

Big Data (CS-3032)

**Kalinga Institute of Industrial Technology
Deemed to be University
Bhubaneswar-751024**

School of Computer Engineering



Strictly for internal circulation (within KIIT) and reference only. Not for outside circulation without permission

3 Credit

Lecture Note

Course Contents



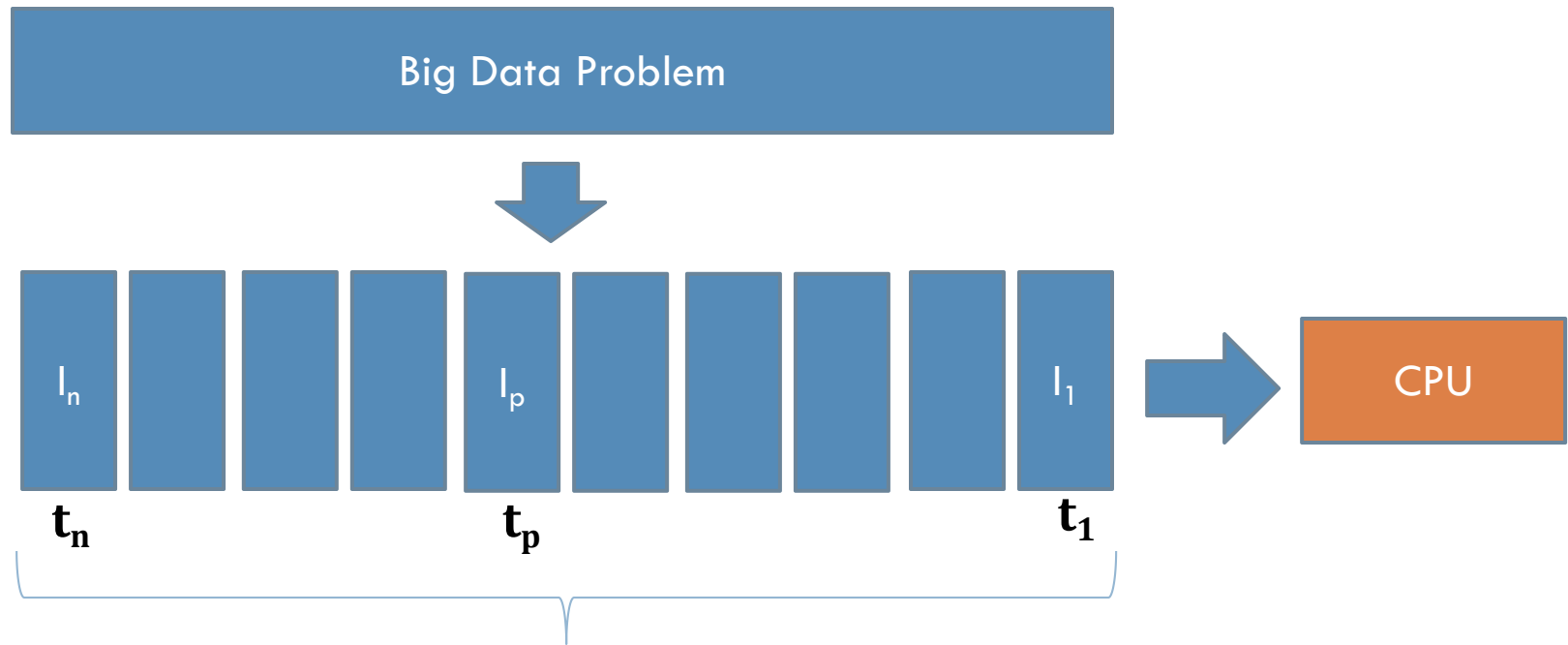
2

Sr #	Major and Detailed Coverage Area	Hrs
5	Big Data Tools Distributed and Parallel Computing for Big Data, Visualizations – Visual data analysis techniques, interaction techniques; Systems and applications. Exploring the Use of Big Data in Business Context, Use of Big Data in Social Networking, Business Intelligence, Product Design and Development	6

