.
- Show info vertically.
- * customer survey data
- * sales
- * profit & loss

attributes → color, Labels, axis

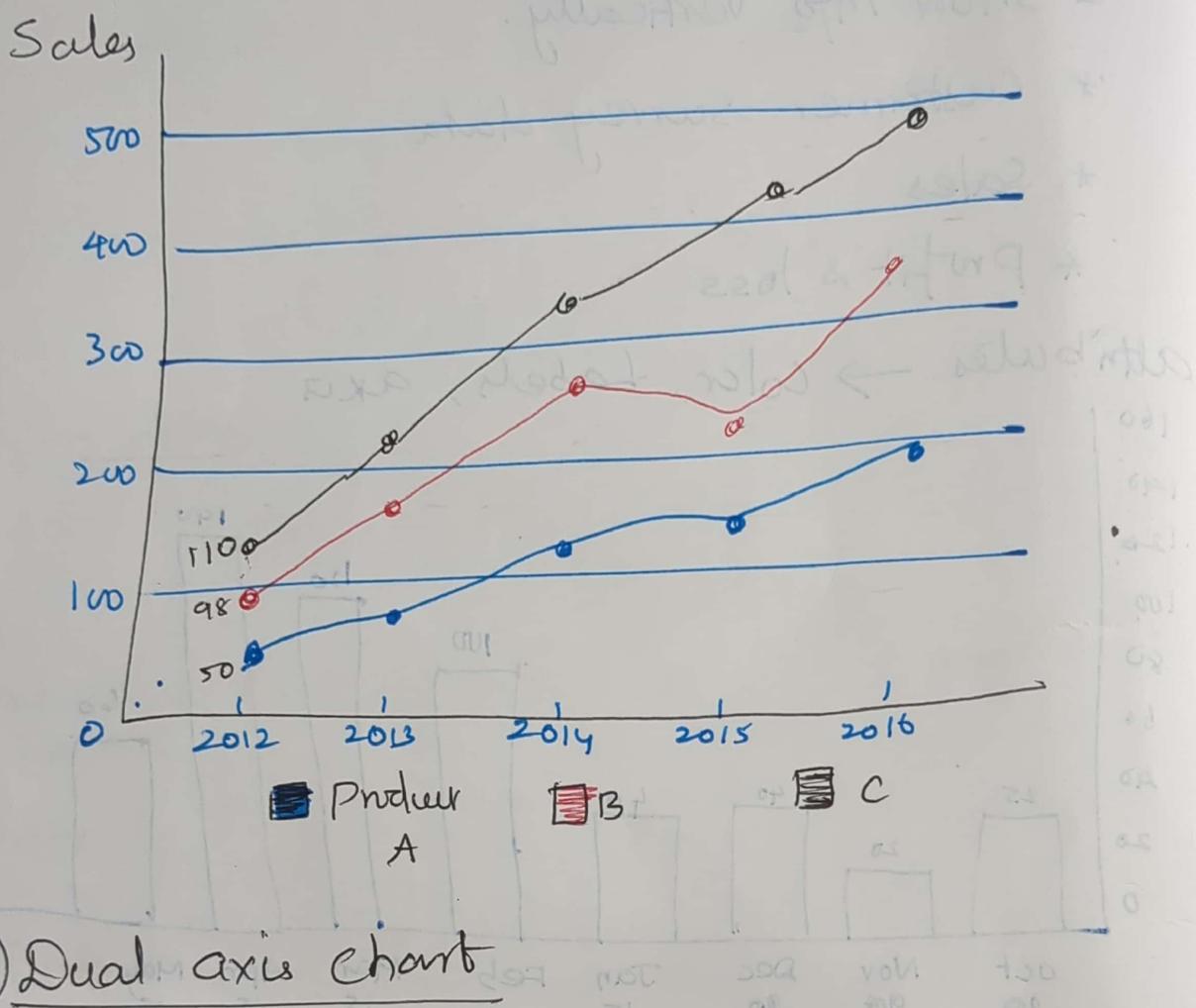


③ Line graph

- Comparison of diff categories of data
- Track changes over short & long periods of time.
- Good for seeing small changes.

