



New Syllabus  
SPPU

As per the New Credit System Syllabus (2020 Course) of  
Savitribai Phule Pune University w.e.f. academic year 2023-2024

# Data Modeling and Visualization

(Code : 417522)

(Compulsory Subject)

Semester VII - Artificial Intelligence and Data Science

Dr. Sujatha Rao

TechKnowledge<sup>TM</sup>  
Publications

# 1

# Introduction to Data Modeling

## Syllabus

**Basic probability :** Discrete and continuous random variables, independence, covariance, central limit theorem, Chebyshev inequality, diverse continuous and discrete distributions. **Statistics, Parameter Estimation, and Fitting a Distribution :** Descriptive statistics, graphical statistics, method of moments, maximum likelihood estimation  
**Data Modeling Concepts** Understand and model subtypes and supertypes Understand and model hierarchical data Understand and model recursive relationships Understand and model historical data.

## 1.1 Basic Probability : Discrete and Continuous Random Variables

### 1.1.1 Random Variable in Statistics

- In probability, a real-valued function, defined over the sample space of a random experiment, is called a **random variable**. That is, the values of the random variable correspond to the outcomes of the random experiment. Random variables could be either discrete or continuous. A random variable's likely values may express the possible outcomes of an experiment, which is about to be performed or the possible outcomes of a preceding experiment whose existing value is unknown. They may also conceptually describe either the results of an "objectively" random process (like rolling a die) or the "subjective" randomness that appears from inadequate knowledge of a quantity.
- The domain of a random variable is a sample space, which is represented as the collection of possible outcomes of a random event. For instance, when a coin is tossed, only two possible outcomes are acknowledged such as heads or tails.

#### Random Variable Definition

- A random variable is a rule that assigns a numerical value to each outcome in a sample space. Random variables may be either discrete or continuous. A random variable is said to be discrete if it assumes only specified values in an interval. Otherwise, it is continuous. We generally denote the random variables with capital letters such as X and Y. When X takes values 1, 2, 3, ..., it is said to have a discrete random variable.
- As a function, a random variable is needed to be measured, which allows probabilities to be assigned to a set of potential values. It is obvious that the results depend on some physical variables which are not predictable. Say, when we toss a fair coin, the final result of happening to be heads or tails will depend on the possible physical conditions. We cannot predict which outcome will be noted.

**Variate**

A variate can be defined as a generalization of the random variable. It has the same properties as that of the random variables without stressing to any particular type of probabilistic experiment. It always obeys a particular probabilistic law.

- A variate is called discrete variate when that variate is not capable of assuming all the values in the provided range.
- If the variate is able to assume all the numerical values provided in the whole range, then it is called continuous variate.

**1.1.1(A) Random Variable Formula**

For a given set of data the mean and variance random variable is calculated by the formula. So, here we will define two major formulas :

## 2. Variance of random variable

1. Mean of random variable
  2. Variance of random variable
1. **Mean of random variable** : If  $X$  is the random variable and  $P$  is the respective probabilities, the mean of a random variable is defined by :

$$\text{Mean } (\mu) = \sum XP$$

where variable  $X$  consists of all possible values and  $P$  consist of respective probabilities.

2. **Variance of Random Variable** : The variance tells how much is the spread of random variable  $X$  around the mean value. The formula for the variance of a random variable is given by :

$$\text{Var}(X) = \sigma^2 = E(X^2) - [E(X)]^2$$

Where  $E(X^2) = \sum X^2P$  and  $E(X) = \sum XP$

**1.1.1(B) Functions of Random Variables**

- Let the random variable  $X$  assume the values  $x_1, x_2, \dots$  with corresponding probability  $P(x_1), P(x_2), \dots$  then the expected value of the random variable is given by :

$$\text{Expectation of } X, E(x) = \sum P(x).$$

- A new random variable  $Y$  can be stated by using a real Borel measurable function  $g : R \rightarrow R$ , to the results of a real-valued random variable  $X$ . That is,  $Y = f(X)$ . The cumulative distribution function of  $Y$  is then given by:

$$F_Y(y) = P(g(X) \leq y)$$

- If function  $g$  is invertible (say  $h = g^{-1}$ ) and is either increasing or decreasing, then the previous relationship can be extended to obtain :

$$F_Y(y) = P(g(X) \leq y) = \begin{cases} P(X \leq h(y)) = F_X(h(y)) & \text{If } h = g^{-1} \text{ increasing} \\ P(X \geq h(y)) = 1 - F_X(h(y)) & \text{If } h = g^{-1} \text{ decreasing} \end{cases}$$

- Now if we differentiate both the sides of the above expressions with respect to  $y$ , then the relation between the probability density functions can be found :

$$f_Y(y) = f_X(h(y))|dh(y)/dy|$$

**1.1.1(C) Random Variable and Probability Distribution**

The probability distribution of a random variable can be

- Theoretical listing of outcomes and probabilities of the outcomes.
- An experimental listing of outcomes associated with their observed relative frequencies.
- A subjective listing of outcomes associated with their subjective probabilities.
- The probability of a random variable  $X$  which takes the values  $x$  is defined as a probability function of  $X$  is denoted by  $f(x) = f(X = x)$
- A probability distribution always satisfies two conditions :

$$f(x) \geq 0$$

$$\sum f(x) = 1$$

**Random Variable Example**

- Ex. 1.1.1 :** Find the mean value for the continuous random variable,  $f(x) = x, 0 \leq x \leq 2$ .

**Soln. :**

Given :  $f(x) = x, 0 \leq x \leq 2$ .

The formula to find the mean value is :

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

$$E(X) = \int_0^2 x f(x) dx$$

$$E(X) = \int_0^2 x \cdot x dx$$

$$E(X) = \int_0^2 x^2 dx$$

$$E(X) = \left(\frac{x^3}{3}\right)_0^2$$

$$E(X) = \left(\frac{2^3}{3}\right) - \left(\frac{0^3}{3}\right)$$

$$E(X) = \left(\frac{8}{3}\right) - (0)$$

$$E(X) = \frac{8}{3}$$

Therefore, the mean of the continuous random variable,  $E(X) = \frac{8}{3}$

### 1.1.1(D) Types of Random Variable

As discussed in the introduction, there are two random variables, such as :

1. Discrete Random Variable
2. Continuous Random Variable

Let's understand these types of variables in detail along with suitable examples below.

#### 1. Discrete Random Variable

- A discrete random variable can take only a finite number of distinct values such as 0, 1, 2, 3, 4, ... and so on. The probability distribution of a random variable has a list of probabilities compared with each of its possible values known as probability mass function.
- In an analysis, let a person be chosen at random, and the person's height is demonstrated by a random variable. Logically the random variable is described as a function which relates the person to the person's height. Now in relation with the random variable, it is a probability distribution that enables the calculation of the probability that the height is in any subset of likely values, such as the likelihood that the height is between 175 and 185 cm, or the possibility that the height is either less than 145 or more than 180 cm. Now another random variable could be the person's age which could be either between 45 years to 50 years or less than 40 or more than 50.
- These are variables whose range is finite or countable. In particular, it means that their values can be listed, or arranged in a sequence. Examples include the number of jobs submitted to a printer, the number of errors, the number of error-free modules, the number of failed components, and so on. Discrete variables don't have to be integers. For example, the proportion of defective components in a lot of 100 can be 0, 1/100, 2/100, ..., 99/100, or 1. This variable assumes 101 different values, so it is discrete, although not an integer.

#### Discrete Example

The number of heads that appear during a series of five coin tosses is a discrete random variable that follows the binomial distribution. We can use that distribution to determine the likelihood of obtaining 0 to 5 heads. The Fig. 1.1.1 displays the probability for each possible outcome.

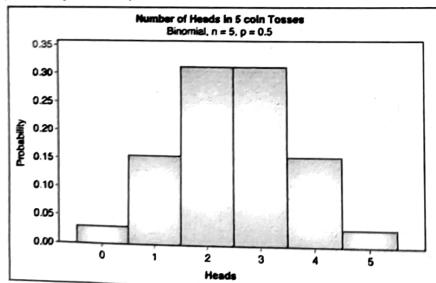


Fig. 1.1.1 : Discrete Distribution of data

### 2. Continuous Random Variable

- On the contrary, continuous random variables assume a whole interval of values. This could be a bounded interval  $(a, b)$ , or an unbounded interval such as  $(-\infty, \infty)$ ,  $(-\infty, b]$ , or  $(-\infty, +\infty)$ . Sometimes, it may be a union of several such intervals. Intervals are uncountable, therefore, all values of a random variable cannot be listed in this case. Examples of continuous variables include various times (software installation time, code execution time, connection time, waiting time, lifetime), also physical variables like weight, height, voltage, temperature, distance, the number of miles per gallon, etc.
- A numerically valued variable is said to be continuous if, in any unit of measurement, whenever it can take on the values a and b. If the random variable X can assume an infinite and uncountable set of values, it is said to be a continuous random variable. When X takes any value in a given interval  $(a, b)$ , it is said to be a continuous random variable in that interval.
- Formally, a continuous random variable is such whose cumulative distribution function is constant throughout. There are no "gaps" in between which would compare to numbers which have a limited probability of occurring.

#### Continuous Example

Body fat percentage is a continuous random variable. In preteen girls, it follows a lognormal distribution. Suppose a researcher needs to find participants with a body fat percentage between 20 and 24 percent. What is the likelihood that the next candidate will fall within that range ?

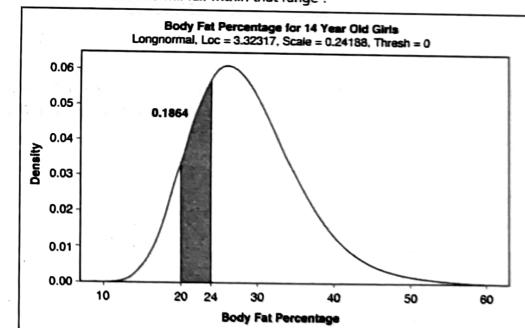


Fig. 1.1.2 : Continuous Distribution of data

### 1.2 Independence of Random Variables

**Definition :** Random variables X and Y are **independent** if

$$P_{X,Y}(x, y) = P_X(x)P_Y(y)$$

for all values of x and y. This means, events  $(X = x)$  and  $(Y = y)$  are independent for all x and y; in other words, variables X and Y take their values independently of each other.

In problems, to show independence of  $X$  and  $Y$ , we have to check whether the joint pmf factors into the product of marginal pmfs for all pairs  $x$  and  $y$ . To prove dependence, we only need to present one counterexample, a pair  $(x, y)$  with  $P(x, y) \neq P_X(x)P_Y(y)$ .

**Ex 1.2.1 :** A program consists of two modules. The number of errors,  $X$ , in the first module and the number of errors,  $Y$ , in the second module have the joint distribution,  $P(0, 0) = P(0, 1) = P(1, 0) = 0.2$ ,  $P(1, 1) = P(1, 2) = P(2, 1) = 0.1$ ,  $P(0, 2) = P(0, 3) = 0.05$ . Find (a) the marginal distributions of  $X$  and  $Y$ , (b) the probability of no errors in the first module, (c) the distribution of the total number of errors in the program. Also, (d) find out if errors in the two modules occur independently.