Traditional Sequential Computing



3

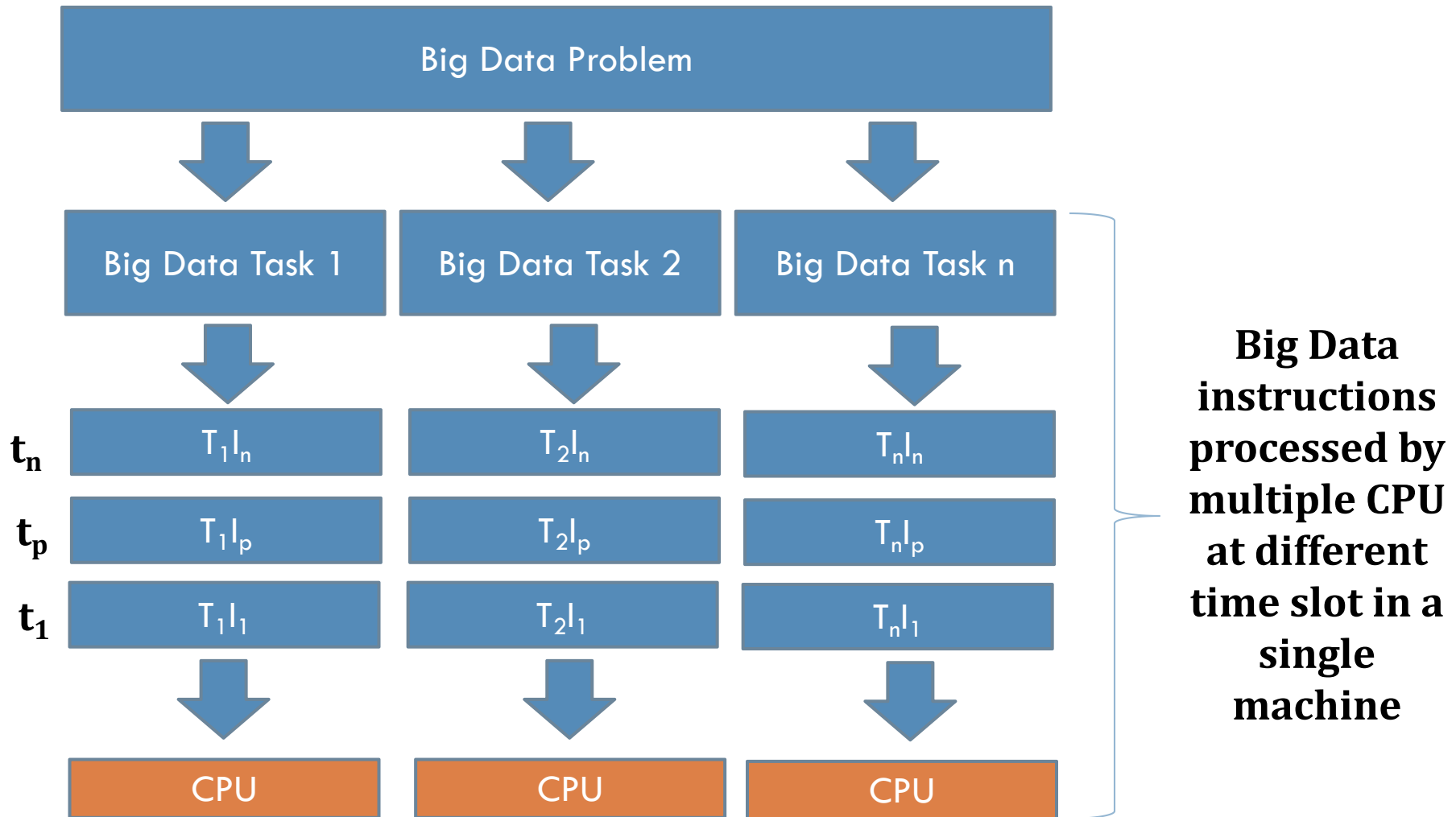


**Big Data Instructions processed by
a CPU at different time slot**

Parallel Computing



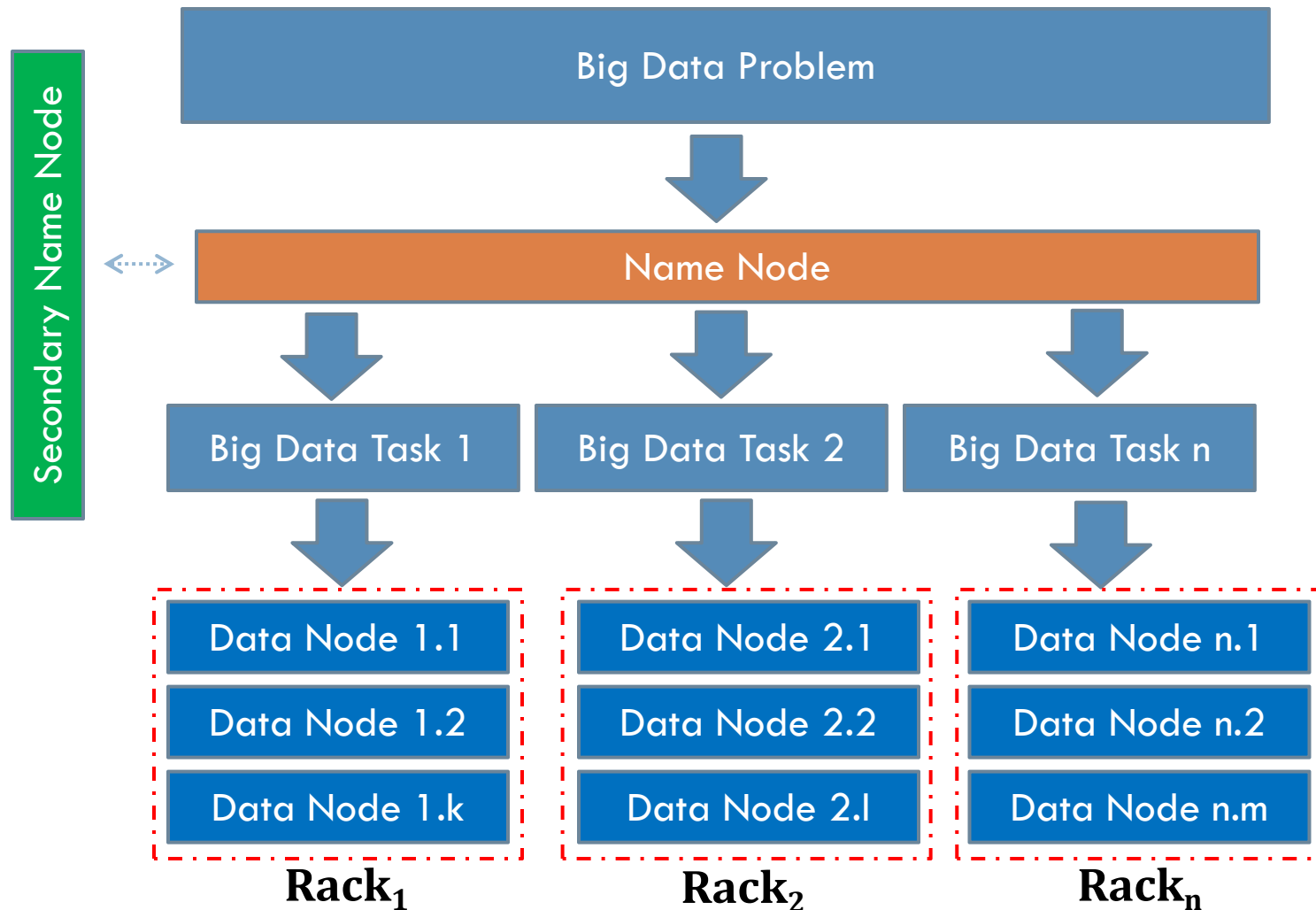
4



Distributed Computing



5



Distributed and Parallel Computing in Big Data



6

- ❑ Multiple computing resources are connected in a network and computing tasks are distributed across the resources. The sharing of tasks increases the speed as well as the efficiency of the system. Because of such reason, the distributed computing is more efficient than traditional methods of computing. It is primarily suitable to process huge amounts of data in a limited time.
- ❑ Another way to improve the processing capability of the system is to add additional computational resources to it. This will help in dividing complex computations into subtasks, which can be handled individually by processing units that are running in parallel. Such system are called parallel systems in which multiple parallel computing resources are involved to carry out computation simultaneously.
- ❑ Organization use both parallel and distributed computing techniques to process Big Data due to the important constraint for business, so called “Time”.

Distributed and Parallel Computing in Big Data cont'd



7

- ❑ Sometimes, computing resources develop technical challenge, and fail to respond. Such situations can be handled by virtualization, where some processing and analytical tasks are delegated to other resources.
- ❑ Another problem that often hampers data storage and processing activities is latency. It is defined as the aggregate delay in the system because of delays in the completion of individual tasks. Such a delay automatically leads to the slowdown in system performance as whole and is often termed as system delay. It affects data management and communication within and across various business units thereby, affecting the productivity and profitability of the organizations. Implementing distributed and parallel computing methodologies helps in handling both latency and data-related problem.

Visualization



8

Visualization is a pictorial representation technique. Anything which is represented in pictorial or graphical form, with the help of diagrams, charts, pictures, flowcharts etc. is known as visualization. Data visualization is a pictorial or visual representation of data with the help of visual aids such as graphs, bar, histograms, tables, pie charts, mind maps etc.

Ways of Representing Visual Data

The data is first analyzed and then the result is visualized. There are 2 ways to visualize a data, namely, Infographics and data visualization.

Infographics – It is the visual representation of information or data rapidly.

Data Visualization – It is the study of representing data or information in a visual form.

Difference: Infographics tell a premeditated story to guide the audience to conclusions (subjective). **Data visualizations** let the audience draw their own conclusions (objective).

An infographic can contain data visualizations but not the other way around.