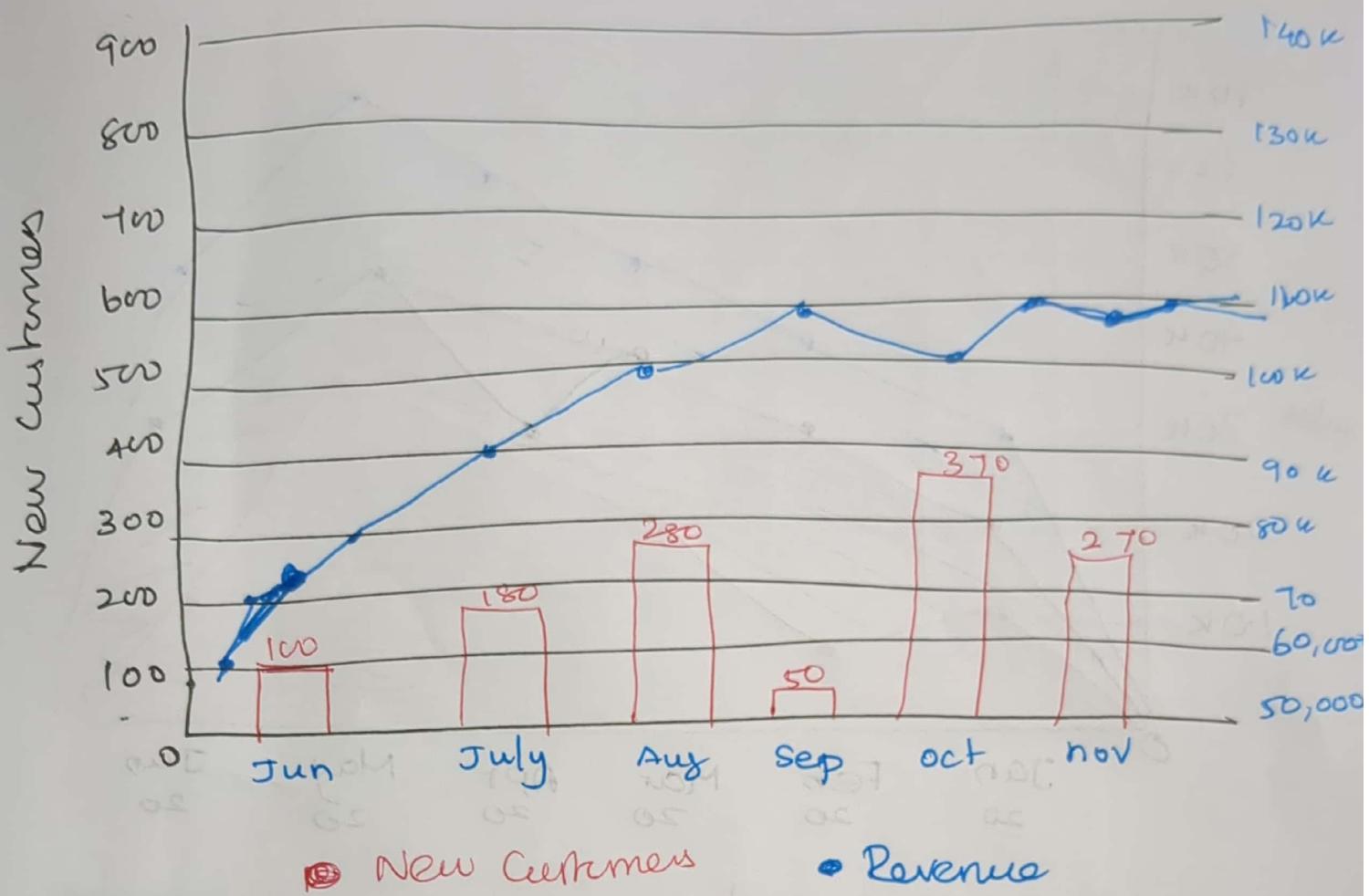
Eg - Track how many chats/emails →
Your team responds per month

- Use solid lines
- Don't plot more than 4 lines.



④ Dual axis chart

- 2 y axis & one x axis
- contains 3 data sets.
- 2 data share one x axis data
- uses
 - * Price & Volume of products
 - * Revenue & units.
 - * Sales & profit.