Soln. :

(a) It is convenient to organize the joint pmf of  $X$  and  $Y$  in a table. Adding row wise and column wise, we get the marginal pmfs,

(x, y)		y				$P_X(x)$
		0	1	2	3	
x	0	0.20	0.20	0.05	0.05	0.50
	1	0.20	0.10	0.10	0.10	0.50
$P_Y(y)$		0.40	0.30	0.15	0.15	1.00

This solves (a).

$$(b) P_X(0) = 0.50.$$

(c) Let  $Z = X + Y$  be the total number of errors. To find the distribution of  $Z$ , we first identify its possible values, then find the probability of each value. We see that  $Z$  can be as small as 0 and as large as 4. Then,

$$P_Z(0) = P(X + Y = 0) = P(X = 0 \cap Y = 0) = P(0, 0) = 0.20,$$

$$P_Z(1) = P(X = 0 \cap Y = 1) + P(X = 1 \cap Y = 0)$$

$$= P(0, 1) + P(1, 0) = 0.20 + 0.20 = 0.40,$$

$$P_Z(2) = P(0, 2) + P(1, 1) = 0.05 + 0.10 = 0.15,$$

$$P_Z(3) = P(0, 3) + P(1, 2) = 0.05 + 0.10 = 0.15,$$

$$P_Z(4) = P(1, 3) = 0.10.$$

It is a good check to verify that  $\sum_z P_Z(z) = 1$

(d) To decide on the independence of  $X$  and  $Y$ , check if their joint pmf factors into a product of marginal pmfs. We see that  $P_{XY}(0, 0) = 0.2$  indeed equals  $P_X(0)P_Y(0) = (0.5)(0.4)$ . Keep checking... Next,  $P_{XY}(0, 1) = 0.2$  whereas  $P_X(0)P_Y(1) = (0.5)(0.3) = 0.15$ . There is no need to check further. We found a pair of  $x$  and  $y$  that violates the formula for independent random variables. Therefore, the numbers of errors in two modules are dependent.

### 1.2.1 What are Independent Events?

- In Probability, the set of outcomes of an experiment is called events. There are different types of events such as independent events, dependent events, mutually exclusive events, and so on.
- If the probability of occurrence of an event  $A$  is not affected by the occurrence of another event  $B$ , then  $A$  and  $B$  are said to be independent events.
- Consider an example of rolling a die. If  $A$  is the event 'the number appearing is odd' and  $B$  be the event 'the number appearing is a multiple of 3', then

$$P(A) = \frac{3}{6} = \frac{1}{2} \text{ and } P(B) = \frac{2}{6} = \frac{1}{3}$$

- Also  $A$  and  $B$  is the event 'the number appearing is odd and a multiple of 3' so that

$$\begin{aligned} P(A \cap B) &= \frac{1}{6} \\ P(A|B) &= \frac{P(A \cap B)}{P(B)} \\ &= \frac{1}{6} \cdot \frac{3}{2} = \frac{1}{4} \end{aligned}$$

- $P(A) = P(A|B) = \frac{1}{2}$ , which implies that the occurrence of event  $B$  has not affected the probability of occurrence of the event  $A$ .
- If  $A$  and  $B$  are independent events, then  $P(A|B) = P(A)$
- Using the Multiplication rule of probability,  $P(A \cap B) = P(B) \cdot P(A|B)$

$$P(A \cap B) = P(B) \cdot P(A)$$

Note : A and B are two events associated with the same random experiment, then A and B are known as independent events if  $P(A \cap B) = P(B) \cdot P(A)$

### 1.3 Covariance

- Expectation, variance, and standard deviation characterize the distribution of a single random variable. Now we introduce measures of association of two random variables.

#### Definition

Covariance  $\sigma_{X,Y} = \text{Cov}(X, Y)$  is defined as

$$\begin{aligned} \text{Cov}(X, Y) &= E((X - EX)(Y - EY)) \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

It summarizes interrelation of two random variables

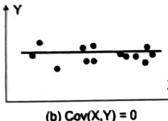
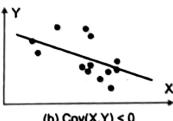
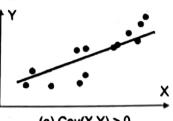


Fig. 1.3.1 : Positive, negative and zero covariance.

## Data Modeling and Visualization

- Covariance is the expected product of deviations of X and Y from their respective expectations. If  $\text{Cov}(X, Y) > 0$ , then positive deviations ( $X - \bar{X}$ ) are more likely to be multiplied by positive ( $Y - \bar{Y}$ ), and negative ( $X - \bar{X}$ ) are more likely to be multiplied by negative ( $Y - \bar{Y}$ ). In short, large X imply large Y, and small X imply small Y.
- These variables are positively correlated as shown in Fig. 1.3.1(a). Conversely,  $\text{Cov}(X, Y) < 0$  means that large X generally correspond to small Y and small X correspond to large Y. These variables are negatively correlated as shown in Fig. 1.3.1(b).
- If  $\text{Cov}(X, Y) = 0$ , we say that X and Y are uncorrelated. We see that independent variables are always uncorrelated. The reverse is not always true. There exist some variables that are uncorrelated but not independent. Notice that adding a constant does not affect the variables' variance or covariance. It shifts the whole distribution of X without changing its variability or degree of dependence of another variable.

Ex. 1.3.1 : Continuing Example 1.2.1 we compute.

x	$P_X(x)$	$x P_X(x)$	$x - \bar{X}$	$(x - \bar{X})^2 P_X(x)$
0	0.5	0	-0.5	0.125
1	0.5	0.5	0.5	0.125
		$u_X = 0.5$		$\sigma_X^2 = 0.25$

and (using the second method of computing variances)

y	$P_Y(y)$	$y P_Y(y)$	$y^2$	$y^2 P_Y(y)$
0	0.4	0	0	0
1	0.3	0.3	1	0.3
2	0.15	0.3	4	0.6
3	0.15	0.45	9	1.35
		$u_Y = 1.05$		$E(Y^2) = 2.25$

Soln. :

$$\text{Var}(X) = 0.25, \text{Var}(Y) = 2.25 - 1.05^2 = 1.1475, \text{Std}(X) = \sqrt{0.25} = 0.5 \text{ and } \text{Std}(Y) = \sqrt{1.1475} = 1.0712$$

$$E(XY) = \sum_{x,y} xy P(x, y) = (1)(1)(0.1) + (1)(2)(0.1) + (1)(3)(0.1) = 0.6$$

(the other five terms in this sum are 0). Therefore,

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) = 0.6 - (0.5)(1.05) = 0.075$$

and

$$\rho = \frac{\text{Cov}(X, Y)}{(\text{Std } X)(\text{Std } Y)} = \frac{0.075}{(0.5)(1.0712)} = 0.1400$$

Thus, the numbers of errors in two modules are *positively and not very strongly correlated*.

## Data Modeling and Visualization

## Population Covariance Formula

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N}$$

## Sample Covariance

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$$

Where,

$x_i$  = data value of x

$y_i$  = data value of y

$\bar{x}$  = mean of x

$\bar{y}$  = mean of y

N = number of data values.

## Covariance of X and Y

- If  $\text{Cov}(X, Y)$  is greater than zero, then we can say that the covariance for any two variables is positive and both the variables move in the same direction.
- If  $\text{Cov}(X, Y)$  is less than zero, then we can say that the covariance for any two variables is negative and both the variables move in the opposite direction.
- If  $\text{Cov}(X, Y)$  is zero, then we can say that there is no relation between two variables.

## 1.3.1 Correlation Coefficient Formula

We have already discussed covariance, which is the evaluation of changes between any two variables. Correlation estimates the depth of the relationship between variables. It is the estimated measure of covariance and is dimensionless. In other words, the correlation coefficient is a constant value always and does not have any units. The relationship between the correlation coefficient and covariance is given by;

$$\text{Correlation, } \rho(X,Y) = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y}$$

Where :

$\rho(X,Y)$  = Correlation between the variables X and Y

$\text{Cov}(X,Y)$  = Covariance between the variables X and Y

$\sigma_X$  = Standard deviation of the X variable

$\sigma_Y$  = Standard deviation of the Y variable

- Based on the value of correlation coefficient, we can estimate the type of correlation between the given two variables. Also, the graphical representation of correlation among two variables is given in the Fig. 1.3.2.

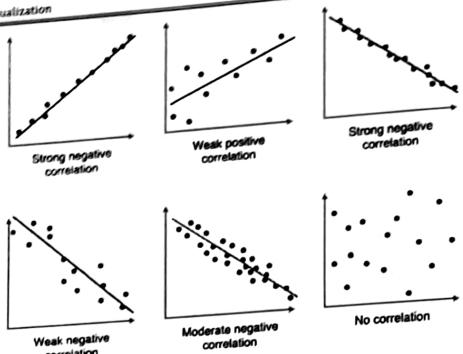


Fig. 1.3.2

### 1.3.2 Difference between Covariance and Correlation

Table 1.3.1 shows the comparison among covariance and correlation in brief.

Table 1.3.1

Sr. No.	Covariance	Correlation
1.	It is a measure to show the extent to which given two random variables change with respect to each other.	It is a measure used to describe how strongly the given two random variables are related to each other.
2.	It is a measure of correlation.	It is defined as the scaled form of covariance.
3.	The value of covariance lies between $-\infty$ and $+\infty$ .	The value of correlation lies between -1 and +1.
4.	It indicates the direction of the linear relationship between the given two variables.	It measures the direction and strength of the linear relationship between the given two variables.

### Covariance and Variance

Covariance and variance both are the terms used in statistics. Variance is the measure of spread of data around its mean value but covariance measures the relation between two random variables.

### Covariance Example

Below example helps in better understanding of the covariance of among two variables.

Ex. 1.3.2 : Calculate the coefficient of covariance for the following data :

X	2	8	18	20	28	30
Y	5	12	18	23	45	50

Soln. :

Number of observations = 6

Mean of X = 17.67

Mean of Y = 25.5

$$\begin{aligned} \text{Cov}(X, Y) &= \left(\frac{1}{6}\right) [(2 - 17.67)(5 - 25.5) + (8 - 17.67)(12 - 25.5) + (18 - 17.67)(18 - 25.5) \\ &\quad + (20 - 17.67)(23 - 25.5) + (28 - 17.67)(45 - 25.5) + (30 - 17.67)(50 - 25.5)] \\ &= 157.83 \end{aligned}$$

### 1.3.3 Covariance Matrix

In statistics and probability theory, a square matrix provides the covariance between each pair of components (or elements) of a given random vector is called a covariance matrix. Any covariance matrix is symmetric and positive semi-definite. The principal diagonal or main diagonal (sometimes a primary diagonal) of this matrix contains variances. That means the covariance of each element with itself. A covariance matrix is also known as the auto-covariance matrix, variance matrix, dispersion matrix, or variance-covariance matrix.

### 1.4 Central Limit Theorem

- The **central limit theorem**, which is a statistical theory, states that when a large sample size has a finite variance, the samples will be normally distributed, and the mean of samples will be approximately equal to the mean of the whole population.
- In other words, the central limit theorem states that for any population with mean and standard deviation, the distribution of the sample mean for sample size N has mean  $\mu$  and standard deviation  $\frac{\sigma}{\sqrt{n}}$ .
- As the sample size gets bigger and bigger, the mean of the sample will get closer to the actual population mean. If the sample size is small, the actual distribution of the data may or may not be normal, but as the sample size gets bigger, it can be approximated by a normal distribution. This statistical theory is useful in simplifying analysis while dealing with stock indexes and much more.
- The CLT can be applied to almost all types of probability distributions. But there are some exceptions. For example, if the population has a finite variance. Also, this theorem applies to independent, identically distributed variables. It can also be used to answer the question of how big a sample you want.
- Remember that as the sample size grows, the standard deviation of the sample average falls because it is the population standard deviation divided by the square root of the sample size. This theorem is an important topic in statistics. In many real-time applications, a certain random variable of interest is a sum of a large number of independent random variables. In such situations, we can use the CLT to justify using the normal distribution.