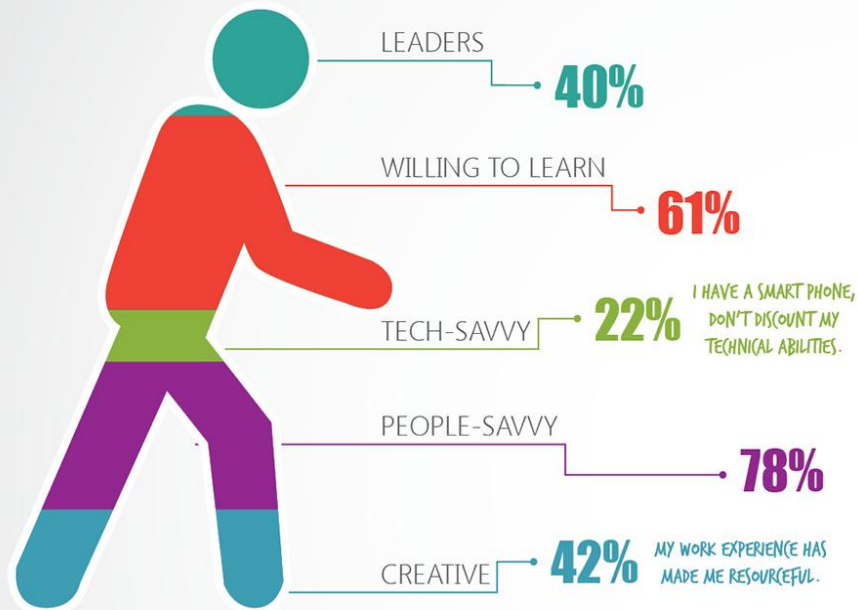
Infographics vs. Data Visualization



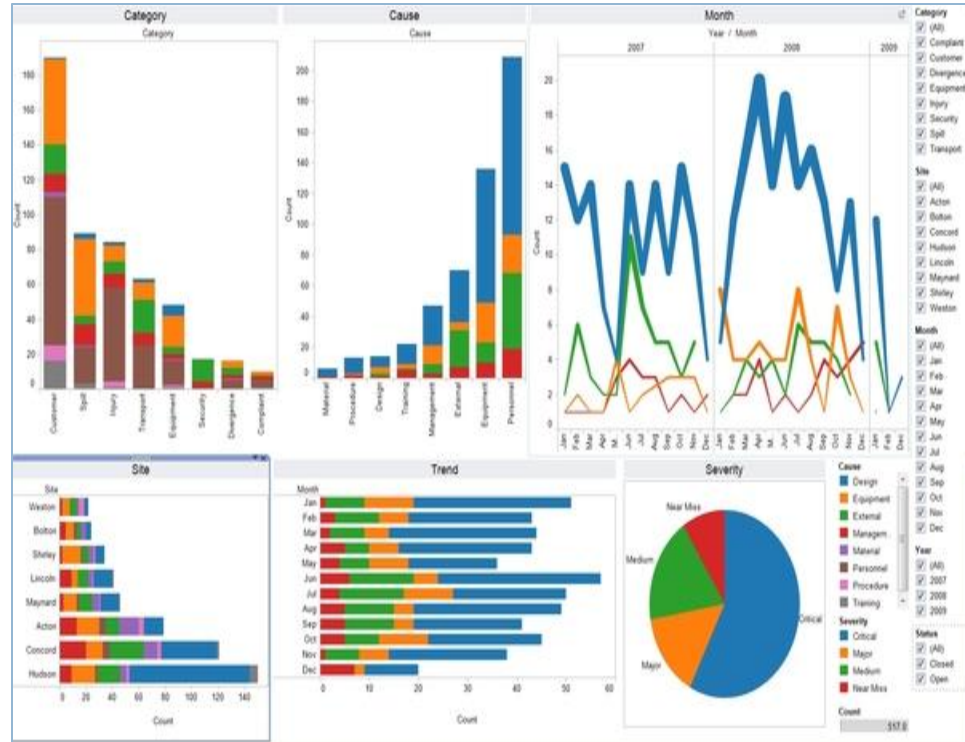
9

Infographics

HOW BABY BOOMERS DESCRIBE THEMSELVES



Data Visualization



Infographics vs. Data Visualization cont'd



10

Infographics are:

- ☐ Best for telling a premeditated story and offer subjectivity
- ☐ Best for guiding the audience to conclusions and point out relationships
- ☐ Created manually for one specific dataset

It is used for Marketing content, Resumes, Blog posts, and Case studies etc.

Data visualizations are:

- ☐ Best for allowing the audience to draw their own conclusions, and offer objectivity
- ☐ Ideal for understanding data at a glance
- ☐ Automatically generated for arbitrary datasets

It is used for Dashboards, Scorecards, Newsletters, Reports, and Editorials etc.

Data Visualization Purpose



11

- ❑ Data presented in the form of graphics can be **analyzed better** than the data presented in words.
- ❑ Patterns, trends, outliers and correlations that **might go undetected in text-based data** can be **exposed and recognized easier** with data visualization software.
- ❑ Data scientists can use data visualizations to make their information **more actionable**. Illustrations, graphs, charts and spreadsheets can turn dull reports into something illuminating, where it's **easier to gather insight and actionable results**.
- ❑ Data Visualization help to **transmit a huge amount of information** to the human brain **at a glance**.
- ❑ Data Visualization **point out key or interesting breakthrough in a large dataset**.

Techniques Used for Data Representation



12

Data can be presented in various forms, which include simple line diagrams, bar graphs tables, metrics etc. Techniques used for a visual representation of the data are as follows:

- ☐ Map
- ☐ Parallel Coordinate Plot
- ☐ Venn Diagram
- ☐ Timeline
- ☐ Euler Diagram
- ☐ Hyperbolic Trees
- ☐ Cluster Diagram
- ☐ Ordinogram
- ☐ Isoline
- ☐ Isosurface
- ☐ Streamline
- ☐ Direct Volume Rendering (DVR)

Map



13

It is generally used to represent the location of different areas of a country and is generally drawn on a plain surface. Google maps is generally widely used for data visualization. Now-a-days it is widely used for finding the location in different domains of country.

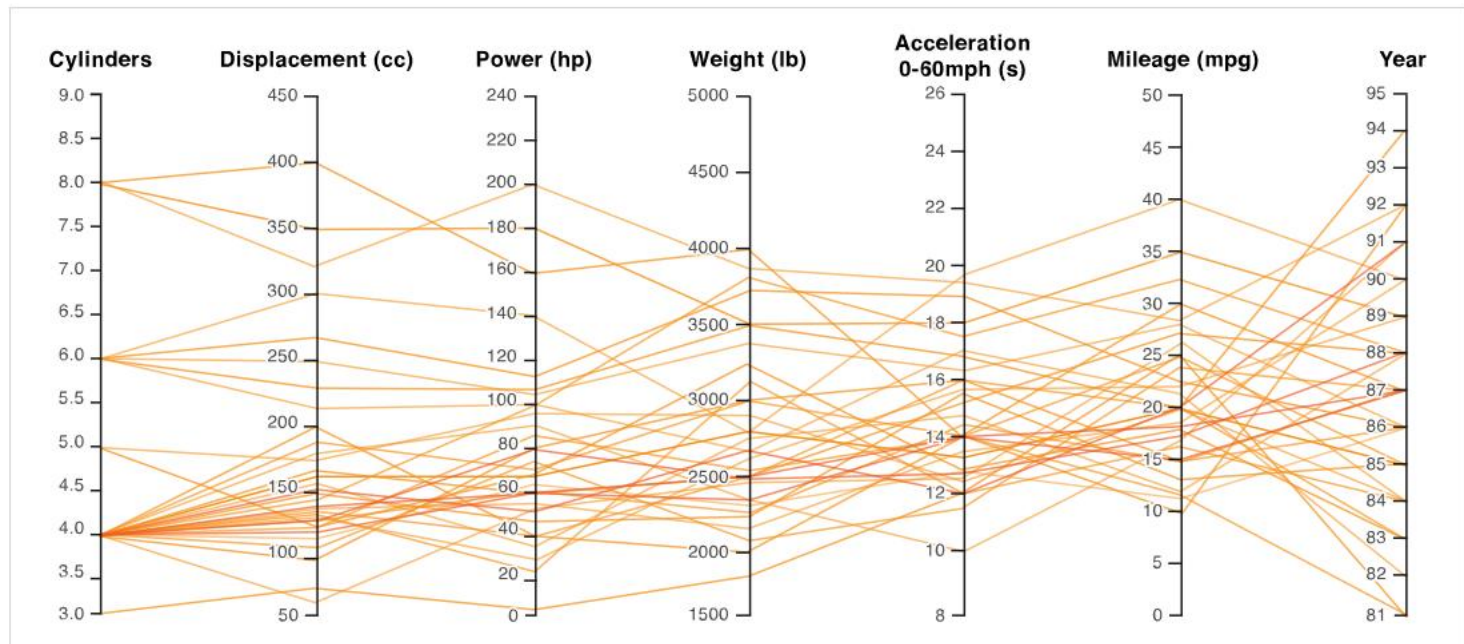


Parallel Coordinate Plot



14

It is a visualization technique of representing multidimensional data. This type of visualisation is used for plotting multivariate, numerical data. Parallel Coordinates Plots are ideal for comparing many variables together and seeing the relationships between them.

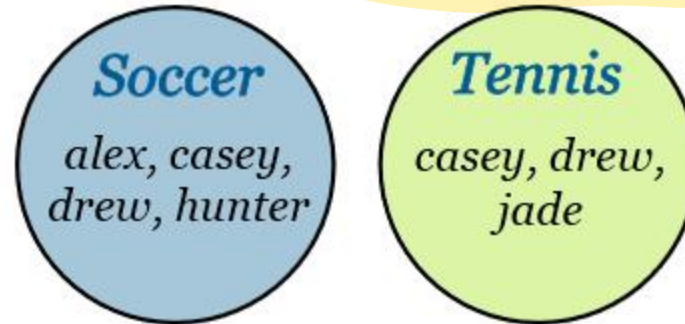


Venn Diagram

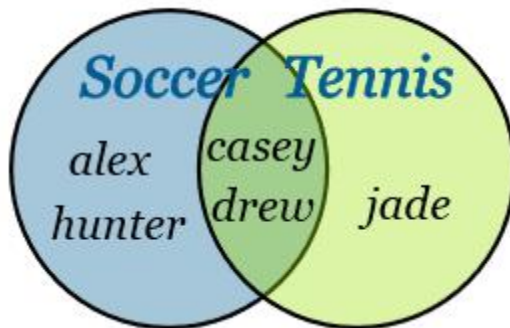


15

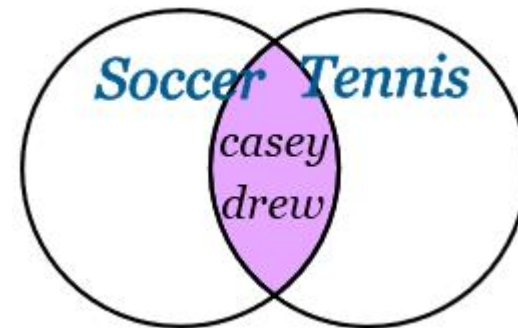
It is used to represent the logical relations between finite collection of sets.