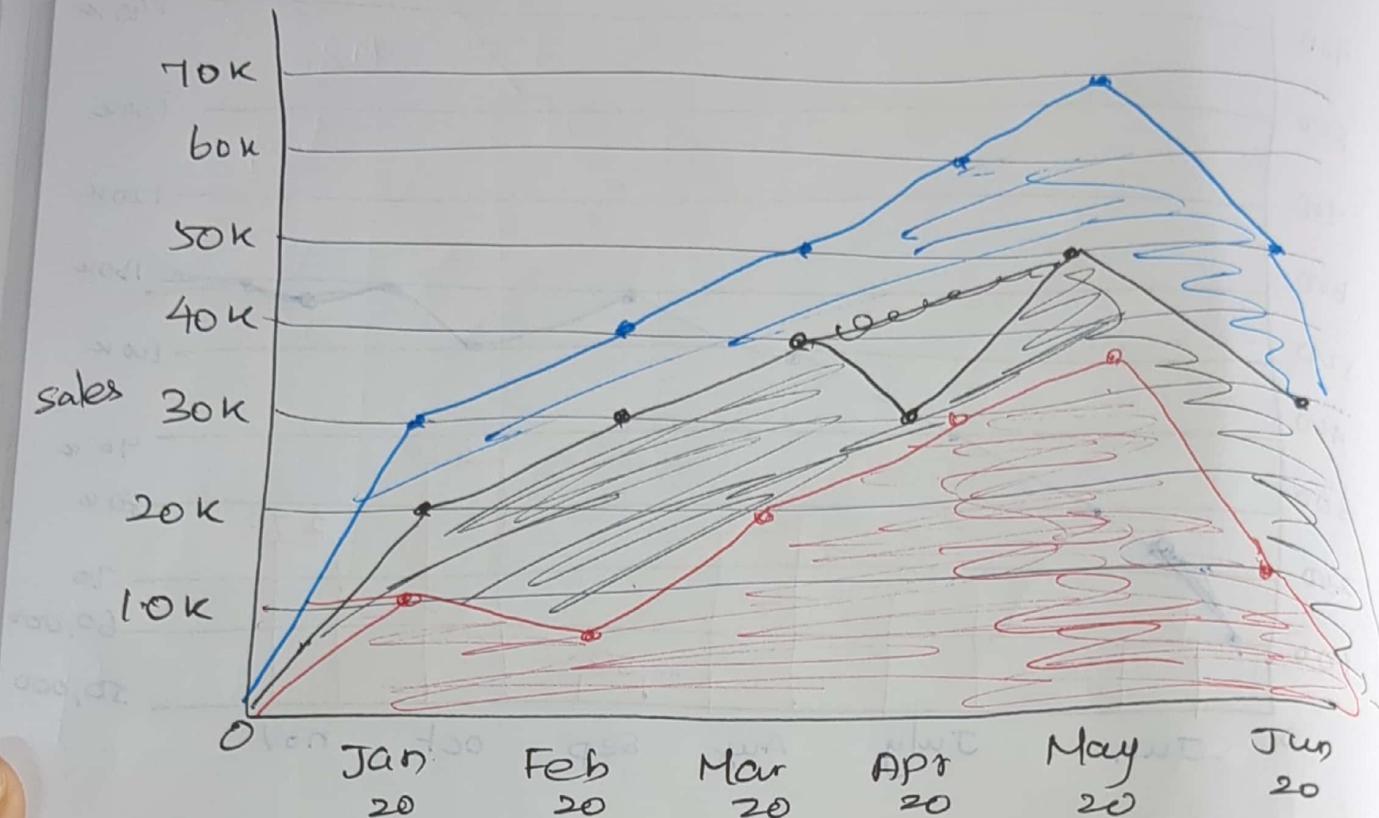


- how many new customers joined in each month.

- Revenue of that customers are bringing.

② Area chart

- combines line chart & bar chart
- Data points are plotted and connected by line segments.
- The area below the line is coloured.
- spot & analyze trends.
- Don't display more than 4 categories.