### Central Limit Theorem Statement

The central limit theorem states that whenever a random sample of size n is taken from any distribution with mean and variance, then the sample mean will be approximately a normal distribution with mean and variance.

## Data Modeling and Visualization

## Assumptions of the Central Limit Theorem

- The sample should be drawn randomly following the condition of randomisation.
- The samples drawn should be independent of each other. They should not influence the other samples.
- When the sampling is done without replacement, the sample size shouldn't exceed 10% of the total population.
- The sample size should be sufficiently large.

## Formula

The formula for the central limit theorem is given below:

## Central Limit Theorem for Sample Means,

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

## Proof

Consider  $x_1, x_2, x_3, \dots, x_n$  are independent and identically distributed with mean  $\mu$  and finite variance  $\sigma^2$ , then any random variable  $Z_n$  as,

$$Z_n = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

Here,

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

Then, the distribution function of  $Z_n$  converges to the standard normal distribution function as  $n$  increases without any bound.

Again, define a random variable  $U_i$  by

$$U_i = \frac{x_i - \mu}{\sigma}$$

$$E(U_i) = 0 \text{ and } V(U_i) = 1$$

Thus, the moment-generating function can be written as

$$M_U(t) = 1 + \frac{t^2}{2} + \frac{t^3}{3!} E(U_i^3) + \dots$$

Also,

$$Z_n = \sqrt{n} \left( \frac{\bar{x} - \mu}{\sigma} \right)$$

$$= \frac{1}{\sqrt{n}} \sum U_i$$

Since  $x_i$  are random independent variables,  $U_i$  are also independent.

## Data Modeling and Visualization

This implies,

$$M_U(t) = \left( 1 + \frac{t^2}{2n} + \frac{t^3}{3!} E(U_i^3) + \dots \right)^n$$

Or

$$\ln M_U(t) = n \ln \left( 1 + \frac{t^2}{2n} + \frac{t^3}{3!} E(U_i^3) + \dots \right)$$

As per Taylor series expansion:

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$$

$$\text{If } x = \frac{t^2}{2n} + \frac{t^3}{3!} E(U_i^3)$$

Then,

$$\ln(M_U(t)) = n \ln(1+x) = n \left( x - \frac{x^2}{2} + \frac{x^3}{3} - \dots \right)$$

Now, Multiply each term by  $n$ , and as  $n \rightarrow \infty$ , all terms but the first go to zero.

$$\lim_{n \rightarrow \infty} \ln(M_U(t)) = \frac{t^2}{2}$$

$$\text{and } \lim_{n \rightarrow \infty} (M_U(t)) = \exp\left(\frac{t^2}{2}\right)$$

Which is the moment-generating function for a standard normal random variable.

## Steps

The steps used to solve the problem of the central limit theorem that are either involving ' $>$ ' ' $<$ ' or "between" are as follows :

- The information about the mean, population size, standard deviation, sample size and a number that is associated with "greater than", "less than", or two numbers associated with both values for a range of "between" is identified from the problem.
- A graph with a centre as the mean is drawn.
- The formula is :
- The z-table is referred to find the 'z' value obtained in the previous step.
- Case 1 : Central limit theorem involving " $>$ ".**  
Subtract the z-score value from 0.5.

**Case 2 :** Central limit theorem involving "**<**".

Add 0.5 to the z-score value.

**Case 3 :** Central limit theorem involving "**between**".

Step 3 is executed.

6) The z-value is found along with the x-bar.

The last step is common to all three cases, that is, to convert the decimal obtained into a percentage.

#### 1.4.1 Solved Examples on Central Limit Theorem

**Ex. 1.4.1 :** 20 students are selected at random from a clinical psychology class; find the probability that their mean GPA is more than 5. If the average GPA scored by the entire batch is 4.91, the standard deviation is 0.72.

Soln. :

Here,

$$\text{Population mean} = \mu = 4.91$$

$$\text{Population standard deviation} = \sigma = 0.72$$

$$\text{Sample size} = n = 20 \text{ (which is less than 30)}$$

Since the sample size is smaller than 30, use the t-score instead of the z-score, even though the population standard deviation is known.

$$\sigma_x = \frac{\sigma}{\sqrt{n}}$$

Substituting the values, we have

$$\begin{aligned}\sigma_x &= \frac{0.72}{\sqrt{20}} \\ &= 0.161\end{aligned}$$

Now, find the t-score :

$$t = \frac{x - \mu}{\sigma_x}$$

For this problem, the raw score  $x = 5$

$$\begin{aligned}t &= \frac{5 - 4.91}{0.161} \\ &= 0.559\end{aligned}$$

Find the probability for the t value using the t-score table. The degree of freedom here would be :

$$Df = 20 - 1 = 19$$

$$P(t \leq 0.559) = 0.7087$$

$$P(t > 0.559) = 1 - 0.7087 = 0.2913$$

Thus, the probability that the score is more than 5 is 9.13 %.

**Ex. 1.4.2 :** The average weight of a water bottle is 30 kg, with a standard deviation of 1.5 kg. If a sample of 45 water bottles is selected at random from a consignment and their weights are measured, find the probability that the mean weight of the sample is less than 28 kg.

Soln. :

$$\text{Population mean} : \mu = 30 \text{ kg}$$

$$\text{Population standard deviation} : \sigma = 1.5 \text{ Kg}$$

$$\text{Sample size} : n = 45 \text{ (which is greater than 30)}$$

Using the z-score, we have

The sample standard deviation :

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

and

$$\sigma_{\bar{x}} = \frac{1.5}{\sqrt{45}}$$

$$= 6.7082$$

Find the z-score for the raw score of  $x = 28 \text{ kg}$

$$z = \frac{x - \mu}{\sigma_{\bar{x}}}$$

$$= (28 - 30) / 6.7082 = -0.2981$$

Using the z-score table OR normal CDF function on a statistical calculator,

$$P(z < -0.2981) = 0.3828$$

Thus, the probability that the weight of the cylinder is less than 28 kg is 38.28%.

**Ex. 1.4.3 :** The record of weights of the female population follows a normal distribution. Its mean and standard deviation are 65 kg and 14 kg, respectively. If a researcher considers the records of 50 females, then what would be the standard deviation of the chosen sample?

Soln. :

$$\text{Mean of the population} : \mu = 65 \text{ kg}$$

$$\text{The standard deviation of the population} : \sigma = 14 \text{ kg}$$

$$\text{Sample size} : n = 50$$

Standard deviation is given by,

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$$= \frac{14}{\sqrt{50}}$$

$$= \frac{14}{7.071}$$

$$\sigma_{\bar{x}} = 1.97$$

## 1.4.2 Applications of Central Limit Theorem

- 1) The sample distribution is assumed to be normal when the distribution is unknown or not normally distributed according to the central limit theorem. This method assumes that the given population is distributed normally. It helps in data analysis.
- 2) The sample mean deviation decreases as we increase the samples taken from the population, which helps in estimating the mean of the population more accurately.
- 3) The sample mean is used to create a range of values which likely includes the population mean.
- 4) The concept of the central limit theorem is used in election polls to estimate the percentage of people supporting a particular candidate as confidence intervals.
- 5) CLT is used in calculating the mean family income in a particular country.
- 6) It is used in rolling many identical, unbiased dice.
- 7) The probability distribution for the total distance covered in a random walk will approach a normal distribution.
- 8) Flipping many coins will result in a normal distribution for the total number of heads (or, equivalently total number of tails).
- 9) By looking at the sample distribution, CLT can tell whether the sample belongs to a particular population.
- 10) It enables us to make conclusions about the sample and population parameters and assists in constructing good machine-learning models.

## 1.5 Chebyshev Inequality

### What is Chebyshev's Inequality?

- Chebyshev's inequality is a probability theory that guarantees that within a specified range or distance from the mean, for a large range of probability distributions, no more than a specific fraction of values will be present. In other words, only a definite fraction of values will be found within a specific distance from the mean of a distribution.
- The formula for the fraction for which no more than a certain number of values can exceed is  $1/K^2$ ; in other words,  $1/K^2$  of a distribution's values can be more than or equal to  $K$  standard deviations away from the mean of the distribution. Further, it also holds that  $1 - (1/K^2)$  of a distribution's values must be within, but not including,  $K$  standard deviations away from the mean of the distribution.
- Chebyshev's inequality is similar to the 68-95-99.7 rule; however, the latter rule only applies to normal distributions. Chebyshev's inequality is broader; it can be applied to any distribution so long as the distribution includes a defined variance and mean.
- Chebyshev's inequality states that within two standard deviations away from the mean contains 75% of the values, and within three standard deviations away from the mean contains 88.9% of the values. It holds for a wide range of probability distributions, not only the normal distribution.
- However, when applied to the normal distribution, Chebyshev's inequality is less precise than the 68-95-99.7 rule; yet, it is important to keep in mind that the theory applies to a far broader range of distributions. It should be noted that standard deviations equal to or less than one are not valid for Chebyshev's inequality formula.

## Chebyshev's Inequality History

Chebyshev's inequality was proven by Pafnuty Chebyshev, a Russian mathematician, in 1867. It was stated earlier by French statistician Irénée-Jules Bienaymé in 1853; however, there was no proof for the theory made with the statement. After Pafnuty Chebyshev proved Chebyshev's inequality, one of his students, Andrey Markov, provided another proof for the theory in 1884.

### Chebyshev's Inequality Statement

Let  $X$  be a random variable with a finite mean denoted as  $\mu$  and a finite non-zero variance, which is denoted as  $\sigma^2$ , for any real number,  $K > 0$ .

$$P(|X - \mu| \geq K) \leq (\sigma^2 / K^2)$$

### Practical Example

Assume that an asset is picked from a population of assets at random. The average return of the population of assets is 12%, and the standard deviation of the population of assets is 5%. To calculate the probability that an asset picked at random from this population, which has a return less than 4% or greater than 20%, Chebyshev's inequality can be applied.

Since there is a limited amount of information and only the mean and standard deviation of a distribution is given, the exact probability of this scenario cannot be determined; thus, Chebyshev's inequality is applied. Below is the application of the theory :

$$\begin{aligned} |X - \mu| &\geq K \\ P(|X - \mu| \geq K) &\leq \frac{\sigma^2}{K^2} = \frac{5\%^2}{8\%^2} \\ P(|X - \mu| \geq K) &\leq \frac{\sigma^2}{K^2} = 39.06\% \end{aligned}$$

Where :

Standard deviation : 5%

Mean : 12%

K : 8%

Thus, the probability of an asset's return to be less than 4% or greater than 20% from the population of assets, which has a mean return of 12% with a standard deviation of 5%, is less than 39.06%, according to Chebyshev's inequality.

## 1.6 Diverse Continuous and Discrete Distributions

### What is Probability ?

Probability denotes the possibility of something happening. It is a mathematical concept that predicts how likely events are to occur. The probability values are expressed between 0 and 1. The definition of probability is the degree to which something is likely to occur. This fundamental theory of probability is also applied to probability distributions.

## What are Probability Distributions ?

A probability distribution is a statistical function that describes all the possible values and probabilities for a random variable within a given range. This range will be bound by the minimum and maximum possible values, but where the possible value would be plotted on the probability distribution will be determined by a number of factors. The mean (average), standard deviation, skewness, and kurtosis of the distribution are among these factors.

## What is a continuous distribution ?

A continuous distribution describes the probabilities of the possible values of a continuous random variable. A continuous random variable is a random variable with a set of possible values (known as the range) that is infinite and uncountable.

Probabilities of continuous random variables ( $X$ ) are defined as the area under the curve of its PDF. Thus, only ranges of values can have a nonzero probability. The probability that a continuous random variable equals some value is always zero.