List your friends that play Soccer OR Tennis



List your friends that play Soccer AND Tennis



- ☐ Draw the Venn Diagram to show people that play Soccer but NOT Tennis
- ☐ Draw the Venn Diagram to show people that play Soccer or play Tennis, but not the both.

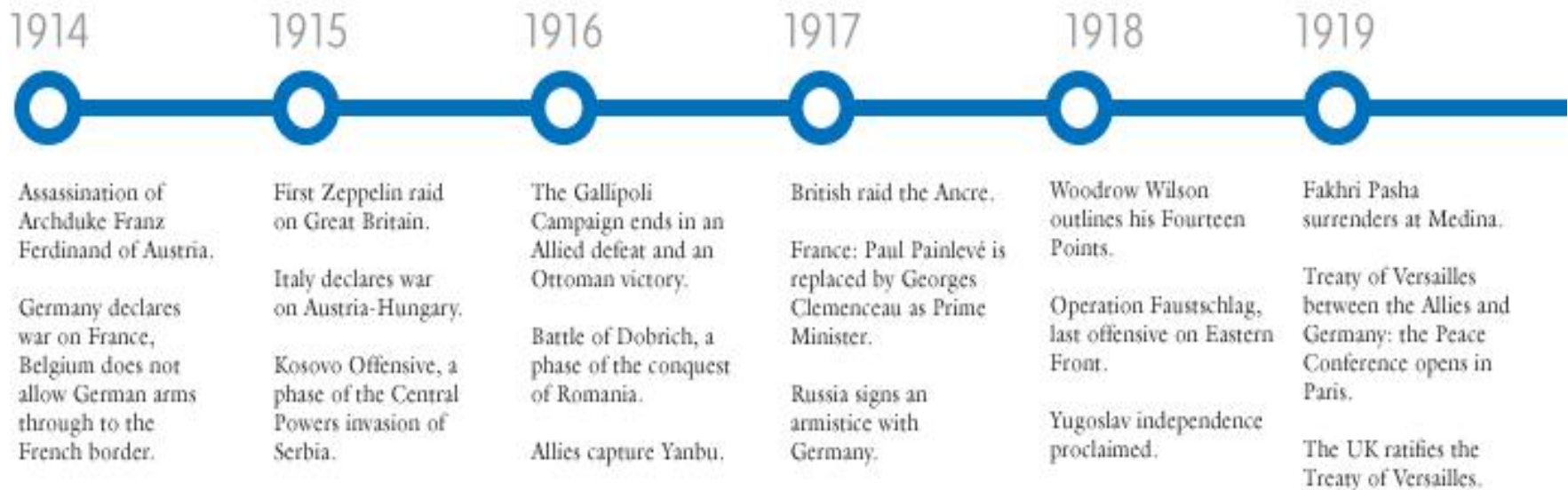
Timeline



16

It is used to represent a chronological display of events.

Timeline of World War I

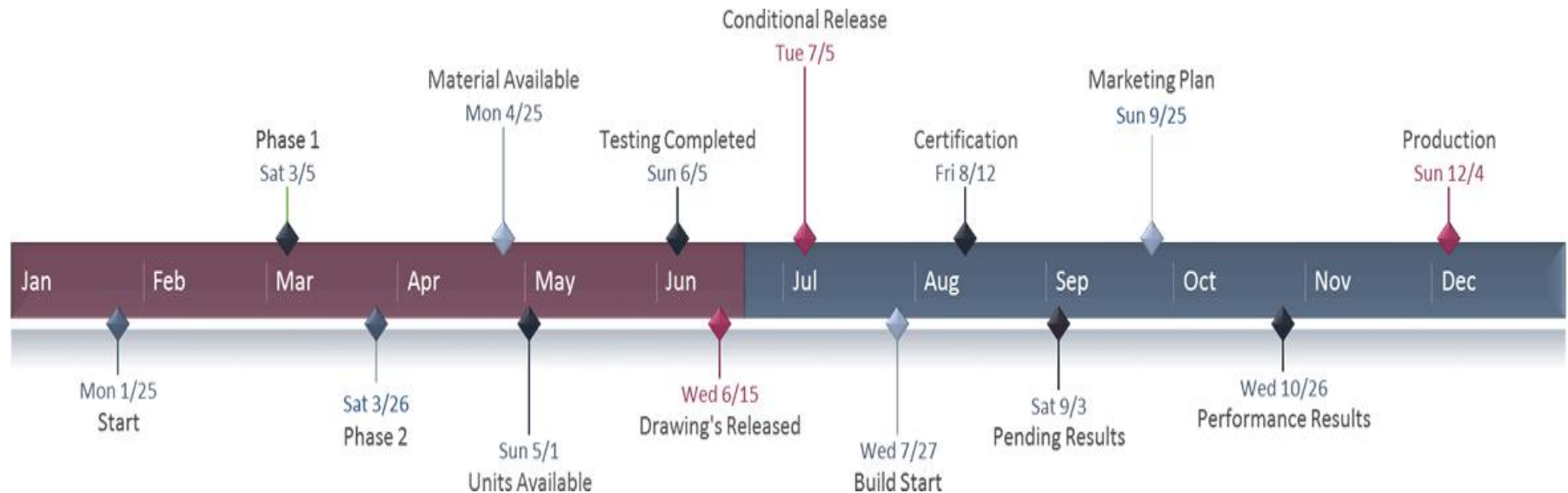


Source: datavizcatalogue.com

Timeline cont'd



17



Source: officetimeline.com

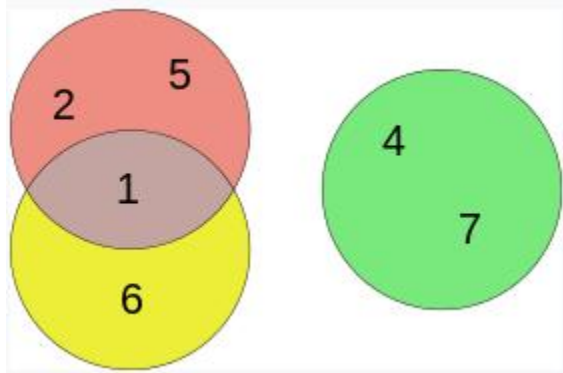
Euler Diagram



18

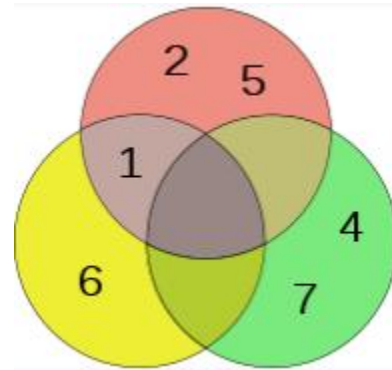
It is a representation of the relationships between sets.

Example: Let's take 3 sets namely $A = \{1, 2, 5\}$, $B = \{1, 6\}$ and $C = \{4, 7\}$. The Euler diagram of the sets looks like:



Source: wikipedia

Draw the Equivalent Venn Diagram



Class Exercise

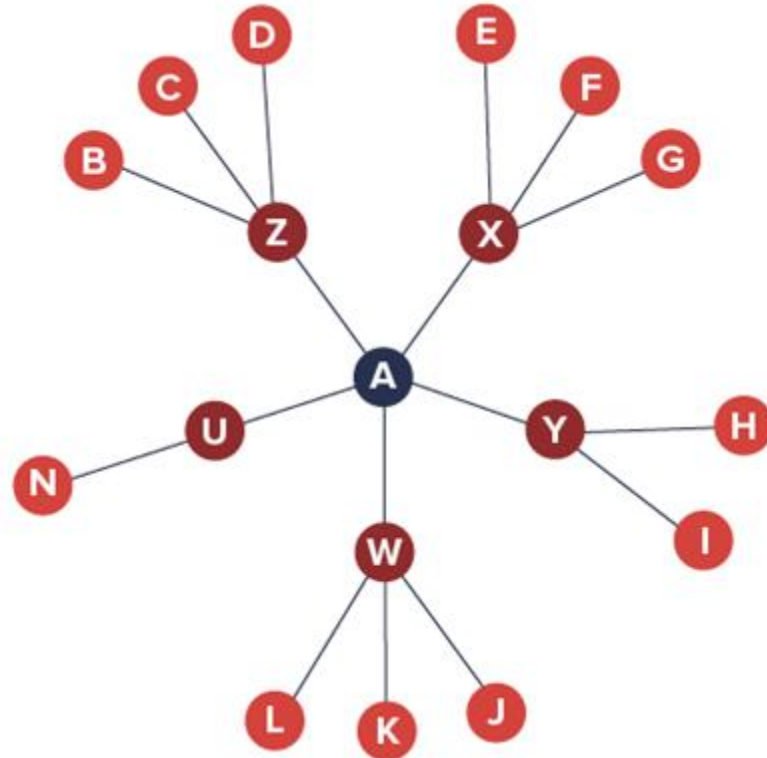
Draw the Euler diagram of the sets, $X = \{1, 2, 5, 8\}$, $Y = \{1, 6, 9\}$ and $Z = \{4, 7, 8, 9\}$. Then draw the equivalent Venn Diagram.