● → Store 1

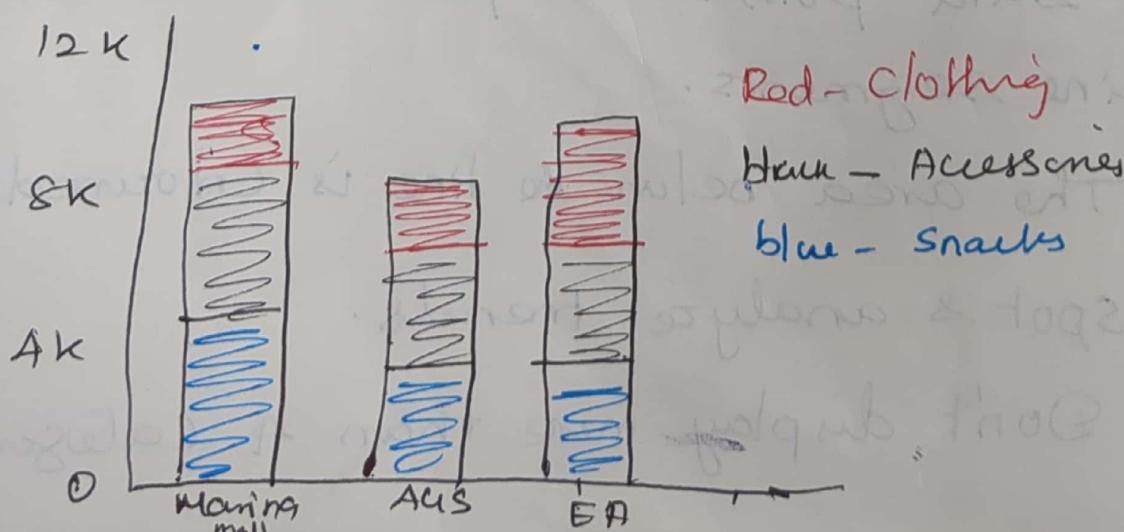
● → Store 2

● → Store 3

Stacked bar chart

- Used to compare many items.

- Extension of bar chart



Red - Clothing

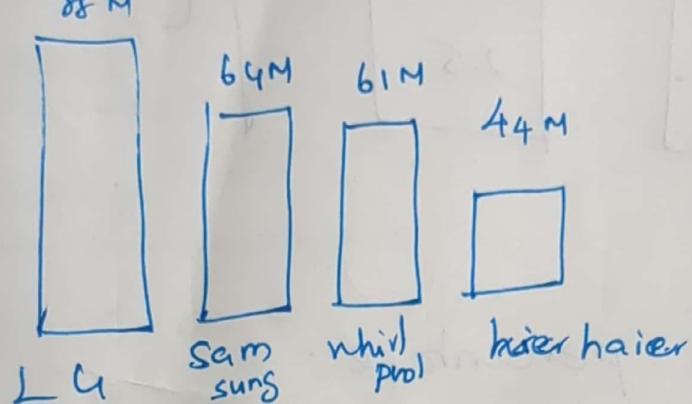
Black - Accessories

Blue - Snacks

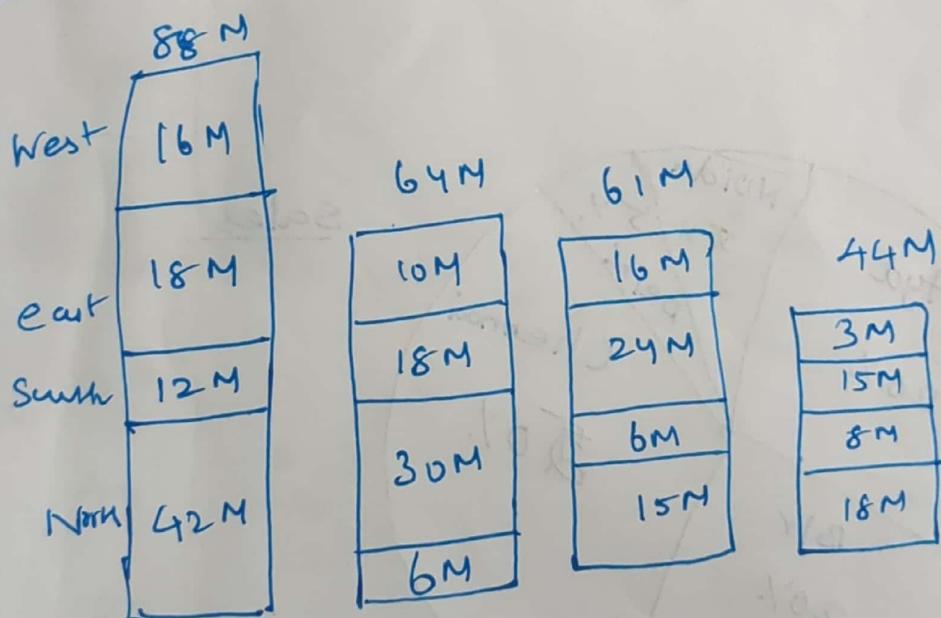
Mekko Chart

- also known as Mari mekko chart.
- Used for growth, market, sales analysis.
- Complex.

