

## Introduction to visualization

### Introduction:-

Data visualization is the representation of data through use of common graphics, such as charts, plots, infographics and even animations. These visual display of information communicate complex data relationships and data-driven insights in a way that is easy to understand.

### Types of Data Visualization

Some of the common visualization techniques such as

1. Tables :- This consists of rows and columns used to compare variables. Tables can show a great deal of information in a structured way.
2. Pie charts and stacked bar charts:-

These graphs are divided into sections that represent parts of a whole. They provide a simple way to organize data and compare the size of each component to one other.

3. Line charts and area charts:-

These visuals show change in one or more quantities by plotting a series of data point over time and are frequently used within predictive analytics.

4. Histograms!- The graph plots a distribution of numbers using a bar chart, representing the quantity of data that falls within a

Particulars songs.

5. Scatter plots:- These visuals are beneficial in revealing the relationship between two variables and they are commonly used with regression data analysis.

6. Heat Maps:- These graphical representations displays are helpful in visualizing behavioral data by location. This can be a location on a map or even a webpage.

7. Tree Maps:- It displays hierarchical data as a set of nested shapes, typically rectangles. Treemaps are great for comparing the proportions between categories via their area size.

→ Visualizing Data - Mapping Data onto Aesthetics

→ Aesthetics and Types of Data

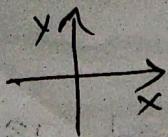
→ Data visualizations map data values into quantifiable features of the resulting graphic. These features are referred as Aesthetics.

→ Aesthetics describe every aspect of a given graphical element.

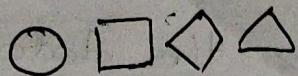
→ All graphical elements have a shape, size and a color.

→ Commonly used aesthetics in data visualization are position, shape, size, color, line width, line type.

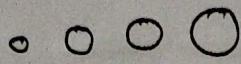
position



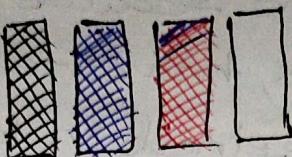
shape



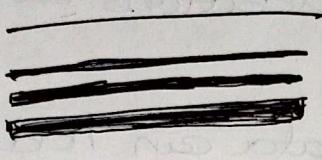
size



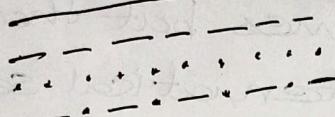
color



Line width



Line type

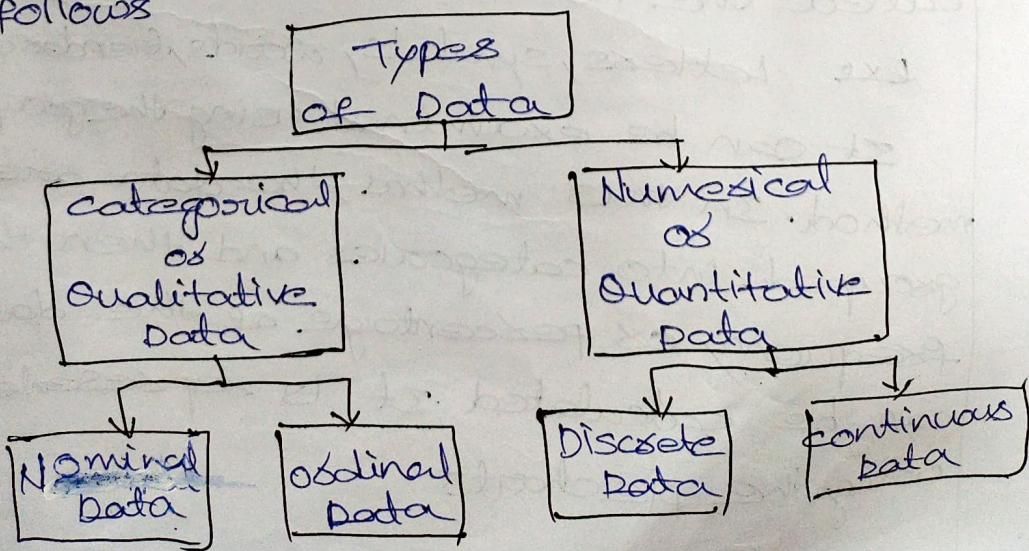


→ Aesthetics will fall into two groups  
They are continuous and discrete.