Hyperbolic Trees



19

A hyperbolic tree (often shortened as hypertree) is an information visualization and graph drawing method inspired by hyperbolic geometry.

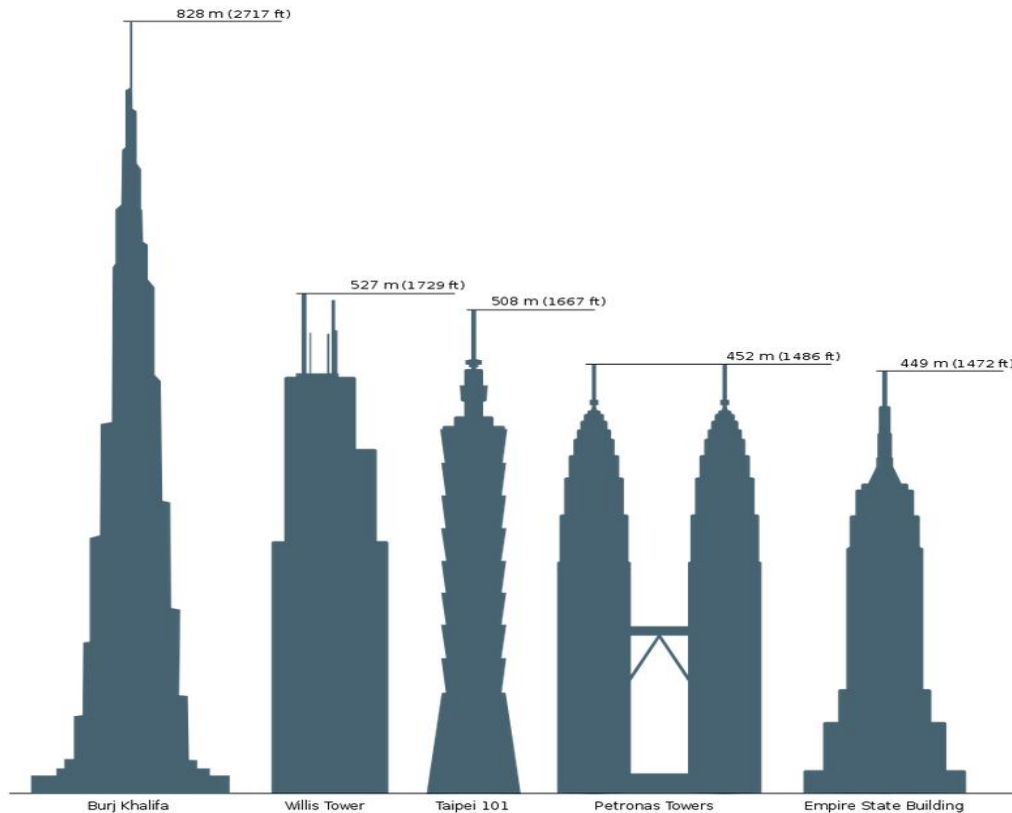


Cluster Diagram

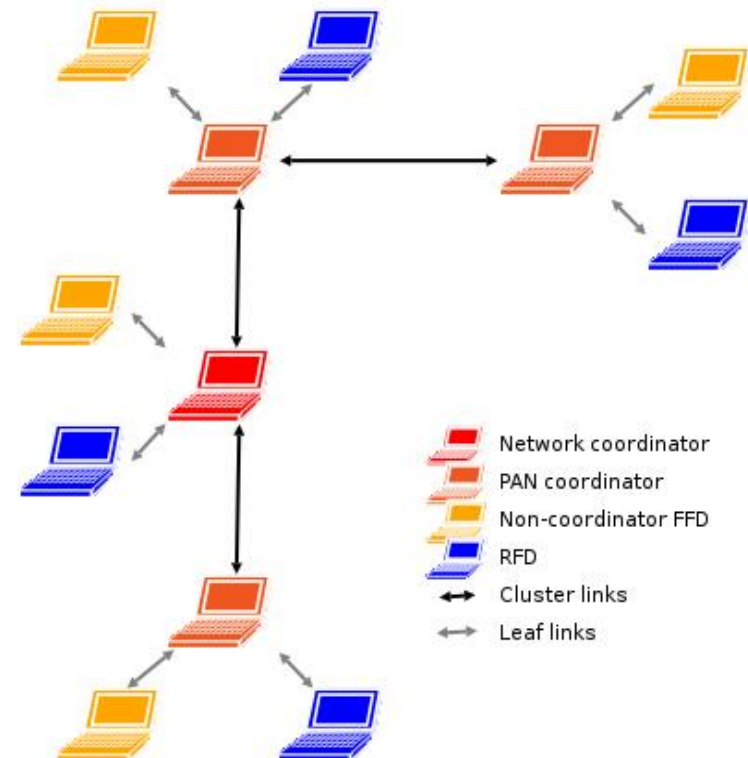


20

A cluster diagram or clustering diagram is a general type of diagram, which represents some kind of cluster. A cluster in general is a group or bunch of several discrete items that are close to each other.



Comparison diagram of sky scraper



Computer network diagram

Ordinogram



21

It is generally used to perform the analysis operation of various sets of **multivariate objects** which are generally used in different domain. Simple two-dimensional graph is an example of ordinogram.

Univariate data – This type of data consists of only one variable. The analysis of univariate data is thus the simplest form of analysis since the information deals with only one quantity that changes. It does not deal with causes or relationships and the main purpose of the analysis is to describe the data and find patterns that exist within it. The example of a univariate data can be height.

Height (cm)	164	167	170	170.4	176.5	180	179.2	165	175
-------------	-----	-----	-----	-------	-------	-----	-------	-----	-----

Multivariate data – This type of data involves two or more than two different variables. The analysis of this type of data deals with causes and relationships and the analysis is done to find out the relationship among the variables.

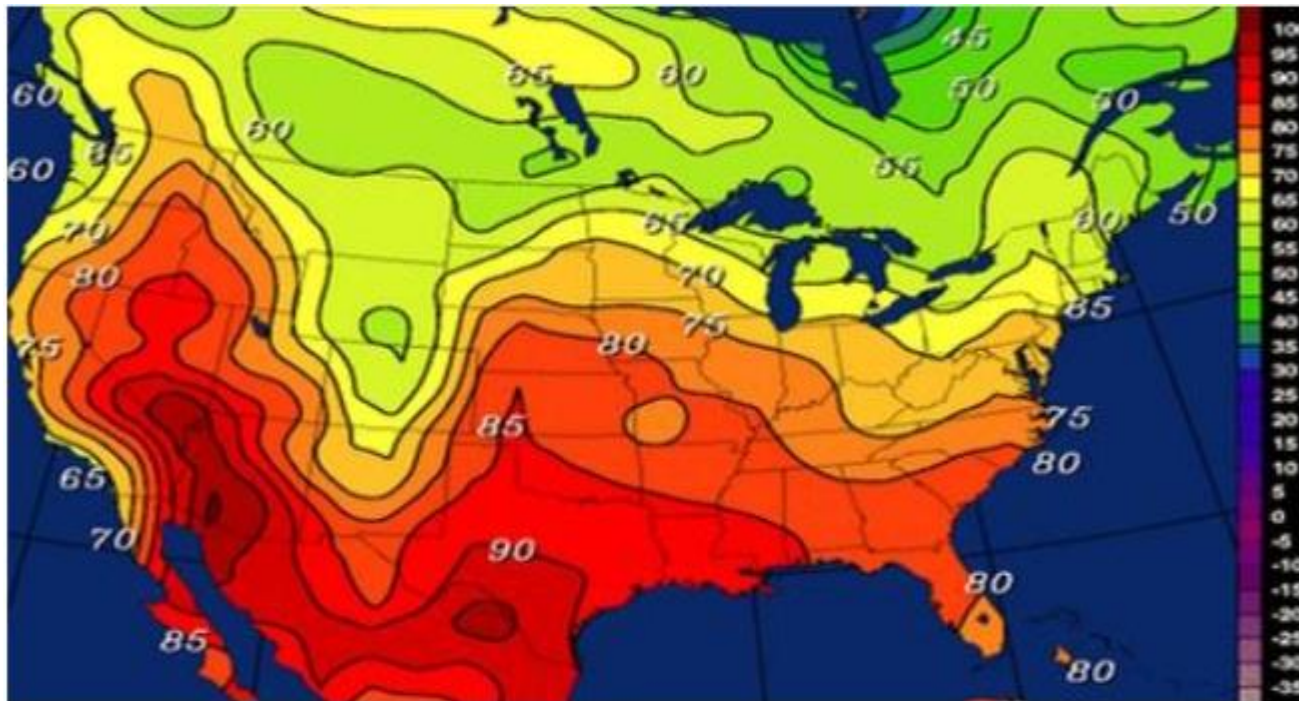
Temperature in Celsius	Ice Cream Sales
20	2000
35	5000