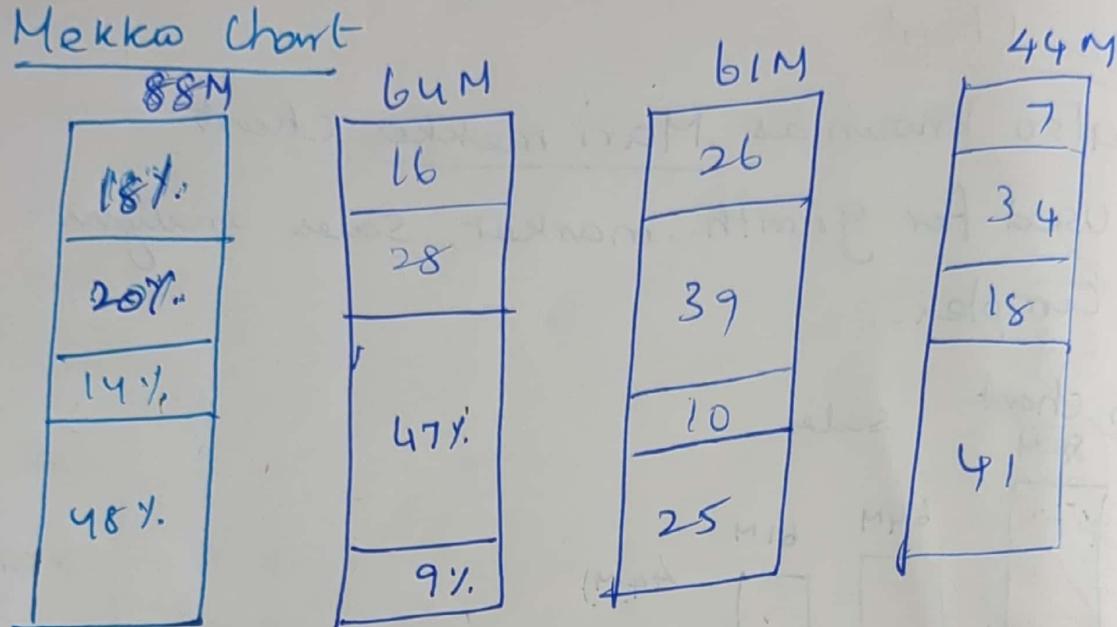
Column chart sales



Stacked Column chart

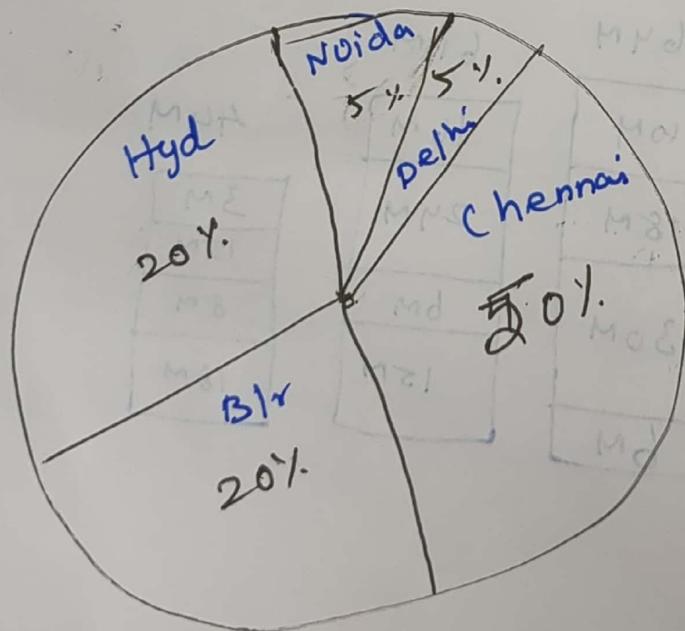


Mekka Chart



Pie Chart

- repr no in percentage.
- Total sum of segments $\rightarrow 100\%$.