→ position, size, line width and color will represent both continuous and discrete data  
→ shape and line type will represent only discrete data.

### Types of Data

We can classify the data as follows



## 1. Qualitative Data

- It is also called as categorical data. It describes the data that fits into the categories.
- It describes the features such as a person's gender, home town etc. The measures are defined in terms of natural language specifications but not in terms of numbers.
- Categorical data can hold numerical values sometimes but those values do not have a mathematical sense.

Ex:- Date of birth, pincode etc.

This qualitative data is mainly divided into two types

1. Nominal Data

2. Ordinal Data

### 1. Nominal Data :-

It helps to label the variables without providing the numerical value. It is also called the nominal scale.

Ex:- letters, symbols, words, genders etc

It can be examined using the grouping method. In this method, the data are grouped into categories and then the frequency or percentage of the data can be calculated. It is represented using the pie charts.

## 2. **Ordinal Data**

It is a type of data that follows a natural order. This variable is mostly found in surveys, economics and so on.

This is commonly represented using a bar chart.

## 2. **Quantitative Data**

It is also known as numerical data which represents the numerical value.  
Ex:- size, weight, height, length.

The quantitative data is divided into two types

1. Discrete Data
2. Continuous Data

1. Discrete Data:- It contains only a finite number of possible values. Those values cannot be subdivided meaningfully.

Ex:- No of students in the class.

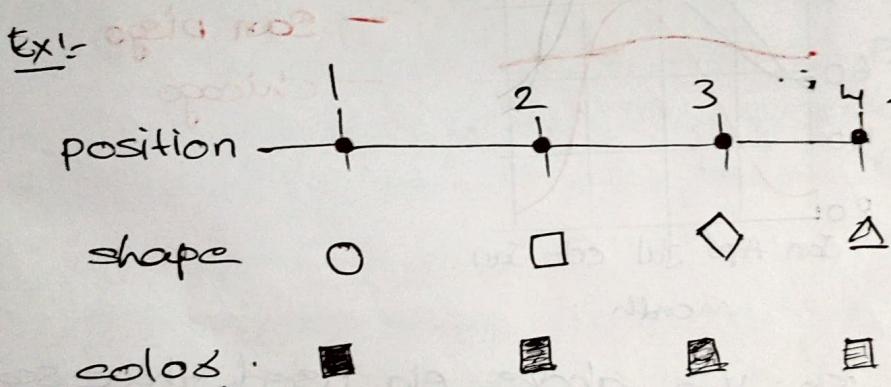
2. Continuous Data:- It is a data that can be calculated. It has an infinite no of probable values that can be selected within a given specific range.

Ex:- Temperature range.

→ Types of Variables encountered in typical data visualization scenarios are

Type of Variable	Examples	Appropriate scale	Description
Quantitative/ numerical continuous	1, 3, 5-7, 83 $1.5 \times 10^2$	continuous	Arbitrary numerical values, i.e integers as well as non-integer values.
Quantitative/ numerical discrete	1, 2, 3, 4	discrete	Numbers in discrete units EX- 0.5, 10, 15
Qualitative/ categorical uncodexed	dog, cat, fish	discrete	Categories without orders. These variables are called as factors
Qualitative/ categorical ordered	good, fair, poor	discrete	Categories with orders. These variables are called as ordered factors
Date or Time	Jan 11 2023 8:09 AM	continuous or discrete	specific days and times.
Text	The quick brown fox jumps over the lazy dog	None or discrete	Free-form text can be treated as categorical if needed.

- scales Map Data Values onto Aesthetics
- To map data values onto aesthetics we need to specify which data values correspond to which aesthetics value.
- The mapping between data values and aesthetics values is created via scales
- A scale defines a unique mapping between data and aesthetics
- A scale must be one to one such that for each specific data value there is exactly one aesthetics value and viceversa.
- If scale is not one to one then the data visualization becomes ambiguous.