Isoline



22

It is basically a 2D data representation of a curved line that generally transfers constantly on the surface of the graph, the plotting of line generally drawn on the basis of data arrangement instead of data visualization.

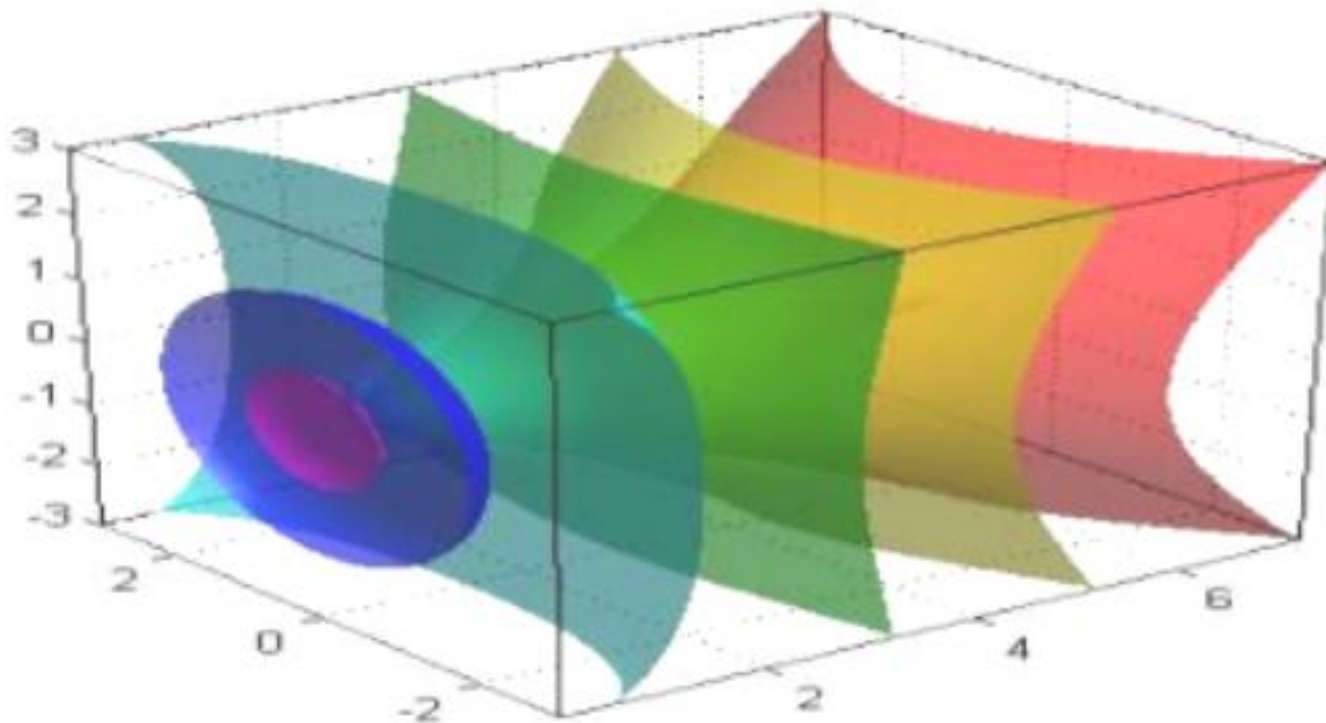


Isosurface



23

It is a 3D representation of an Isoline. Isosurfaces are designed to present points that are bound by a constant value in a volume of space i.e. in a domain that covers 3D space.

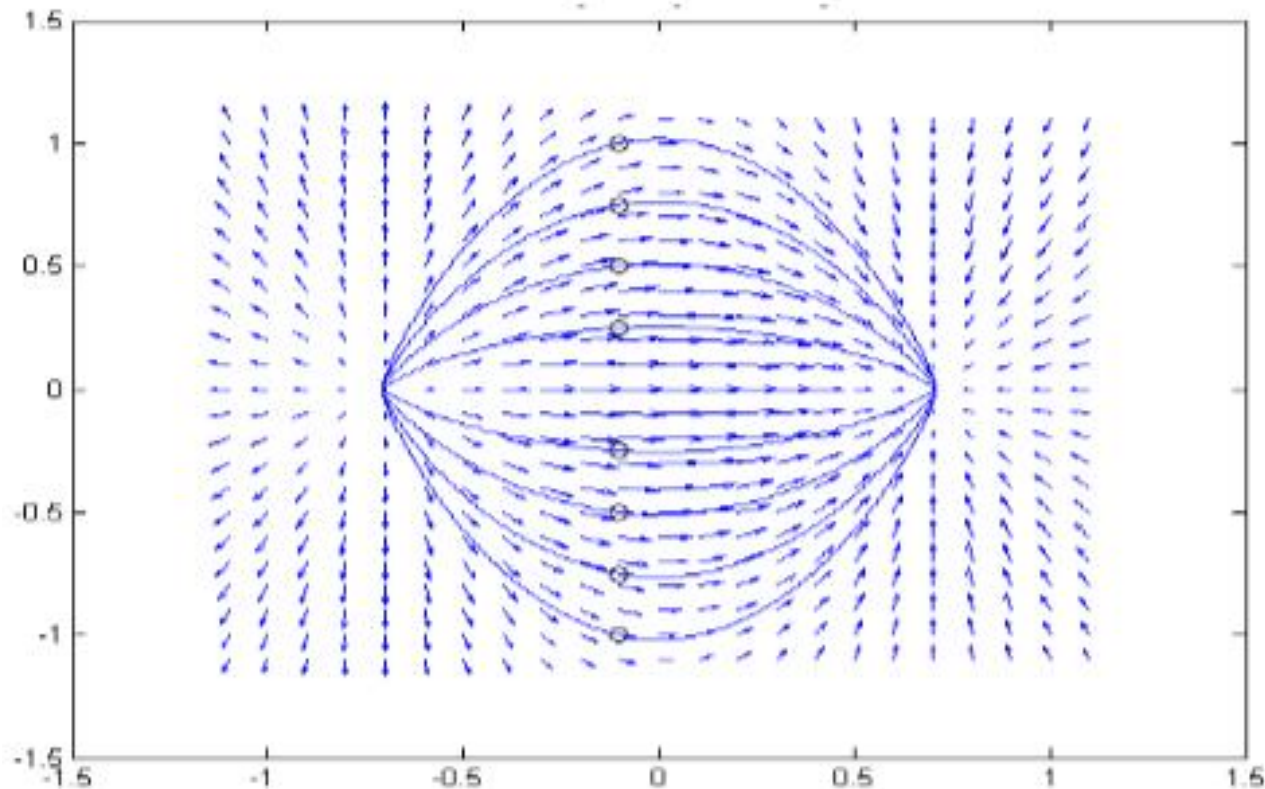


Streamline



24

It is a field that is generated from the description of velocity vector field of the data flow.