Sales

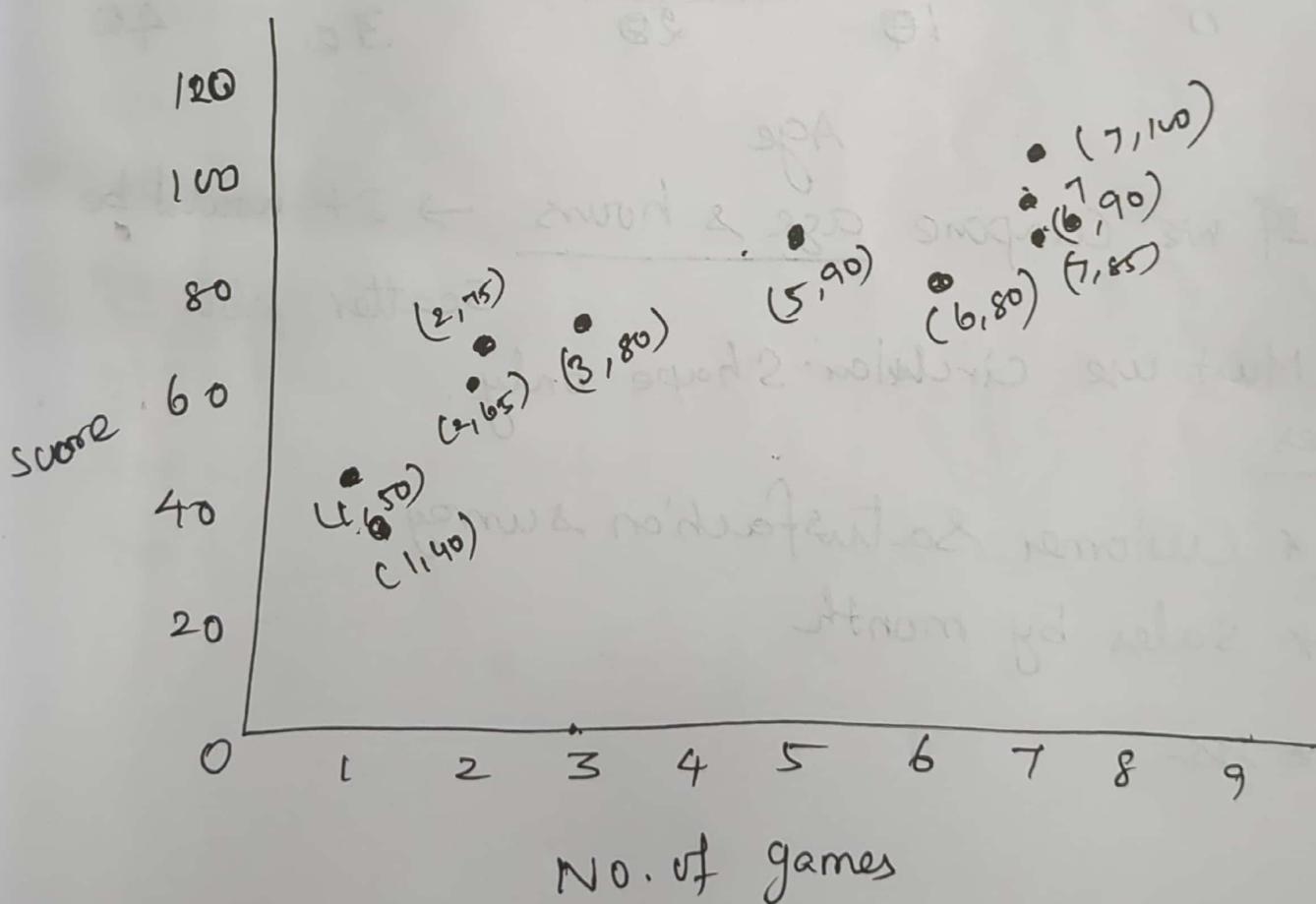
Scatter plot chart

- Shows relationship b/w 2 variables.
- Use scatter plot chart, when we have to compare 2 large set of numerical values.