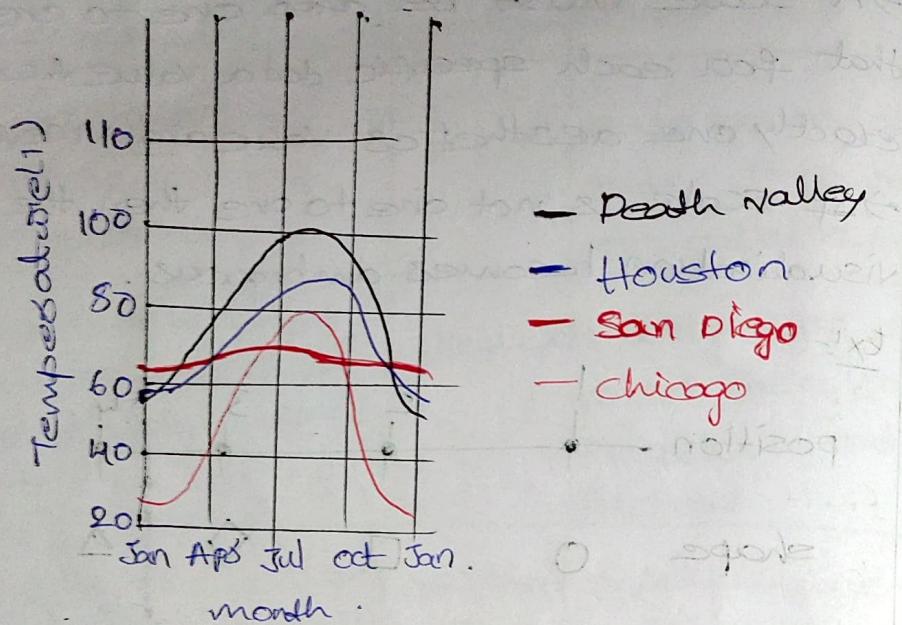


scale links data values to aesthetics. Here the numbers 1 through 4 have been mapped onto a position scale, a shape scale and a color scale. For each scale, each number corresponds to a unique position, shape or color and viceversa.

Let us take the dataset of daily temperature normals for four selected locations in U.S.

we map temperature onto the y-axis and day of the year onto the x-axis, location onto colors, and visualize these aesthetics with solid lines.

The result is a standard line plot showing the temperature normals at the four city locations as they change during the year.

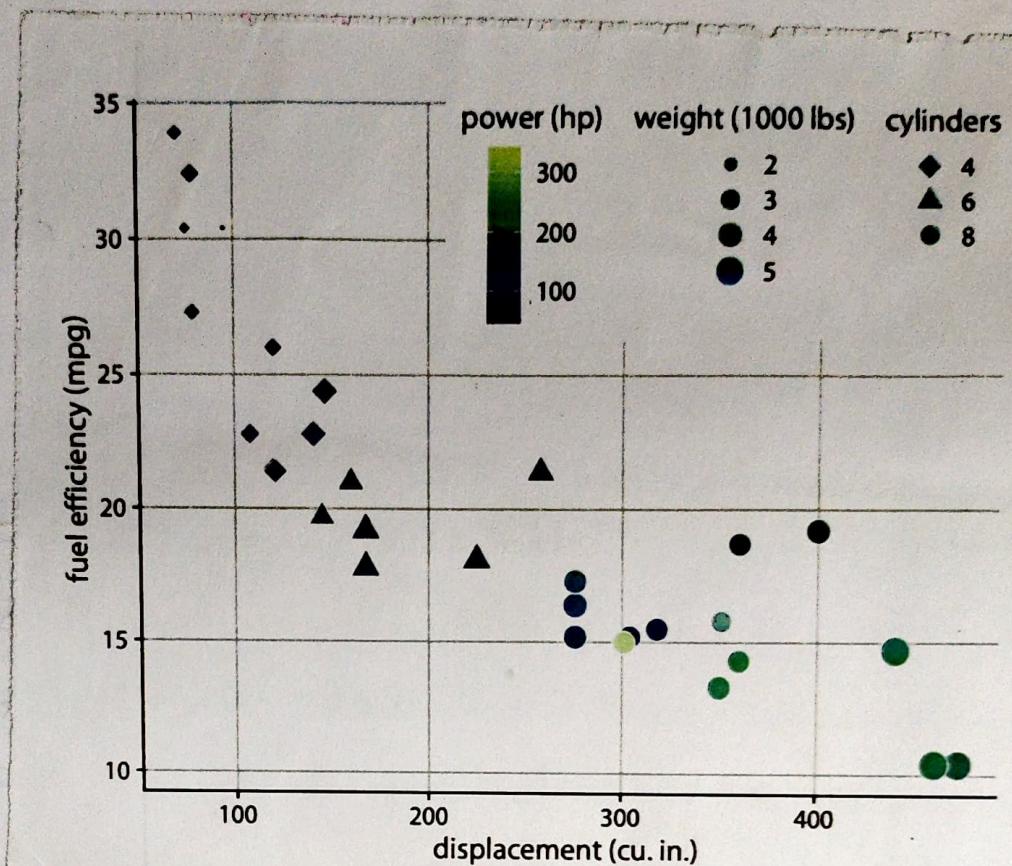


In the above fig. used three scales in total, two position scales and one color scale.