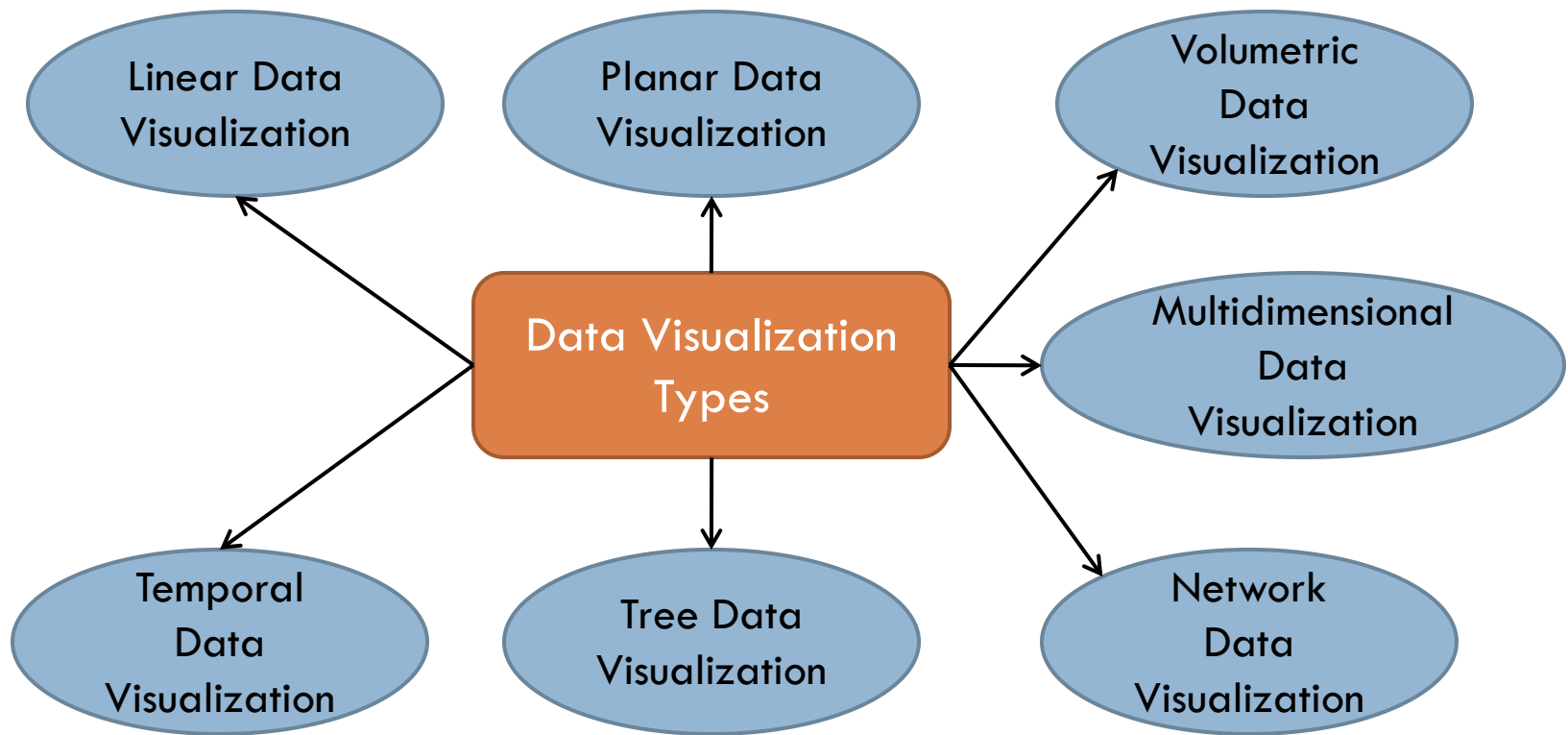


Types of Data Visualization



25

Data visualization can be done in different ways such as:

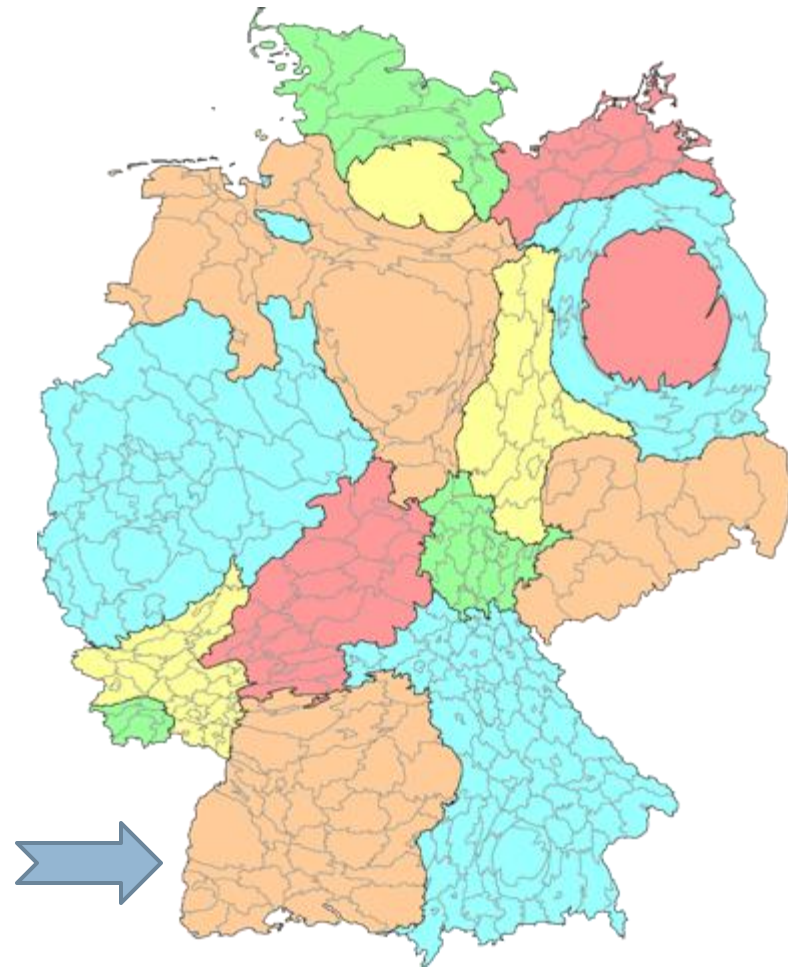


Types of Data Visualization cont'd



26

- ❑ **Linear Data Visualization:** Data always represented in list format. Basically it's not considered as a visualization technique rather is a data organization technique. No tool is used to visualize the data. It is also called as 1D data visualization.
- ❑ **Planar Data Visualization:** Data generally take in the form of images, diagrams or charts over a plane surface. The best example of this type of data visualization is Cartogram and dot distribution. A cartogram is a map in which some thematic mapping variable – such as travel time, or population is substituted for land area or distance. Some tools used to build planar data visualization are GeoCommons, Polymaps, Google Maps, Tableau Public etc.