Now, look at some varieties of the continuous probability distribution.

## Normal Distribution

Normal Distribution is one of the most basic continuous distribution types. Gaussian distribution is another name for it. Around its mean value, this probability distribution is symmetrical. It also demonstrates that data close to the mean occurs more frequently than data far from it. Here, the mean is 0, and the variance is a finite value.

In the example, you generated 100 random variables ranging from 1 to 50. After that, you created a function to define the normal distribution formula to calculate the probability density function. Then, you have plotted the data points and probability density function against X-axis and Y-axis, respectively.

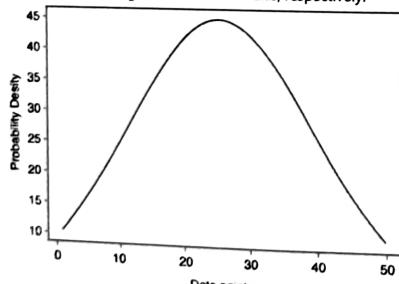


Fig. 1.6.1 : PDF

## Continuous Uniform Distribution

In continuous uniform distribution, all outcomes are equally possible. Each variable has the same chance of being hit as a result. Random variables are spaced evenly in this symmetric probabilistic distribution, with a  $1/(b-a)$  probability.

## Log-Normal Distribution

The random variables whose logarithm values follow a normal distribution are plotted using this distribution. Take a look at the random variables  $X$  and  $Y$ . The variable represented in this distribution is  $Y = \ln(X)$ , where  $\ln$  denotes the natural logarithm of  $X$  values.

The size distribution of rain droplets can be plotted using log normal distribution.

## Exponential Distribution

- In a Poisson process, an exponential distribution is a continuous probability distribution that describes the time between events (success, failure, arrival, etc.).
- You can see in the below example how to get random samples of exponential distribution and return Numpy array samples by using the `numpy.random.exponential()` method.
- In Statistics, the **probability distribution** gives the possibility of each outcome of a random experiment or event. It provides the probabilities of different possible occurrences. Also read, events in probability, here.
- To recall, the **probability is a measure of uncertainty of various phenomena**. Like, if you throw a dice, the possible outcomes of it, is defined by the probability. This distribution could be defined with any random experiments, whose outcome is not sure or could not be predicted. Let us discuss now its definition, function, formula and its types here, along with how to create a table of probability based on random variables.

## Example of the distribution of weights

The continuous normal distribution can describe the distribution of weight of adult males. For example, you can calculate the probability that a man weighs between 160 and 170 pounds.

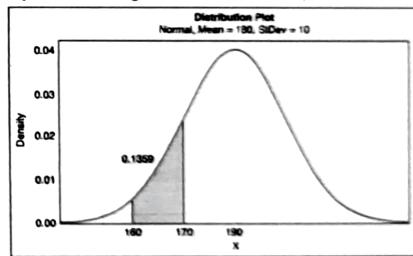


Fig. 1.6.2 : Distribution plot of the weight of adult males

The shaded region under the curve in this example represents the range from 160 and 170 pounds. The area of this range is 0.136; therefore, the probability that a randomly selected man weighs between 160 and 170 pounds is 13.6%. The entire area under the curve equals 1.0.

However, the probability that  $X$  is exactly equal to some value is always zero because the area under the curve at a single point, which has no width, is zero. For example, the probability that a man weighs exactly 190 pounds to infinite precision is zero. You could calculate a nonzero probability that a man weighs more than 190 pounds, or less than 190 pounds, or between 189.9 and 190.1 pounds, but the probability that he weighs exactly 190 pounds is zero.

### What is a discrete distribution?

A discrete distribution describes the probability of occurrence of each value of a discrete random variable. A discrete random variable is a random variable that has countable values, such as a list of non-negative integers.

With a discrete probability distribution, each possible value of the discrete random variable can be associated with a non-zero probability. Thus, a discrete probability distribution is often presented in tabular form.

### Example of the number of customer complaints

With a discrete distribution, unlike with a continuous distribution, you can calculate the probability that  $X$  is exactly equal to some value. For example, you can use the discrete Poisson distribution to describe the number of customer complaints within a day. Suppose the average number of complaints per day is 10 and you want to know the probability of receiving 5, 10, and 15 customer complaints in a day.

x	P(X = x)
5	0.037833
10	0.125110
15	0.034718

You can also view a discrete distribution on a distribution plot to see the probabilities between ranges

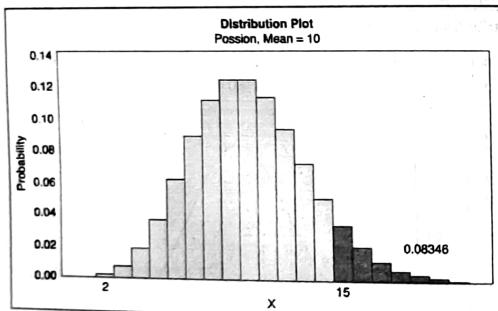


Fig. 1.6.3 : Distribution plot of the number of customer complaints

The shaded bars in this example represent the number of occurrences when the daily customer complaints is 15 or more. The height of the bars sums to 0.08346; therefore, the probability that the number of calls per day is 15 or more is 8.35%.

The most commonly used families of discrete distributions.

#### 1.6.1 Bernoulli Distribution

- The simplest random variable (excluding non-random ones!) takes just two possible values. Call them 0 and 1.
- Definition :** A random variable with two possible values, 0 and 1, is called a Bernoulli variable, its distribution is Bernoulli distribution, and any experiment with a binary outcome is called a Bernoulli trial.

- This distribution is named after a Swiss mathematician Jacob Bernoulli (1654-1705) who discovered not only Bernoulli but also Binomial distribution. Good or defective components, parts that pass or fail tests, transmitted or lost signals, working or malfunctioning hardware, benign or malicious attachments, sites that contain or do not contain a keyword, girls and boys, heads and tails, and so on, are examples of Bernoulli trials. All these experiments fit the same Bernoulli model, where we shall use generic names for the two outcomes: "successes" and "failures." These are nothing but commonly used generic names; in fact, successes do not have to be good, and failures do not have to be bad.
- If  $P(1) = p$  is the probability of a success, then  $P(0) = q = 1 - p$  is the probability of a failure. We can then compute the expectation and variance as

$$E(X) = \sum_x P(x) = (0)(1-p) + (1)p = p$$

$$\begin{aligned} \text{Var}(X) &= \sum_x (x - p)^2 P(x) \\ &= (0 - p)^2 (1-p) + (1 - p)^2 p \\ &= p(1-p)(p+1-p) = p(1-p) \end{aligned}$$

Bernoulli distribution

$p = \text{Probability of success}$
$P(x) = \begin{cases} q = 1 - p & \text{if } x = 0 \\ p & \text{if } x = 1 \end{cases}$
$E(X) = p$
$\text{Var}(X) = pq$

- In fact, we see that there is a whole family of Bernoulli distributions, indexed by a parameter  $p$ . Every  $p$  between 0 and 1 defines another Bernoulli distribution. The distribution with  $p = 0.5$  carries the highest level of uncertainty because  $\text{Var}(X) = pq$  is maximized by  $p = q = 0.5$ . Distributions with lower or higher  $p$  have lower variances. Extreme parameters  $p = 0$  and  $p = 1$  define non-random variables 0 and 1, respectively; their variance is 0.

#### 1.6.2 Binomial Distribution

- Now consider a sequence of independent Bernoulli trials and count the number of successes in it. This may be the number of defective computers in a shipment, the number of updated files in a folder, the number of girls in a family, the number of e-mails with attachments, etc
- Definition :** A variable described as the number of successes in a sequence of independent Bernoulli trials has Binomial distribution. Its parameters are  $n$ , the number of trials, and  $p$ , the probability of success.

Remark : "Binomial" can be translated as "two numbers," bi meaning "two" and nom meaning "a number," thus reflecting the concept of binary outcomes.

- Binomial probability mass function is

$$P(x) = P(X = x) = \binom{n}{x} p^x q^{n-x}, x = 0, 1, \dots, n \quad \dots(1.6.1)$$

which is the probability of exactly  $x$  successes in  $n$  trials. In this formula,  $p^x$  is the probability of  $x$  successes, probabilities being multiplied due to independence of trials. Also,  $q^{n-x}$  is the probability of the remaining  $(n - x)$  trials being failures. Finally,  $\binom{n}{x} = \frac{n!}{x!(n-x)!}$  is the number of elements of the sample space  $\Omega$  that form the event  $\{X = x\}$ . This is the number of possible orderings of  $x$  successes and  $(n - x)$  failures among  $n$  trials, and it is computed as  $C(n, x)$ . Due to a somewhat complicated form of Equation (1.6.1), practitioners use a table of Binomial distribution.

$$E(X) = E(X_1 + \dots + X_n) = E(X_1) + \dots + E(X_n) = p + \dots + p = np$$

and

$$\text{Var}(X) = \text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n) = npq$$

Binomial distribution

$n$ = Number of trials
$P$ = Probability of success
$P(X) = \binom{n}{x} p^x q^{n-x}$
$E(X) = np$
$\text{Var}(X) = npq$

### 1.6.3 Geometric Distribution

- Again, consider a sequence of independent Bernoulli trials. Each trial results in a "success" or a "failure."
  - Definition :** The number of Bernoulli trials needed to get the first success has Geometric distribution.
  - Example :** A search engine goes through a list of sites looking for a given key phrase. Suppose the search terminates as soon as the key phrase is found. The number of sites visited is Geometric.
- A hiring manager interviews candidates, one by one, to fill a vacancy. The number of candidates interviewed until one candidate receives an offer has Geometric distribution.
- Geometric random variables can take any integer value from 1 to infinity, because one needs at least 1 trial to have the first success, and the number of trials needed is not limited by any specific number. (For example, there is no guarantee that among the first 10 coin tosses there will be at least one head.) The only parameter is  $p$ , the probability of a "success." Geometric probability mass function has the form  $P(x) = P(\text{the 1st success occurs on the } x^{\text{th}} \text{ trial}) = (1 - p)^{x-1} p$ ,  $x = 1, 2, \dots$ , which is the probability of  $(x - 1)$  failures followed by one success.

Geometric distribution

$P$ = Probability of success
$P(X) = (1 - p)^{x-1} p$ , $x = 1, 2, \dots$
$E(X) = \frac{1}{P}$
$\text{Var}(X) = \frac{1-p}{P^2}$

### 1.6.4 Poisson Distribution

The next distribution is related to a concept of rare events or Poisson events. Essentially it means that two such events are extremely unlikely to occur simultaneously or within a very short period of time. Arrivals of jobs, telephone calls, e-mail messages, traffic accidents, network blackouts, virus attacks, errors in software, floods, and earthquakes are examples of rare events.

#### Definition

The number of rare events occurring within a fixed period of time has Poisson distribution. This distribution bears the name of a famous French mathematician Simeon-Denis Poisson (1781–1840).

Poisson distribution	$\lambda$ = Frequency, average number of events
	$P(X) = e^{-\lambda} \frac{\lambda^x}{x!}$ , $x = 1, 2, \dots$
	$E(X) = \lambda$
	$\text{Var}(X) = \lambda$

**Ex. 1.6.1 :** Customers of an internet service provider initiate new accounts at the average rate of 10 accounts per day. (a) What is the probability that more than 8 new accounts will be initiated today? (b) What is the probability that more than 16 accounts will be initiated within 2 days?