This is a typical no of scales for a basic visualization, but we can use more than three scale at once.

Let us consider the following figure, which represent fuel efficiency versus displacement. Here we uses five scales, two position scales, one color scale, one size scale and one shape scale and.

all the scales represent a different variable from the dataset.



## → color scales

These are three fundamental uses for color in data visualizations. They are

1. colors as a tool to distinguish.
2. colors to represent data values
3. colors as a tool to highlight.

### 1. color as a Tool to Distinguish:-

→ We use colors as a means to distinguish discrete items or groups that have an intrinsic order, such as different countries on a map or different manufacturers of a certain product.

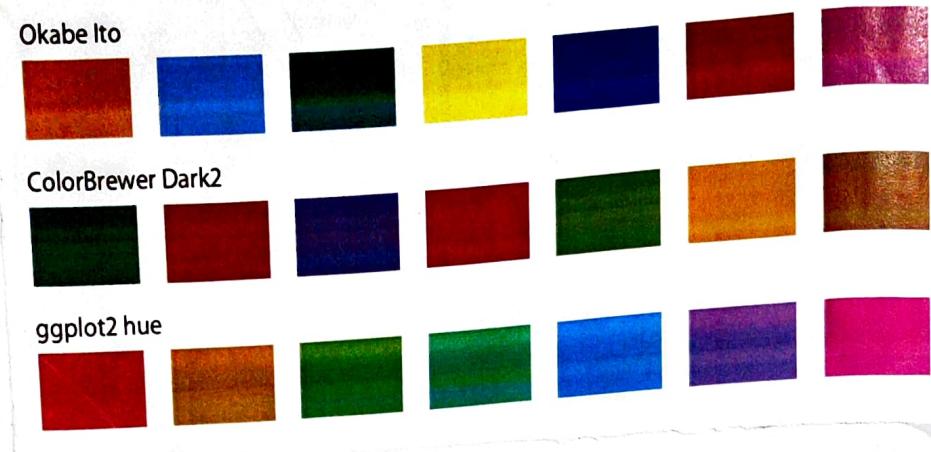
→ We use a qualitative color scale. Such a scale contains a finite set of specific colors that are chosen to look clearly distinct from each other while also being equivalent to each other.

→ The second condition requires that no one color should stand out relative to the others.

→ The colors should not create the impression of an order, as would be the case with a sequence of colors that get successively lighter. Such colors would create an apparent order among the items being colored, which by definition have no order.

→ Among colors that show differences in brightness, some colors often have higher contrast than others.

The figure below shows three types of qualitative color scales: the scales are Okabe Ito, colorBrewer Dark2 scales, ggplot2 hue scales.



### Colors to Represent Data Values

→ colors can also be used to represent quantitative data values, such as income, temperature or speed.  
→ we use a sequential color scale: such a scale contains a sequence of colors that clearly indicates which values are larger or smaller than which other ones and how distant two specific values are from each other.

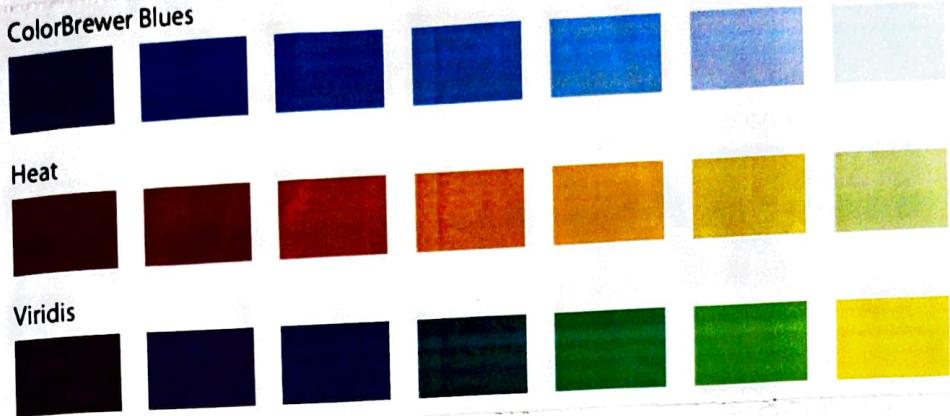
→ the second condition implies that the color scale needs to be perceived to vary uniformly across its entire range.