Q	NO. of games	3	5	2	6	7	1	2	7	1	7
	scores	80	90	75	80	90	50	65	85	40	100

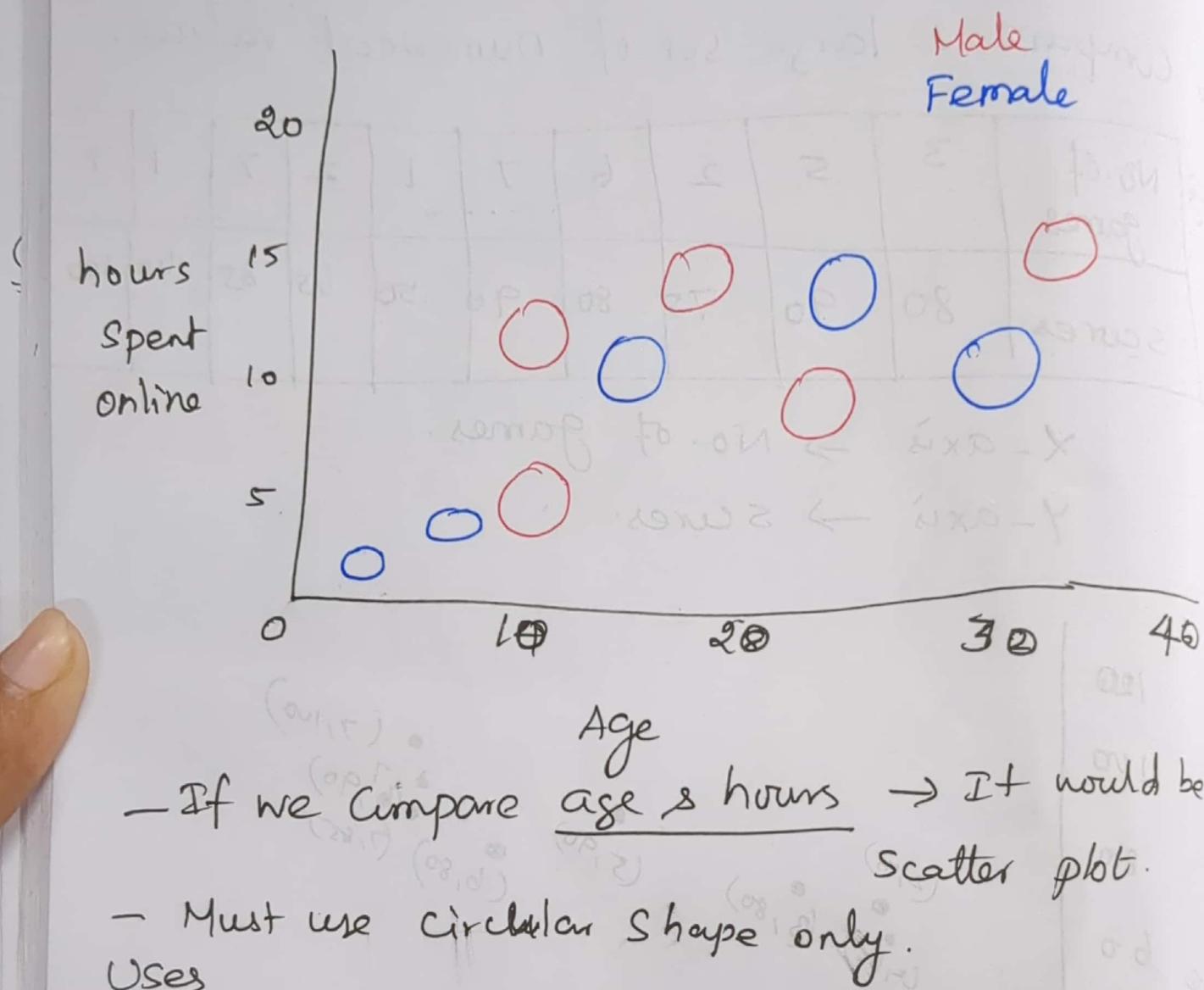
x-axis → No. of games.

y-axis → scores.



Bubble Chart

- Similar to scatter plot, but there is a third dataset shown by bubble.



- If we compare age & hours → It would be scatter plot.

- Must use circular shape only.

Uses

* Customer satisfaction survey.

* Sales by month

* etc

Waterfall chart

- Shows how an initial value changes with intermediate values.

- It shows the impact of overall result.