Soln. :

- (a) New account initiations qualify as rare events because no two customers open accounts simultaneously. Then the number  $X$  of today's new accounts has Poisson distribution with parameter  $\lambda = 10$ .
- $$\therefore P(X > 8) = 1 - F_8(8) = 1 - 0.333 = 0.667.$$
- (b) The number of accounts,  $Y$ , opened within 2 days does not equal  $2X$ . Rather,  $Y$  is another Poisson random variable whose parameter equals 20. Indeed, the parameter is the average number of rare events, which, over the period of two days, doubles the one-day average. With  $\lambda = 20$ ,
- $$P(Y > 16) = 1 - F_{16}(16) = 1 - 0.221 = 0.779.$$

Thus the covariance of these two variables is denoted by  $\text{Cov}(X, Y)$ . The formula is given below for both population covariance and sample covariance.

### 1.6.5 Probability Distribution of Random Variables

A random variable has a probability distribution, which defines the probability of its unknown values. Random variables can be discrete (not constant) or continuous or both. That means it takes any of a designated finite or countable list of values, provided with a probability mass function feature of the random variable's probability distribution or can take any numerical value in an interval or set of intervals. Through a probability density function that is representative of the random variable's probability distribution or it can be a combination of both discrete and continuous.

Two random variables with equal probability distribution can yet vary with respect to their relationships with other random variables or whether they are independent of these. The recognition of a random variable, which means, the outcomes of randomly choosing values as per the variable's probability distribution function, are called random variates.

which is the probability of exactly  $x$  successes in  $n$  trials. In this formula,  $p^x$  is the probability of  $x$  successes, probabilities being multiplied due to independence of trials. Also,  $q^{n-x}$  is the probability of the remaining  $(n - x)$  trials being failures. Finally,  $\binom{n}{x} = \frac{n!}{x!(n-x)!}$  is the number of elements of the sample space  $\Omega$  that form the event  $\{X = x\}$ . This is the number of possible orderings of  $x$  successes and  $(n - x)$  failures among  $n$  trials, and it is computed as  $C(n, x)$ . Due to a somewhat complicated form of Equation (1.6.1), practitioners use a table of Binomial distribution.

$$E(X) = E(X_1 + \dots + X_n) = E(X_1) + \dots + E(X_n) = p + \dots + p = np$$

and

$$\text{Var}(X) = \text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n) = npq$$

Binomial distribution

$n$	= Number of trials
$P$	= Probability of success
$P(X) = \binom{n}{x} p^x q^{n-x}$	
$E(X) = np$	
$\text{Var}(X) = npq$	

### 1.6.3 Geometric Distribution

- Again, consider a sequence of independent Bernoulli trials. Each trial results in a "success" or a "failure."
- Definition :** The number of Bernoulli trials needed to get the first success has Geometric distribution.
- Example :** A search engine goes through a list of sites looking for a given key phrase. Suppose the search terminates as soon as the key phrase is found. The number of sites visited is Geometric.
- A hiring manager interviews candidates, one by one, to fill a vacancy. The number of candidates interviewed until one candidate receives an offer has Geometric distribution.
- Geometric random variables can take any integer value from 1 to infinity, because one needs at least 1 trial to have the first success, and the number of trials needed is not limited by any specific number. (For example, there is no guarantee that among the first 10 coin tosses there will be at least one head.) The only parameter is  $p$ , the probability of a "success." Geometric probability mass function has the form  $P(x) = P(\text{the 1st success occurs on the } x^{\text{th}} \text{ trial}) = (1 - p)^{x-1} p$ ,  $x = 1, 2, \dots$ , which is the probability of  $(x - 1)$  failures followed by one success.

Geometric distribution

$P$	= Probability of success
$P(X) = (1 - p)^{x-1} p$	, $x = 1, 2, \dots$
$E(X) = \frac{1}{P}$	
$\text{Var}(X) = \frac{1-p}{P^2}$	

### 1.6.4 Poisson Distribution

The next distribution is related to a concept of rare events or Poisson events. Essentially it means that two such events are extremely unlikely to occur simultaneously or within a very short period of time. Arrivals of jobs, telephone calls, e-mail messages, traffic accidents, network blackouts, virus attacks, errors in software, floods, and earthquakes are examples of rare events.

#### Definition

The number of rare events occurring within a fixed period of time has Poisson distribution. This distribution bears the name of a famous French mathematician Simeon-Denis Poisson (1781–1840).

Poisson distribution

$\lambda$	= Frequency, average number of events
$P(X) = e^{-\lambda} \frac{\lambda^x}{x!}$	, $x = 1, 2, \dots$
$E(X) = \lambda$	
$\text{Var}(X) = \lambda$	

- Ex. 1.6.1 :** Customers of an internet service provider initiate new accounts at the average rate of 10 accounts per day. (a) What is the probability that more than 8 new accounts will be initiated today? (b) What is the probability that more than 16 accounts will be initiated within 2 days?

Soln. :

- (a) New account initiations qualify as rare events because no two customers open accounts simultaneously. Then the number  $X$  of today's new accounts has Poisson distribution with parameter  $\lambda = 10$ .  
 $\therefore P(X > 8) = 1 - F_8(8) = 1 - 0.333 = 0.667$ .
- (b) The number of accounts,  $Y$ , opened within 2 days does not equal  $2X$ . Rather,  $Y$  is another Poisson random variable whose parameter equals 20. Indeed, the parameter is the average number of rare events, which, over the period of two days, doubles the one-day average. With  $\lambda = 20$ ,  
 $P(Y > 16) = 1 - F_{16}(16) = 1 - 0.221 = 0.779$ .

Thus the covariance of these two variables is denoted by  $\text{Cov}(X, Y)$ . The formula is given below for both population covariance and sample covariance.

### 1.6.5 Probability Distribution of Random Variables

A random variable has a probability distribution, which defines the probability of its unknown values. Random variables can be discrete (not constant) or continuous or both. That means it takes any of a designated finite or countable list of values, provided with a probability mass function feature of the random variable's probability distribution or can take any numerical value in an interval or set of intervals. Through a probability density function that is representative of the random variable's probability distribution or it can be a combination of both discrete and continuous.

Two random variables with equal probability distribution can yet vary with respect to their relationships with other random variables or whether they are independent of these. The recognition of a random variable, which means, the outcomes of randomly choosing values as per the variable's probability distribution function, are called random variates.

**Probability Distribution Formulas**

<b>Binomial Distribution</b>	$P(X) = {}^n C_x a^x b^{n-x}$ Where, a = probability of success b = probability of failure n = number of trials x = random variable denoting success
<b>Cumulative Distribution Function</b>	$F_X(x) = \int_{-\infty}^x f_X(t) dt$
<b>Discrete Probability Distribution</b>	$P(x) = \frac{n!}{r!(n-r)!} \cdot p^r (1-p)^{n-r}$ $P(x) = C(n, r) \cdot p^r (1-p)^{n-r}$

**Solved Examples**

**Ex. 1.6.2 :** A coin is tossed twice. X is the random variable of the number of heads obtained. What is the probability distribution of  $X$ ?

**Soln. :**

First write, the value of  $X = 0, 1$  and  $2$ , as the possibility are there that

No head comes

One head and one tail comes

And head comes in both the coins

Now the probability distribution could be written as;

$$P(X = 0) = P(\text{Tail + Tail}) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

$$P(X = 1) = P(\text{Head + Tail}) \text{ or } P(\text{Tail + Head}) = \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2}$$

$$P(X = 2) = P(\text{Head + Head}) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

We can put these values in tabular form;

X	0	1	2
P(X)	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

**Ex. 1.6.3 :** The weight of a pot of water chosen is a continuous random variable. The following table gives the weight in kg of 100 containers recently filled by the water purifier. It records the observed values of the continuous random variable and their corresponding frequencies. Find the probability or chances for each weight category.

Weight W	Number of Containers
0.900 – 0.925	1
0.925 – 0.950	7
0.950 – 0.975	25
0.975 – 1.000	32
1.000 – 1.025	30
1.025 – 1.050	5
<b>Total</b>	<b>100</b>

We first divide the number of containers in each weight category by 100 to give the probabilities.

Weight W	Number of Containers	Probability
0.900 – 0.925	1	0.01
0.925 – 0.950	7	0.07
0.950 – 0.975	25	0.25
0.975 – 1.000	32	0.32
1.000 – 1.025	30	0.30
1.025 – 1.050	5	0.05
<b>Total</b>	<b>100</b>	<b>1.00</b>

**1.7 Statistics, Parameter Estimation and Fitting a Distribution****1.7.1 Statistics**

Statistics is a branch of mathematics that deals with the collection, analysis, interpretation, presentation, and organization of data. It provides methods and techniques to extract meaningful information from data, make inferences, and draw conclusions about populations based on sample data. There are two main branches of statistics :

- a. **Descriptive Statistics :** Descriptive statistics involve summarizing and presenting data in a meaningful way. Measures such as mean, median, mode, standard deviation, and percentiles are commonly used to describe the central tendency and variability of the data.
- b. **Inferential Statistics :** Inferential statistics use sample data to make inferences or predictions about the larger population from which the sample was drawn. Techniques like hypothesis testing, confidence intervals, and regression analysis fall under inferential statistics.

## 1.7.2 Parameter Estimation

Parameter estimation is a fundamental aspect of statistics that involves determining the unknown parameters of a statistical model based on observed data. In many real-world scenarios, data is generated from a specific statistical distribution with unknown parameters. The goal of parameter estimation is to find the best estimates for these unknown parameters, often denoted as  $\theta$ , so that the statistical model closely represents the observed data.

There are two main approaches to parameter estimation :

- Method of Moments** : In this approach, the parameters of the distribution are estimated by equating sample moments (such as the mean, variance, etc.) to the theoretical moments of the distribution.
- Maximum Likelihood Estimation (MLE)** : MLE is a powerful and widely used method for parameter estimation. It involves finding the parameter values that maximize the likelihood function, which measures the probability of obtaining the observed data given a particular set of parameter values.

Parameter estimation is the process of computing a model's parameter values from measured data. You can apply parameter estimation to different types of mathematical models, including statistical models, parametric dynamic models, and data-based Simulink® models.

### 1. Statistical Models

- Engineers and scientists apply parameter estimation to statistical models to estimate :
  - Parameters of a probability distribution, such as the mean and standard deviation of a normal distribution
  - Regression coefficients of a regression model
- The Statistics and Machine Learning Toolbox™ supports these and similar parameter estimation tasks with more than 40 different probability distributions, including normal, Weibull, gamma, generalized Pareto, and Poisson. The toolbox also supports linear and nonlinear regression.

### 2. Dynamic Models

- Engineers apply parameter estimation to dynamic models to compute :
  - Coefficients of transfer functions, including ARX, ARMAX, Box-Jenkins, and output-error models
  - Entries in state-space matrices
  - Coefficients of ODEs or well-structured systems with parameter constraints (grey-box system identification)
  - Regression coefficients, saturation levels, or dead-zone limits for nonlinear dynamic systems, including nonlinear ARX and Hammerstein-Wiener
- The System Identification Toolbox™ supports parameter estimation for linear and nonlinear parametric dynamic models.

### 3. Data-Based Simulink Models

- Engineers developing Simulink models can apply parameter estimation to develop an accurate plant model for control system design, or to create a digital twin.

- To create an accurate plant model for control design using Simulink :
  - Collect input-output test data from the plant
  - Use optimization to estimate the model's parameter values, so the simulated model output matches the measured plant output
  - You can use Simulink Design Optimization™ to interactively preprocess test data, automatically estimate model parameters, and validate estimation results.
- To create a digital twin of a current hardware asset :
  - Model the machine's dynamics in Simulink or Simscape™
  - Automatically update parameters of a model when new data from the physical asset becomes available
  - You can develop algorithms to estimate digital twin parameters with Simulink Design Optimization (8 . 37). You can also deploy the algorithm on premises, cloud, or edge systems using Simulink Compiler™.