The different types of sequential color scales.

→ the colorBrewer Blues scale is a monochromatic scale that varies from dark to light blue.

→ The Heat and Viridis scales are multihue

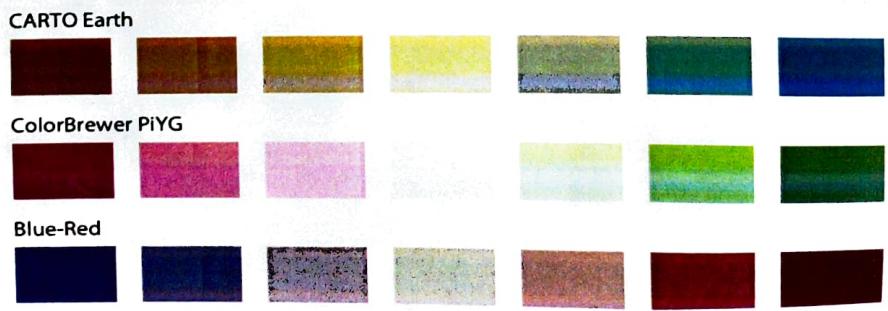
scales that vary from dark red to light yellow and from dark blue via green to light yellow respectively.



→ When we want to show how the data values vary across geographic regions. In this case, we can draw a map of the geographic regions and color them by the data values. Such maps are called choropleths.

→ When the dataset containing both positive and negative numbers. The appropriate color scale in this situation :-

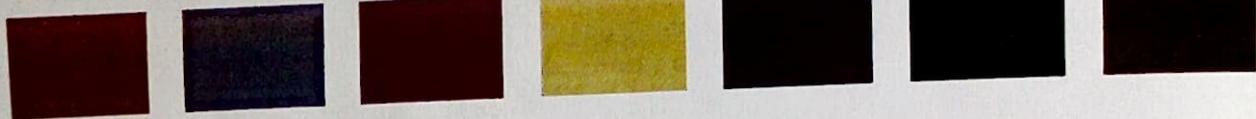
→ The diverging scales need to be balanced, so that the progression from light colors in the center to dark colors on the outside is approximately the same in either direction.



### 3-color as a Tool to Highlight:-

- colors can also be an effective tool to highlight specific elements in the data.
- These may be specific categories or values in the dataset that carry key information about the data.
- This can be achieved with accent color scales, which are color scales that contain both a set of subdued colors and a matching set of stronger, darker and/or more saturated colors.

Okabe Ito Accent



Grays with accents



ColorBrewer Accent



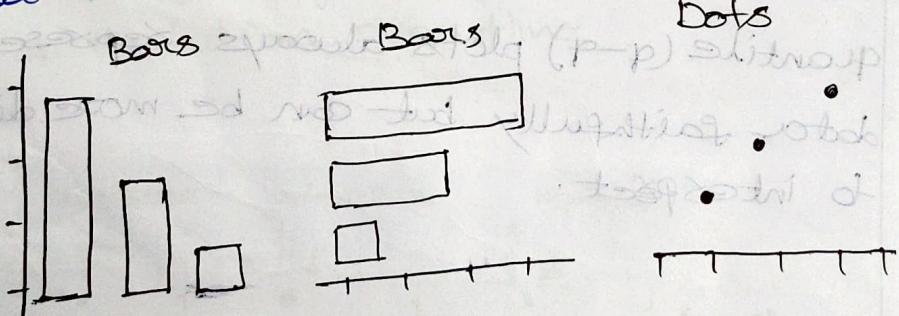
## → Directory of visualizations:

The directory of visualizations are mainly categorized into

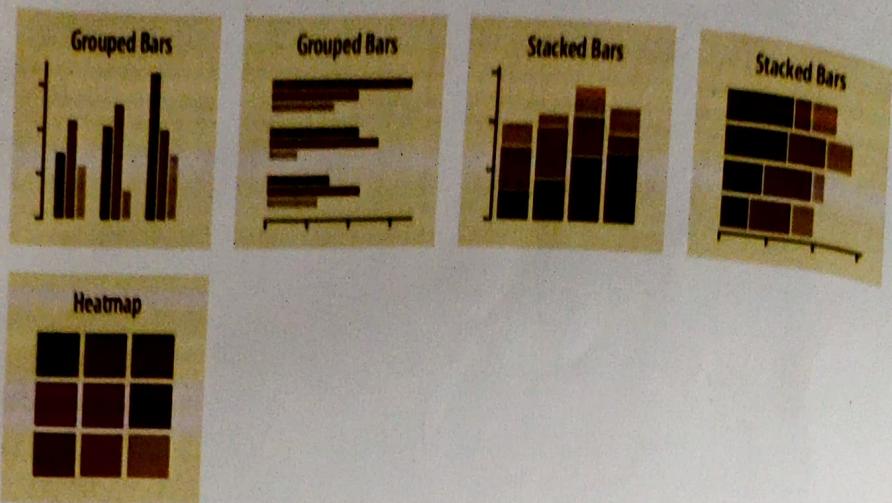
1. Amounts
2. Distributions
3. Proportions
4. X-y Relationships
5. Geospatial Data
6. Uncertainty.

### 1. Amounts:

The most common approach to visualizing amounts is using bars, either vertically or horizontally arranged. However, instead of using bars, we can also place dots at the locations where the corresponding bar could end. This is called a dot plot.



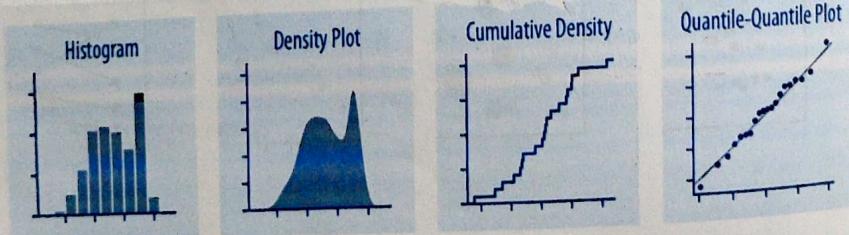
If there are two or more sets of categories for which we want to show amounts, we can group or stack the bars. We can also map the categories onto the x and y axes and show amounts by colors via a heatmap.



## Distributions:

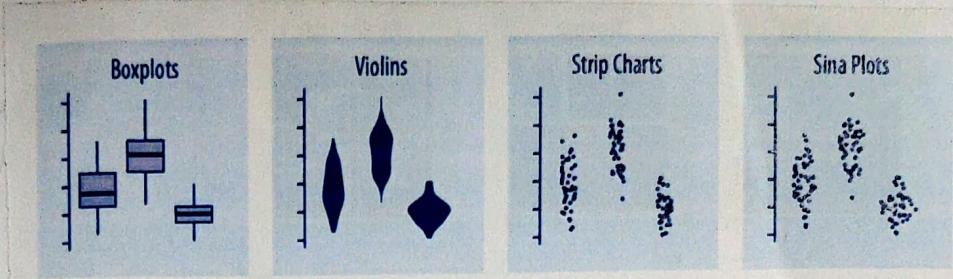
→ histograms and density plots provide the most intuitive visualizations of a distribution but both require arbitrary parameter choices and can be misleading.

→ cumulative densities and quantile-quantile ( $q-q$ ) plots always represent the data faithfully but can be more difficult to interpret.



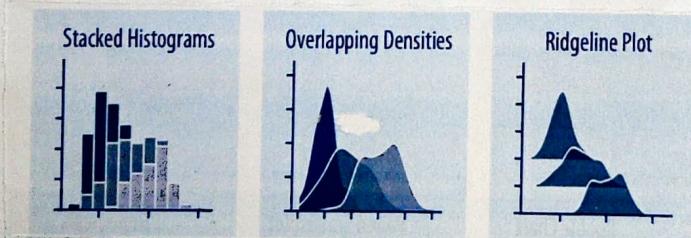
Pareto  
generated in 0.12 seconds

→ Boxplots, violin plots, strip charts, and sine plots are useful when we want to visualize many distributions at once and/or if we are primarily interested in overall shifts among the distributions.