e.g

cash flow begin month \rightarrow 1000

" " end " \rightarrow 1050

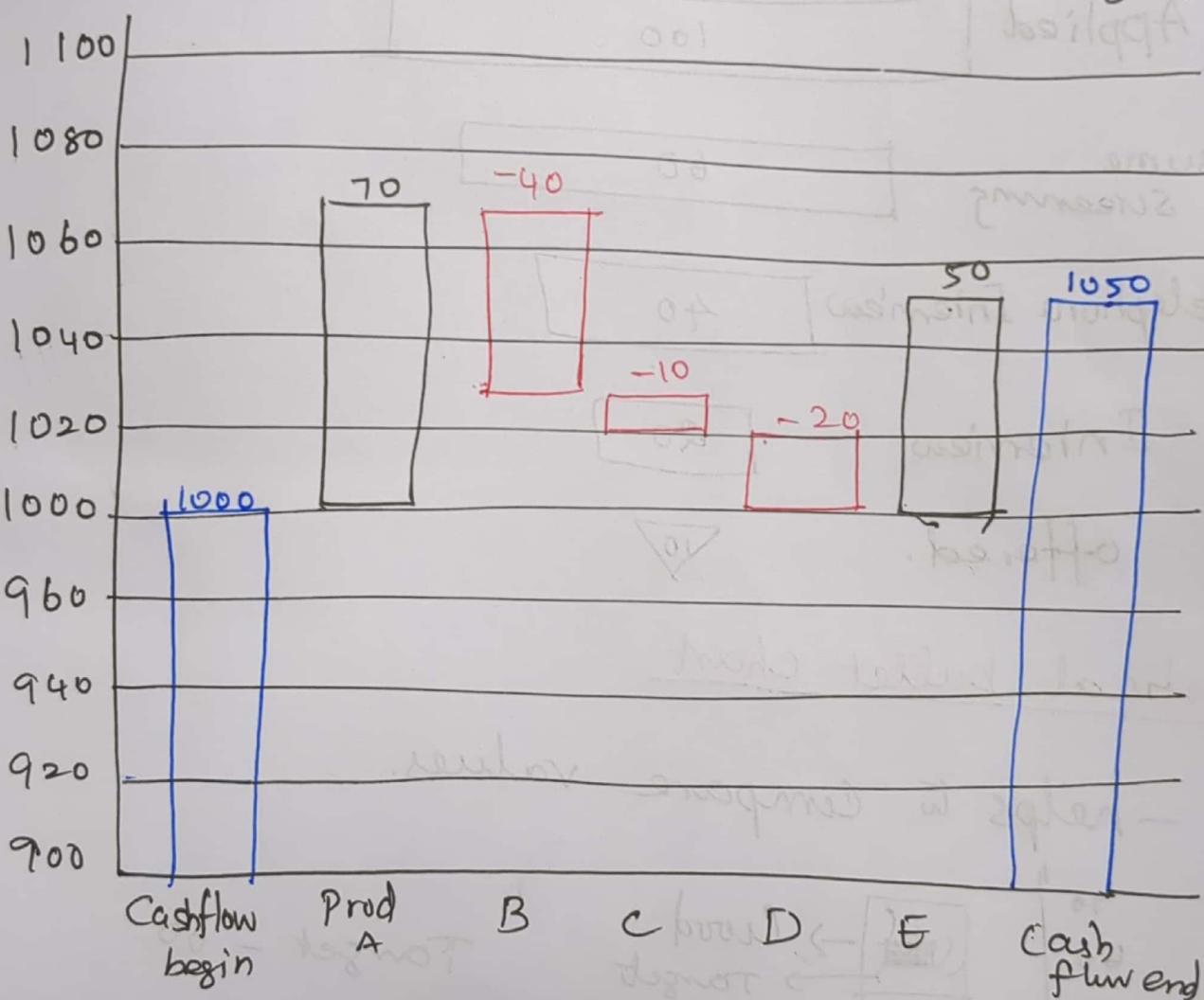
Prod A \rightarrow 70

B \rightarrow -40

C \rightarrow -10

D \rightarrow -20

E \rightarrow 50

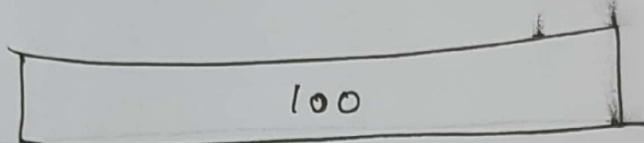


Funnel Chart

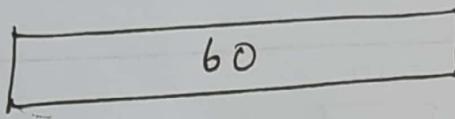
- how data moves through a process.
- Used in sales, recruitment, real estate -

Recruitment process

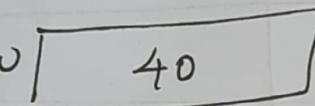
Applied



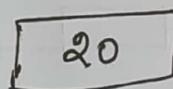
Resume screening



Telephone Interview



Interview



Offered.



Vertical bullet chart

- helps to compare values.