### 1.7.2(A) Methods of Parameter Estimation

The techniques used for parameter estimation are called estimators. Some estimators are :

- Probability Plotting** : A method of finding parameter values where the data is plotted on special plotting paper and parameters are derived from the visual plot
- Rank Regression (Least Squares)** : A method of finding parameter values that minimizes the sum of the squares of the residuals.
- Maximum Likelihood Estimation** : A method of finding parameter values that, given a set of observations, will maximize the likelihood function.
- Bayesian Estimation Methods** : A family of estimation methods that tries to minimize the posterior expectation of what is called the utility function. In practice, what this means is that existing knowledge about a situation is formulated, data is gathered, and then posterior knowledge is used to update our beliefs.

#### Qualities of Estimators

If the value an estimator estimates for the parameter,  $\theta'$ , always converges to the actual parameter value  $\theta$  as the quantity of data used for parameter estimation increases, we say an estimator is **consistent**.

The bias of an estimator is the deviation of the expectation from the actual true value. If, for a given estimator, the bias is zero, we say that that estimator is **unbiased**.

### 1.7.3 Fitting a Distribution

Fitting a distribution is the process of selecting an appropriate probability distribution that best describes a given dataset. As mentioned earlier, this process involves choosing candidate distributions and estimating their parameters based on the observed data.

Common steps involved in fitting a distribution include :

- Data Selection and Preprocessing** : Gather relevant data and prepare it for analysis by handling missing values, outliers, and transforming the data if necessary.
- Candidate Distributions** : Select a set of candidate distributions based on the characteristics of the data and the underlying phenomenon.

- Data Modeling and Visualization**
- c. **Parameter Estimation** : Estimate the parameters of each candidate distribution using methods like the method of moments or maximum likelihood estimation.
  - d. **Goodness of Fit Testing** : Evaluate how well each candidate distribution fits the data using goodness-of-fit tests or graphical methods.
  - e. **Model Selection** : Choose the best-fitting distribution among the candidates based on the goodness-of-fit results.
  - f. **Interpretation and Inference** : Use the selected distribution and estimated parameters for predictions, simulations, or gaining insights into the phenomenon.

Fitting a distribution is an iterative process that may require trying different distributions and refining the parameter estimates until a satisfactory fit is achieved. It is an essential step in many statistical analyses and modeling tasks to understand the underlying data generating process and make accurate predictions or decisions based on the data.

### 1.7.3(A) Distribution Fitting

- A statistical distribution is the theoretical frequency of the occurrence of values that a variable can take. In the Simulation Fitting node, a set of theoretical statistical distributions is compared to each field of data. The distributions that are available for fitting are described in the topic Distributions.
- The parameters of the theoretical distribution are adjusted to give the best fit to the data according to a measurement of the goodness of fit; either the Anderson-Darling criterion or the Kolmogorov-Smirnov criterion. The results of the distribution fitting by the Simulation Fitting node show which distributions were fitted, the best estimates of the parameters for each distribution, and how well each distribution fits the data.
- During distribution fitting, correlations between fields with numeric storage types, and contingencies between fields with a categorical distribution, are also calculated. The results of the distribution fitting are used to create a Simulation Generate node.
- Before any distributions are fitted to your data, the first 1000 records are examined for missing values. If there are too many missing values, distribution fitting is not possible. If so, you must decide whether either of the following options are appropriate :
  - Use an upstream node to remove records with missing values.
  - Use an upstream node to impute values for missing values.
- Distribution fitting does not exclude user-missing values. If your data have user-missing values and you want those values to be excluded from distribution fitting then you should set those values to system missing.
- The role of a field is not taken into account when the distributions are fitted. For example, fields with the role Target are treated the same as fields with roles of Input, None, Both, Partition, Split, Frequency, and ID.

## 1.8 Descriptive Statistics

### Definition of Descriptive Statistics

- Descriptive statistics refers to a branch of statistics that involves summarizing, organizing, and presenting data meaningfully and concisely. It focuses on describing and analyzing a dataset's main features and characteristics without making any generalizations or inferences to a larger population.

- Data Modeling and Visualization**
- The primary goal of descriptive statistics is to provide a clear and concise summary of the data, enabling researchers or analysts to gain insights and understand patterns, trends, and distributions within the dataset. This summary typically includes measures such as central tendency (e.g., mean, median, mode), dispersion (e.g., range, variance, standard deviation), and shape of the distribution (e.g., skewness, kurtosis).
  - Descriptive statistics also involves a graphical representation of data through charts, graphs, and tables, which can further aid in visualizing and interpreting the information. Common graphical techniques include histograms, bar charts, pie charts, scatter plots, and box plots.
  - By employing descriptive statistics, researchers can effectively summarize and communicate the key characteristics of a dataset, facilitating a better understanding of the data and providing a foundation for further statistical analysis or decision-making processes.

### Descriptive Statistics Examples

**Ex. 1.8.1 :** Exam Scores Suppose you have the following scores of 20 students on an exam :

85, 90, 75, 92, 88, 79, 83, 95, 87, 91, 78, 86, 89, 94, 82, 80, 84, 93, 88, 81

To calculate descriptive statistics.

Soln. :

- **Mean** : Add up all the scores and divide by the number of scores.  

$$\text{Mean} = (85 + 90 + 75 + 92 + 88 + 79 + 83 + 95 + 87 + 91 + 78 + 86 + 89 + 94 + 82 + 80 + 84 + 93 + 88 + 81) / 20 = 1770 / 20 = 88.5$$
- **Median** : Arrange the scores in ascending order and find the middle value.  

$$\text{Median} = 86 \text{ (middle value)}$$
- **Mode** : Identify the score(s) that appear(s) most frequently.  

$$\text{Mode} = 88$$
- **Range** : Calculate the difference between the highest and lowest scores.  

$$\text{Range} = 95 - 75 = 20$$
- **Variance** : Calculate the average of the squared differences from the mean.  

$$\text{Variance} = [(85 - 88.5)^2 + (90 - 88.5)^2 + \dots + (81 - 88.5)^2] / 20 = 33.25$$
- **Standard Deviation** : Take the square root of the variance. Standard Deviation =  $\sqrt{33.25} = 5.77$

### 1.8.1 Types of Descriptive Statistics

- Descriptive statistics break down into several types, characteristics, or measures. Some authors say that there are two types. Others say three or even four.
- Datasets consist of a distribution of scores or values. Statisticians use graphs and tables to summarize the frequency of every possible value of a variable, rendered in percentages or numbers. For instance, if you held a poll to determine people's favorite Beatle, you'd set up one column with all possible variables (John, Paul, George, and Ringo), and another with the number of votes.
- Statisticians depict frequency distributions as either a graph or as a table.

## 1.8.2 Measures of Central Tendency

Measures of central tendency estimate a dataset's average or center, finding the result using three methods: mean, mode, and median.

- Mean :** The mean is also known as "M" and is the most common method for finding averages. You get the mean by adding all the response values together, and dividing the sum by the number of responses, or " $N$ ". For instance, say someone is trying to figure out how many hours a day they sleep in a week. So, the data set would be the hour entries (e.g., 6, 8, 7, 10, 8, 4, 9), and the sum of those values is 52. There are seven responses, so  $N = 7$ . You divide the value sum of 52 by  $N$ , or 7, to find  $M$ , which in this instance is 7.3.
- Mode :** The mode is just the most frequent response value. Datasets may have any number of modes including "zero." You can find the mode by arranging your dataset's order from the lowest to highest value and then looking for the most common response. So, in using our sleep study from the last part: 4, 6, 7, 8, 8, 9, 10. As you can see, the mode is eight.
- Median :** Finally, we have the median, defined as the value in the precise center of the dataset. Arrange the values in ascending order (like we did for the mode) and look for the number in the set's middle. In this case, the median is eight.

### Variability (Also Called Dispersion)

The measure of variability gives the statistician an idea of how spread out the responses are. The spread has three aspects - range, standard deviation, and variance.

- Range :** Use range to determine how far apart the most extreme values are. Start by subtracting the dataset's lowest value from its highest value. Once again, we turn to our sleep study : 4, 6, 7, 8, 8, 9, 10. We subtract four (the lowest) from ten (the highest) and get six. There's your range.
- Standard Deviation :** This aspect takes a little more work. The standard deviation ( $s$ ) is your dataset's average amount of variability, showing you how far each score lies from the mean. The larger your standard deviation, the greater your dataset's variable. Follow these six steps :

1. List the scores and their means.
2. Find the deviation by subtracting the mean from each score.
3. Square each deviation.
4. Total up all the squared deviations.
5. Divide the sum of the squared deviations by  $N - 1$ .
6. Find the result's square root.

Raw Number/Data	Deviation from Mean	Deviation Squared
4	$4 - 7.3 = -3.3$	10.89
6	$6 - 7.3 = -1.3$	1.69
7	$7 - 7.3 = -0.3$	0.09

Raw Number/Data	Deviation from Mean	Deviation Squared
8	$8 - 7.3 = 0.7$	0.49
8	$8 - 7.3 = 0.7$	0.49
9	$9 - 7.3 = 1.7$	2.89
10	$10 - 7.3 = 2.7$	7.29
M = 7.3	Sum = 0.9	Square sums = 23.83

When you divide the sum of the squared deviations by  $N - 1$  : 23.83/6, you get 3.971, and the square root of that result is 1.992. As a result, we now know that each score deviates from the mean by an average of 1.992 points.

**Variance :** Variance reflects the dataset's degree spread. The greater the degree of data spread, the larger the variance relative to the mean. You can get the variance by just squaring the standard deviation. Using the above example, we square 1.992 and arrive at 3.971.

### 1.8.3 Univariate Descriptive Statistics

Univariate descriptive statistics examine only one variable at a time and do not compare variables. Rather, it allows the researcher to describe individual variables. As a result, this sort of statistic is also known as descriptive statistics. The patterns identified in this sort of data may be explained using the following :

- Measures of central tendency (mean, mode, and median)
- Data dispersion (standard deviation, variance, range, minimum, maximum, and quartiles) (standard deviation, variance, range, minimum, maximum, and quartiles)
- Tables of frequency distribution
- Pie graphs
- Frequency polygon histograms
- Bar graphs

### 1.8.4 Bivariate Descriptive Statistics

- When using bivariate descriptive statistics, two variables are concurrently analyzed (compared) to see whether they are correlated. Generally, by convention, the independent variable is represented by the columns, and the rows represent the dependent variable.'
- There are numerous real-world applications for bivariate data. For example, estimating when a natural occurrence will occur is quite valuable. Bivariate data analysis is a tool in the statistician's toolbox. Sometimes, something as simple as projecting one parameter against the other on a Two-dimensional plane can better understand what the information is trying to convince you. For example, the scatterplot below demonstrates the link between the period between eruptions at Old Faithful and the eruption's duration.

### 1.8.5 Main Purpose of Descriptive Statistics

Descriptive statistics can be useful for two things :-

- 1) providing basic information about variables in a dataset and
- 2) highlighting potential relationships between variables. Graphical/Pictorial Methods are measures of the three most common descriptive statistics that can be displayed graphically or pictorially. It is used to summarise data. Descriptive statistics only make statements about the data set used to calculate them; they never go beyond your data.

#### Scatter Plots

- A scatter plot employs dots to indicate values for two separate numeric variables. Each dot's location on the horizontal and vertical axes represents a data point's values. Scatter plots are being used to monitor relationships between variables.
- The main purposes of scatter plots are to examine and display relationships between two numerical variables. The points in a scatter plot document the values of individual points and trends when the data is obtained as a whole. Identification of correlational links is prevalent with scatter plots. In these situations, we want to know what a good vertical value prediction would be given a specific horizontal value.
- This can lead to overplotting when there are many data points to plot. When data points are overlaid to the point where it is difficult to see the connections between them and the variables, this is known as overplotting. It might be difficult to discern how densely-packed data points are when lots of them are in a tiny space.