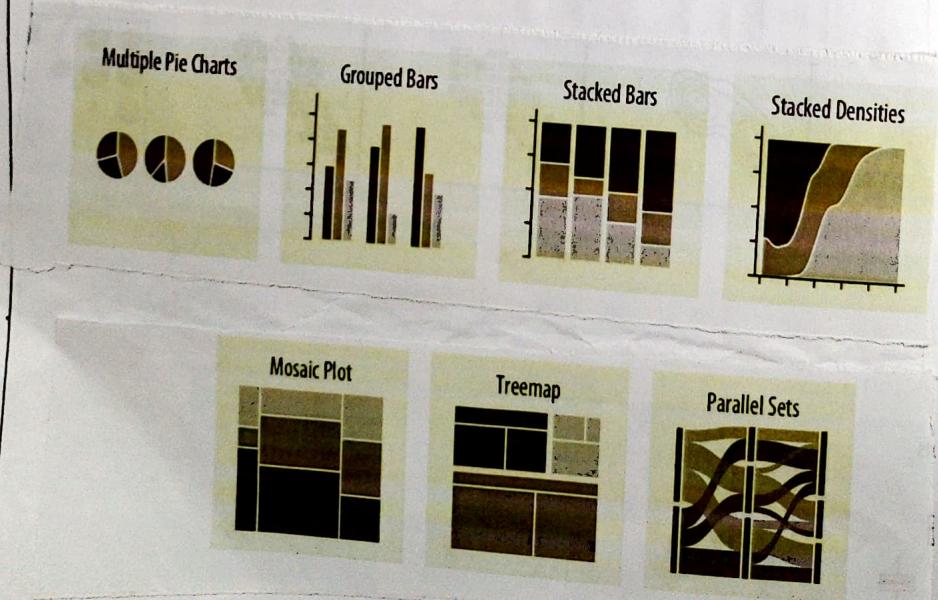
→ stacked histograms and overlapping densities allow a more in-depth comparison of a smaller number of distributions.

→ Ridgeline plots can be a useful alternative to violin plots and are often useful when visualizing very large amount of distribution in distribution over time.



### 3. Proportions:

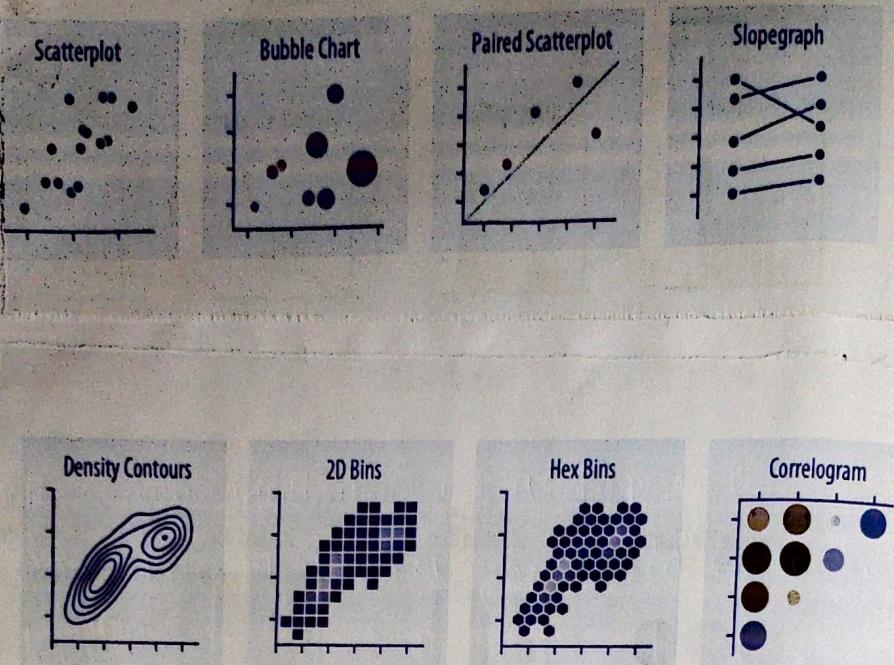
- proportions can be visualized as pie charts, side by side bars or stacked bars.
- when we visualize proportions with bars, the bars can be arranged either vertically or horizontally.
- pie charts emphasize that the individual parts add up to a whole and highlight simple fractions
- the individual pieces are more easily compared in side-by-side bars.
- stacked bar look can be useful when comparing multiple set of proportions.
- stacked densities are appropriate when the proportions change along a continuous variable.