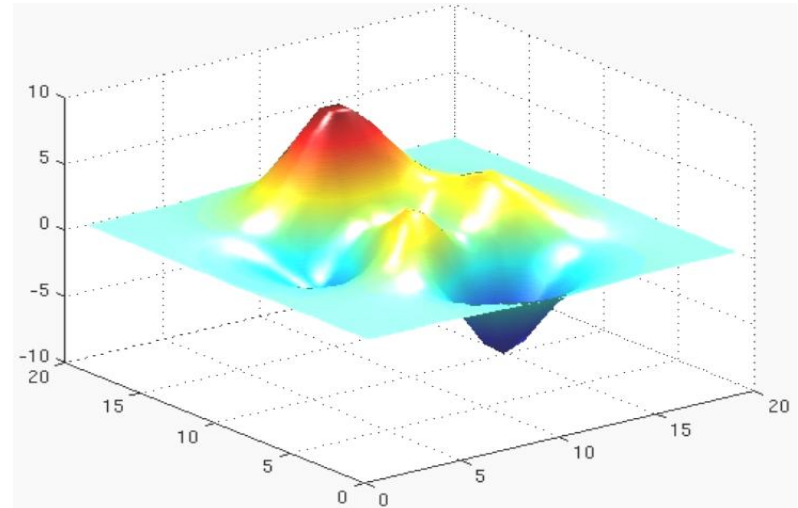
Germany-population-cartogram

Types of Data Visualization cont'd

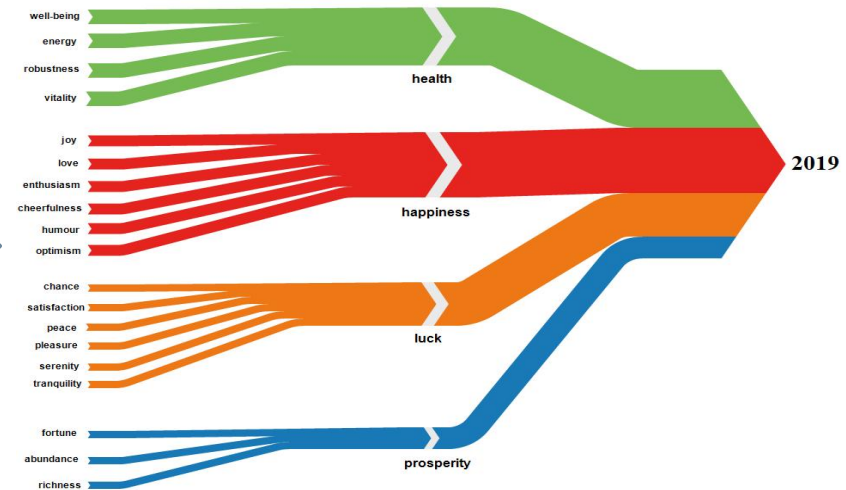


27

❑ **Volumetric Data Visualization:** the presentation of data generally involves exactly with three dimensions to present simulations, surface and volume rendering etc. and commonly used scientific studies. Basic tools used for it are AC3D, AutoQ3D, TrueSpace etc.



❑ **Temporal Data Visualization:** Sometimes, visualizations are time dependent so to visualize the dependence of analyses of time, the temporal data visualization is used which include Gantt chart, Time series and Sanky diagram etc.

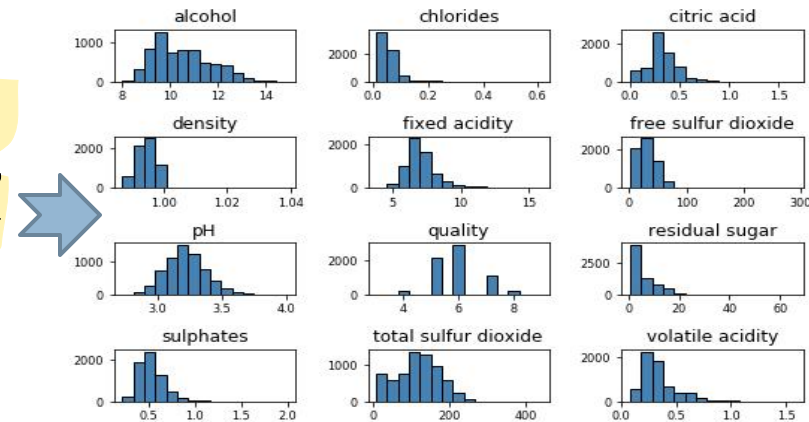


Types of Data Visualization cont'd

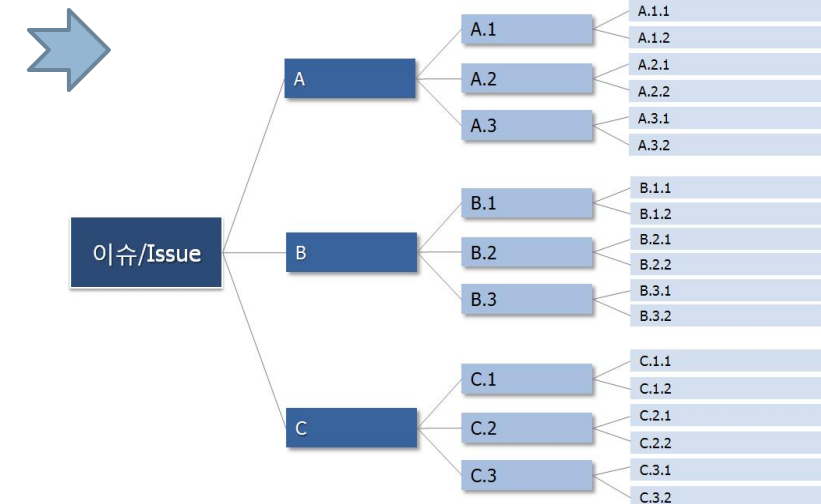


28

❑ **Multidimensional Data visualization:** Numerous dimension are generally used to represent the data. Generally pie charts, histograms, bar charts etc are generally used. Many Eyes, Google Charts, Tableau Public, etc. are some tools used to create such visualization.



❑ **Tree/Hierarchical Data visualization:** Sometimes, data relationships need to be shown in the form of hierarchies and to represent it, tree or hierarchical data visualization. Examples include hyperbolic tree, wedge-stack graph, etc. Google Charts, d3, etc. are some tools used to create such visualization.

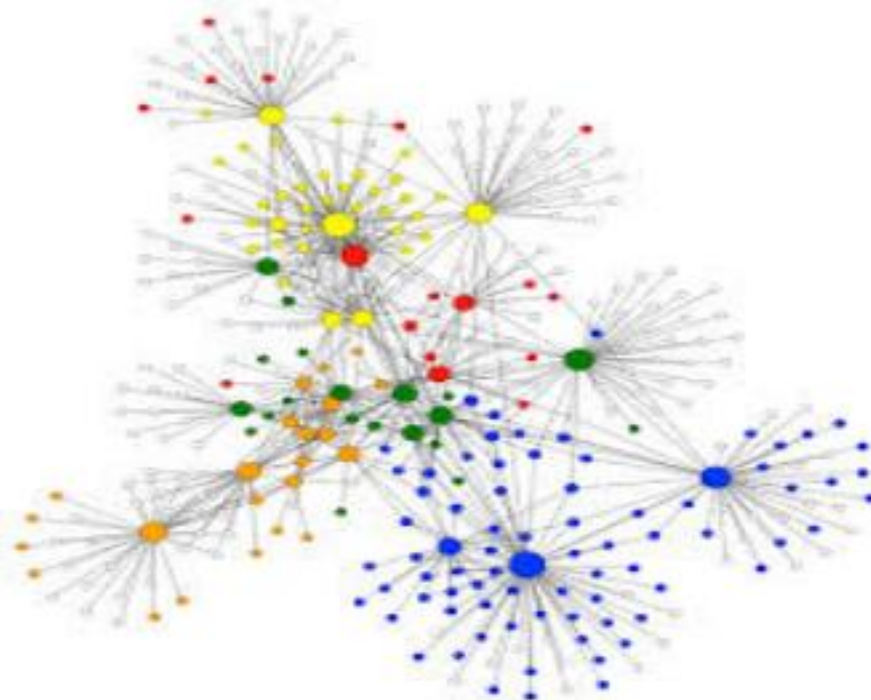


Types of Data Visualization cont'd



29

- ❑ **Network Data Visualization:** This approach is generally used to represent the relations that are too complex in the form of hierarchies. Some of the basic tools used for network data visualization are hive plot, Pajek, Gephi, NodeXL, Google Fusion Tables, Many Eyes, d3/Protovis etc.



Social Network Visualization

Visualization Interaction Techniques



30

The following interaction techniques are used in information visualization to overcome various limitations such as maximum amount of information is limited by the resolution.

- ❑ **Zooming:** It is one of the basic interaction techniques of information visualizations. It allows the user to specify the scale of magnification and increasing or decreasing the magnification of an image by that scale. This allows the user focus on a specific area and information outside of this area is generally discarded.
- ❑ **Filtering:** It is one of the basic interaction techniques often used in information visualization used to limit the amount of displayed information through filter criteria.
- ❑ **Details on demand:** This technique allows interactively selecting parts of data to be visualized more detailed and additional information on a point-by-point basis.
- ❑ **Overview-plus-Detail:** Two graphical presentation, wherein one shows a rough overview of the complete information space and neglects details, and the other one shows a small portion of the information space and visualizes details. Both are either shown sequentially or in parallel.

Application of Data Visualization



31

There are 3 ways to use data visualization in a company.

- 1. Internal Communication:** Any key data that influences decision-making is prime for data visualization. This is specifically true for the information delivered to higher-ups such as boss or other key stakeholders. Examples are presentation, reports or financial statements.
- 2. Client Reporting:** With data visualization, results reporting to clients or customers is more impactful.
- 3. Marketing Content:** Public-facing content for thought leadership or promotion is more credible with data. Content such as blogs, whitepapers, infographics etc. can be beneficial.

**THANK
YOU!**