### 1.8.6 Difference between Descriptive Statistics and Inferential Statistics

- So, what's the difference between the two statistical forms? We've already touched upon this when we mentioned that descriptive statistics doesn't infer any conclusions or predictions, which implies that inferential statistics do so.
- Inferential statistics takes a random sample of data from a portion of the population and describes and makes inferences about the entire population. For instance, in asking 50 people if they liked the movie they had just seen, inferential statistics would build on that and assume that those results would hold for the rest of the movie-going population in general.
- Therefore, if you stood outside that movie theater and surveyed 50 people who had just seen *Rocky 20: Enough Already!* and 38 of them disliked it (about 76 percent), you could extrapolate that 76% of the rest of the movie-watching world will dislike it too, even though you haven't the means, time, and opportunity to ask all those people.

## 1.9 Graphical Statistics

- Graphical statistics Despite highly developed theory and methodology of Statistics, when it comes to analysis of real data, experienced statisticians will often follow a very simple advice : Before you do anything with a data set, look at it!

- A quick look at a sample may clearly suggest :
  - a probability model, i.e., a family of distributions to be used;
  - statistical methods suitable for the given data; 230 Probability and Statistics for Computer Scientists
  - presence or absence of outliers;
  - presence or absence of heterogeneity;
  - existence of time trends and other patterns;
  - relation between two or several variables.
- There is a number of simple and advanced ways to visualize data. This section introduces
  1. histograms,
  2. stem-and-leaf plots,
  3. boxplots,
  4. time plots,
  5. scatter plots.
- Each graphical method serves a certain purpose and discovers certain information about data.

#### 1. Histogram

- A histogram shows the shape of a pmf or a pdf of data, checks for homogeneity, and suggests possible outliers. To construct a histogram, we split the range of data into equal intervals, "bins," and count how many observations fall into each bin.
- A frequency histogram consists of columns, one for each bin, whose height is determined by the number of observations in the bin. A relative frequency histogram has the same shape but a different vertical scale. Its column heights represent the proportion of all data that appeared in each bin.

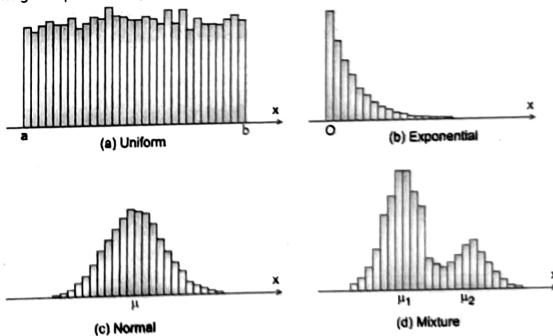


Fig. 1.9.1 : Histograms of various samples

**2. Boxplot**

The main descriptive statistics of a sample can be represented graphically by a boxplot. To construct a boxplot, we draw a box between the first and the third quartiles, a line inside for a median, and extend whiskers to the smallest and the largest observations.

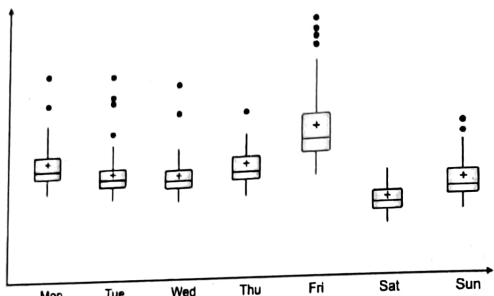


Fig. 1.9.2 : Parallel boxplots of internet traffic

**3. Scatter plots and time plots**

Scatter plots are used to see and understand a relationship between two variables. These can be temperature and humidity, experience and salary, age of a network and its speed, number of servers and the expected response time, and so on. To study the relationship, both variables are measured on each sampled item. For example, temperature and humidity during each of n days, age and speed of n networks, or experience and salary of n randomly chosen computer scientists are recorded. Then, a scatter plot consists of n points on an (x, y)-plane, with x- and y-coordinates representing the two recorded variables.

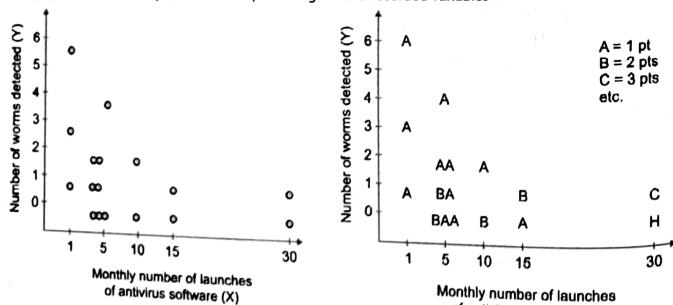


Fig. 1.9.3 : Scatter plots

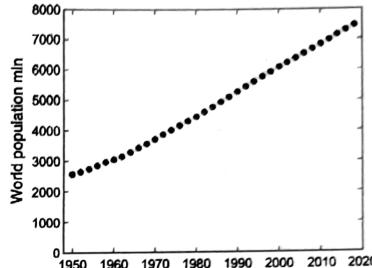


Fig. 1.9.4 : Time Plot

**4. Stem-and-leaf plot**

Stem-and-leaf plots are similar to histograms although they carry more information. Namely, they also show how the data are distributed within columns.

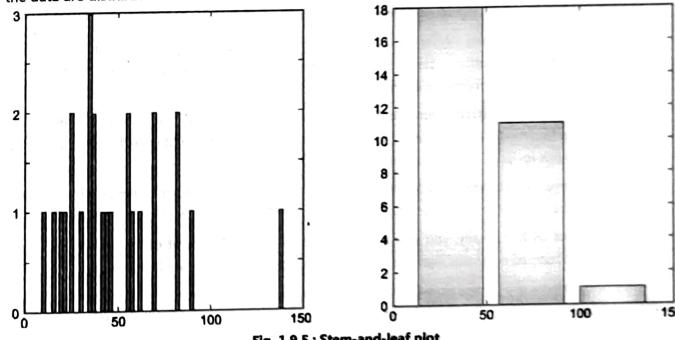


Fig. 1.9.5 : Stem-and-leaf plot

To construct a stem-and-leaf plot, we need to draw a stem and a leaf. The first one or several digits form a stem, and the next digit forms a leaf. Other digits are dropped; in other words, the numbers get rounded. For example, a number 239 can be written as

23 | 9

with 23 going to the stem and 9 to the leaf, or as

2 | 3

with 2 joining the stem, 3 joining the leaf, and digit 9 being dropped. In the first case, the leaf unit equals 1 while in the second case, the leaf unit is 10, showing that the (rounded) number is not 23 but 230.

## 1.10 Method of Moments

### Definition of Method of Moments

- The method of moments is a way to estimate population parameters, like the population mean or the population standard deviation. The basic idea is that you take known facts about the population, and extend those ideas to a sample.
- The method of moments is a statistical technique used for estimating the parameters of a probability distribution based on a sample of data. The basic idea behind the method of moments is to equate the sample moments (such as mean, variance, etc.) with their corresponding population moments, which are expressed in terms of the unknown parameters of the distribution. By solving these equations, we can obtain estimates for the unknown parameters.
- Here's a general outline of how the method of moments works :
  - Choose a probability distribution :** First, you need to make an assumption about the underlying probability distribution that generated the data. This assumption is typically based on prior knowledge or the characteristics of the data.
  - Calculate sample moments :** Next, you calculate the sample moments from your data. For example, the first moment is the sample mean, the second moment is the sample variance, and so on.
  - Express population moments :** After choosing a distribution, you express the population moments (mean, variance, etc.) of that distribution in terms of its parameters. For example, for a normal distribution, the mean and variance are population moments expressed in terms of the mean ( $\mu$ ) and variance ( $\sigma^2$ ) parameters.
  - Equate sample and population moments :** Set up equations by equating the sample moments with their corresponding population moments, using the estimates for the moments obtained from the data in step 2 and the expressions for population moments from step 3.
  - Solve for parameter estimates :** Solve the equations from step 4 to obtain estimates for the unknown parameters of the distribution. These estimates are known as the method of moments estimators.
  - Assess goodness of fit :** Finally, you may want to assess the goodness of fit of the chosen distribution by comparing the observed moments with the moments predicted by the estimated parameters.
- The method of moments is a straightforward and intuitive approach to parameter estimation, but its accuracy and efficiency can vary depending on the distributional assumption and the sample size. In some cases, more advanced estimation methods like Maximum Likelihood Estimation (MLE) or Bayesian estimation might be preferred. However, the method of moments remains a valuable tool in the statistician's toolkit, especially for quick and simple estimation tasks.
- For example, it's a fact that within a population :

$$\text{Expected value } E(x) = \mu$$

For a sample, the estimator  $\hat{\mu}$  is just the sample mean is  $\bar{X}$ :

The formula for the sample mean is :

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

This is the first moment condition.

- The second moment condition involves the variance. The population variance is  $\text{Var}(x) = \sigma^2$ , so we just need to use the method of moments to estimate the variance in the sample. Here's how the formula is derived :
  - Use the fact that the population variance  $\text{Var}(x) = \sigma^2$  is the same as :  $E(x - \mu)^2 = \sigma^2$ .
  - As in the first moment, replace the population expectation by the sample equivalent (the sample mean). However, while the first moment was the sample mean of  $x$ , this time it's  $x - \mu$ , giving the formula.
- The same principle is used to derive higher moments like skewness and kurtosis :

### Skewness

$$\alpha_3 = \frac{1}{100} \sigma^{-3} \sum f(X_i - \bar{X})^3$$

### Kurtosis

$$\alpha_4 = \frac{1}{100} \sigma^{-4} \sum f(X_i - \bar{X})^4$$

- The above method is probably the most widely used method of moments. There is another method, which uses sample moments about the mean instead of sample moments about the origin.

### Advantages and Disadvantages

#### Advantages

- Method of moments is simple (compared to other methods like the maximum likelihood method) and can be performed by hand.

#### Disadvantages

- The parameter estimates may be inaccurate. This is more frequent with smaller samples and less common with large samples.
- The method may not result in sufficient statistics. In other words, it may not take into account all of the relevant information in the sample.

## 1.11 Maximum Likelihood Estimation

- Maximum likelihood estimation is a method that determines values for the parameters of a model. The parameter values are found such that they maximise the likelihood that the process described by the model produced the data that were actually observed.
- The above definition may still sound a little cryptic so let's go through an example to help understand this.
- Let's suppose we have observed 10 data points from some process. For example, each data point could represent the length of time in seconds that it takes a student to answer a specific exam question. These 10 data points are shown in the Fig. 1.11.1.

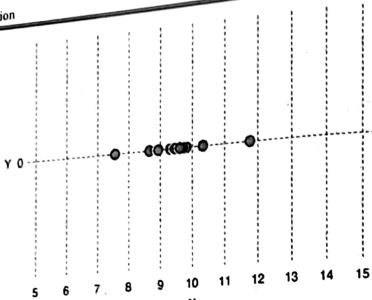


Fig. 1.11.1 : Maximum Likelihood Estimation

The 10 (hypothetical) data points that we have observed :

- We first have to decide which model we think best describes the process of generating the data. This part is very important. At the very least, we should have a good idea about which model to use. This usually comes from having some domain expertise but we won't discuss this here.
- For these data we'll assume that the data generation process can be adequately described by a Gaussian (normal) distribution. Visual inspection of the Fig. 1.11.1 suggests that a Gaussian distribution is plausible because most of the 10 points are clustered in the middle with few points scattered to the left and the right (Making this sort of decision on the fly with only 10 data points is ill-advised but given that I generated these data points we'll go with it).