- Mosaic plots assume that every level of one grouping variable can be combined with every level of another grouping variable.
- Treemaps work well even if the subdivisions of one group are entirely distinct from the subdivisions of another.
- parallel set work when there are more than two grouping variables.

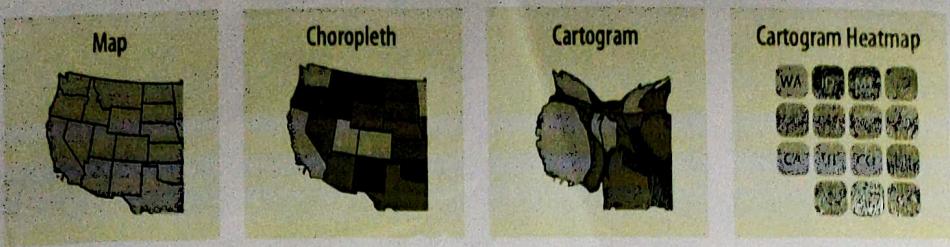
#### 4. X-y Relationships :-

- Scatterplots represent the archetypical visualization when we want to show one quantitative variable relative to another.
- We have three quantitative variables, we can map one onto the dot size creating a variant of the scatterplot called a bubble chart.
- Paired scatterplots are used for paired data, where the variables along the x and y axes are measured in the same units.
- Paired data can also be shown as a slopegraph of paired points connected by straight lines.
- For large numbers of points, regular scatterplots can become uninformative due to overplotting. In this case we use contourlines, 2D bins or hexbins which may provide an alternative.



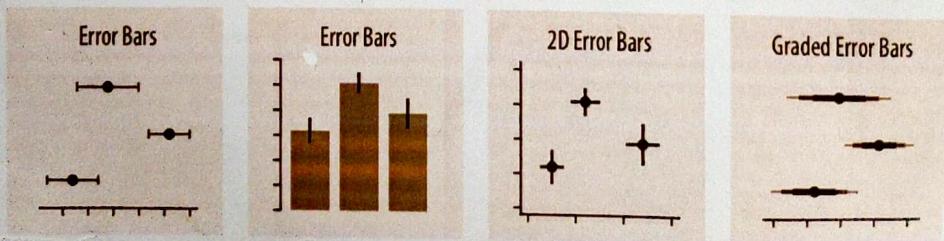
## 5) Geospatial Data

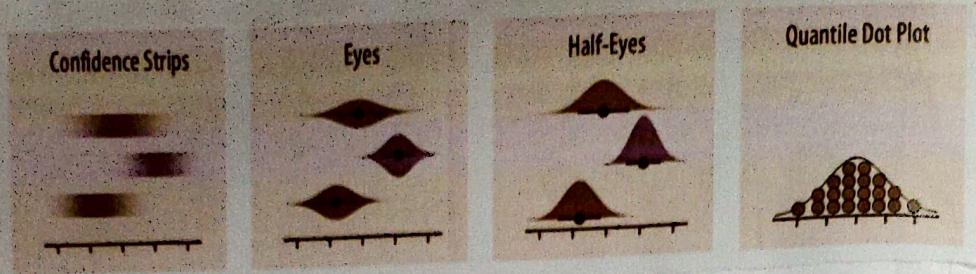
- The primary mode of showing geospatial data is in the form a map.
- A Map makes coordinates on the globe and projects them onto a flat surface such that shapes and distances on the globe are approximately represented by shapes and distances.
- We can show data values in different regions by coloring those regions in the map, according to the data. such map is called as choropleth.
- To distort the different regions according to some other quantity or simplify each region into a square, such visualizations are called cartograms.



## 6) Uncertainty

- Error bars are meant to indicate the range of likely values for some estimate or measurement.
- Graded error bars show multiple ranges at the same time, where each range corresponds to a different degree of confidence.
- Confidence strips provide a visual sense of uncertainty but are difficult to read accurately.
- Eyes and half eyes combine error bars with approaches to visualize distributions.
- A quantile dot plot can serve as an alternative visualization of an uncertainty distribution.





- 1) Define data visualization and its properties?
- 2) Explain different types of co ordinated systems with an example?
- 3) Explain the use of colors and color scales in data visualization?
- 4) Write different types of visualization methods?