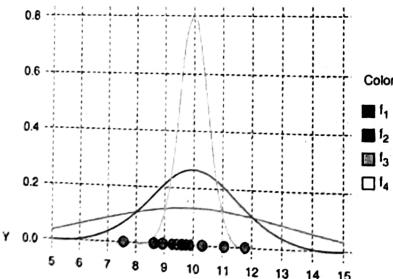


Fig. 1.11.2 : Gaussian distribution

- The 10 data points and possible Gaussian distributions from which the data were drawn.  $f_1$  is normally distributed with mean 10 and variance 2.25 (variance is equal to the square of the standard deviation), this is likelihood is to find the parameter values that give the distribution that maximise the probability of observing the data.
- The true distribution from which the data were generated was  $f_1 \sim N(10, 2.25)$ , which is the blue curve in the figure.

### Calculating the Maximum Likelihood Estimates

- Now that we have an intuitive understanding of what maximum likelihood estimation is we can move on to learning how to calculate the parameter values. The values that we find are called the maximum likelihood estimates (MLE).
- Again we'll demonstrate this with an example. Suppose we have three data points this time and we assume that they have been generated from a process that is adequately described by a Gaussian distribution. These points are 9, 9.5 and 11. How do we calculate the maximum likelihood estimates of the parameter values of the Gaussian distribution  $\mu$  and  $\sigma$ ?
- What we want to calculate is the total probability of observing all of the data, i.e. the joint probability distribution of all observed data points. To do this we would need to calculate some conditional probabilities, which can get very difficult. So it is here that we'll make our first assumption. The assumption is that each data point is generated independently of the others. This assumption makes the maths much easier. If the events (i.e. the process that generates the data) are independent, then the total probability of observing all of the data is the product of observing each data point individually (i.e. the product of the marginal probabilities).
- The probability density of observing a single data point  $x$ , that is generated from a Gaussian distribution is given by :

$$P(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

- The semi colon used in the notation  $P(x; \mu, \sigma)$  is there to emphasise that the symbols that appear after it are parameters of the probability distribution. So it shouldn't be confused with a conditional probability (which is typically represented with a vertical line e.g.  $P(A|B)$ ).

In our example the total (joint) probability density of observing the three data points is given by :

$$P(9, 9.5, 11; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9-\mu)^2}{2\sigma^2}\right) \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9.5-\mu)^2}{2\sigma^2}\right) \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(11-\mu)^2}{2\sigma^2}\right)$$

We just have to figure out the values of  $\mu$  and  $\sigma$  that results in giving the maximum value of the above expression.

- If you've covered calculus in your maths classes then you'll probably be aware that there is a technique that can help us find maxima (and minima) of functions. It's called differentiation.
- All we have to do is find the derivative of the function, set the derivative function to zero and then rearrange the equation to make the parameter of interest the subject of the equation. And viola, we'll have our MLE values for our parameters. I'll go through these steps now but I'll assume that the reader knows how to perform differentiation on common functions. If you would like a more detailed explanation then just let me know in the comments.

## 1.12 Data Modeling Concepts

### What is a Data Model?

- A Data Model is a graphical technique to construct your database that allows you to start with the big ideas. It provides a uniform framework for representing real-world entities and their attributes. These entities have been reduced to their most basic elements and are frequently linked in some way. Everything is expressed in a fairly straightforward manner, including entities, their properties, and relationships between things. Data models also assist you in describing how entity data is organized, stored, and queried.

**Data Modeling and Visualization**

- The Data Model gives us an idea of how the final system will look after it has been fully implemented. It specifies the data items as well as the relationships between them. In a database management system, data models are used to show how data is stored, connected, accessed, and changed. We portray the information using a set of symbols and language so that members of the organization may communicate and understand it. For more understanding on the different types of data models and all that they entail, take a look at this blog from our archives.

**Why Create a Data Model?**

Creating and maintaining a Data model for your database has numerous advantages. A couple of these having a Data model from a business standpoint:

- A Data model creates a common communication layer that makes it easier for the Data Architect and the business folks to communicate. It ensures that all the intricacies can be discussed in detail because it is a visual model with a shared terminology.
- High-quality documentation ensures that the code that is implemented is also of high quality. The code builds on the preceding documentation phase's decisions, reducing the number of errors.
- Developers can spend more time on feature development because there is less code that needs to be changed. As a result, the amount of time spent coding is reduced, and some costs are eliminated.
- Creating a Data Model allows you to determine the scope of the project. Complexity is decreased and tough topics are comprehended because of the Data model's visual depiction.
- The following are the technical advantages of having a Data Model :
  - Technical Layer :** A technical layer is attached to a Data model which contains all the technical information (specified by the Data Architect), allowing developers to concentrate on implementation rather than interpretation.
  - Lesser Mistakes :** There are fewer mistakes. Fewer mistakes are made on the data and application side as a result of the Data model's clarity and precision. Developers can concentrate on feature development rather than database design.
  - Database Structure Optimization :** The Database structure can be optimized right from the start before any data is entered. This decreases the amount of data that needs to be moved (i.e. to improve performance after the database is in production).
  - Data Risk Reduction :** Risks to data are reduced. Data Architects and Database Administrators can develop backup and restore procedures if they have a better understanding of the data's size. In a disaster recovery scenario, having strategies and safeguards in place decreases risks.
- Data Modeling is the process of cleansing and organizing data into a visual representation or plan that aids in the mapping out of database relationships and operations. Various types of data models serve as a blueprint for creating an optimal database, regardless of its specific contents.
- A Data Modeler completes the process by directly working with data entities and characteristics to determine their relationships and construct a suitable model. Data architects also work on various types of data models, focusing on the creation of physical blueprints.

**1.12.1 Types of Data Models**

- This article talks about the different types of Data Models and the process of Data Modeling along with a few Data Modeling Tools. We're also going to take a look at the types of data modeling later in the article to provide you with a comprehensive understanding of the topic.

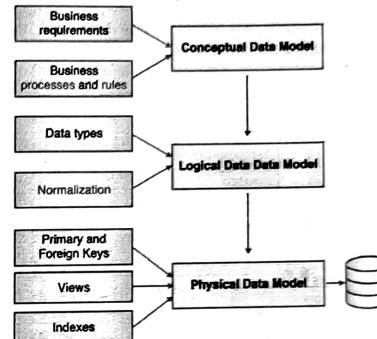


Fig. 1.12.1 : Types of Data Models

Database and information system design, like any other design process, starts with a high level of abstraction and gradually gets more concrete and specific. Based on the level of abstraction they give, there are three types of Data models. The approach will begin with a conceptual model and progress to a logical model before arriving at a physical model. The parts that follow go through each of the types of Data models in greater detail.

- 1. Conceptual Data Model :** Conceptual Data Model is the first in types of Data Model. They're also known as **'Domain Models'**, and they provide a big-picture view of the system's contents, how it'll be organized, and which business rules will be involved. Typically, Conceptual Models are generated as part of the early project requirements gathering process. Entity classes (which define the types of items that are significant for the company to represent in the Data model), their characteristics and limitations, their relationships, and applicable security and Data Integrity requirements are typically included. Any notation is usually straightforward.

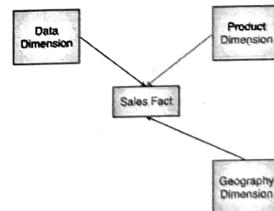


Fig. 1.12.2 : Example of Conceptual Data Model

Data Modeling and Visualization

- 2. Logical Data Model :** The Logical Data Model is the second one in the types of Data models. They are less abstract and provide more information on the concepts and relationships in the domain. It follows one of the formal Data Modelling notation systems. Logical data models show the relationships between entities as well as Data properties such as Data types and their respective lengths. Technical system needs are not specified in Logical Data Models.
- In agile or DevOps techniques, the logical data Modeling stage is frequently skipped. In highly procedural implementation contexts or for projects that are data-oriented by nature, such as Data warehouse design or reporting system development, Logical Data Models might be valuable.

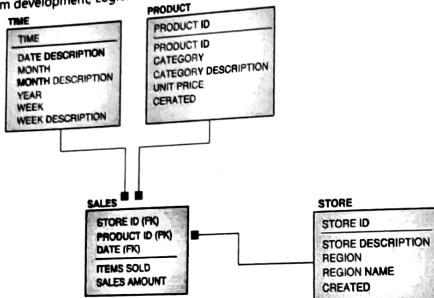


Fig. 1.12.3 : Logical Data Model

- 3. Physical Data Model :** Physical Data Model is the last one in the types of Data models. They define the format in which data will be physically stored in a database. As a result, they are the least abstract of the bunch. Physical Data Models provide a finished design that can be implemented as a Relational Database, complete with associative tables that show the associations between entities. The primary keys and foreign keys will be utilized to keep those relationships up to date. Physical Data Models might incorporate attributes related to a Database Management System (DBMS), such as performance tweaking.

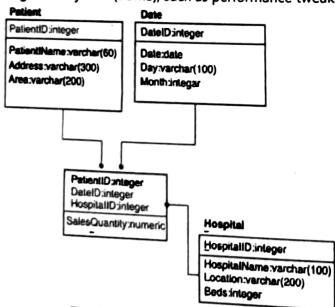


Fig. 1.12.4 : Physical Data Model

**1.13 Data Modeling Process**

Data Modeling as a discipline invites stakeholders to examine Data processing and storage in great depth. Different conventions govern which symbols are used to represent data, how models are built up, and how business needs are communicated in Data Modeling methodologies. All approaches provide structured workflows that comprise a list of actions that must be completed in sequential order. An example of a workflow is shown as follows :

- Determine the entities. The identification of the items, events, or concepts represented in the Data set to be modeled is the first step in the Data modeling process. Each entity should be logically different from the others while maintaining a sense of unity.
- Determine the most important properties of each entity. One or more distinct traits called attributes, separate each entity type from the others. For example, a "customer" object might have a first name, last name, phone number, and salutation, but an "address" entity might have a street name and number, as well as a city, state, nation, and zip code.
- Determine the relationships that exist between entities. In the first draught of a Data model, the nature of each entity's relationships with the others will be established. In the above example, each client "resides" at a certain address. If the model was expanded to include an object named "orders," each order would be shipped to and paid to an address. These links are often defined using a single modeling language (UML).
- Complete the attribute-to-entity mapping. This ensures that the model appropriately reflects the company's data usage intentions. There are a variety of formal Data modeling paradigms in use. Object-oriented developers commonly employ analysis and design patterns, although stakeholders from other business areas may utilize different patterns.
- Assign keys as appropriate, and choose a level of normalization that strikes a balance between the need to decrease redundancy and the necessity for performance. Normalization is a method of structuring the various types of Data models (and the databases they represent) in which numerical identifiers, known as keys, are allocated to groupings of data to indicate relationships between them without having to repeat the data.
- For example, if each client is given a key, that key can be linked to their address and order history without having to duplicate the information in the customer database. Normalization reduces the amount of storage space required by a database, but at the expense of query performance.
- Complete and test the different types of Data models. As business demands evolve, Data modeling is an iterative process that should be repeated and updated.

**1.13.1 Types of Data Modeling**

As enterprises' data storage needs have expanded, Data Modelling has evolved alongside Database Management Systems, with model types becoming more complicated. Listed below are a few different types of models :

- Hierarchical Data Models :** In a treelike style, Hierarchical Data Models represent one-to-many relationships. Each record in this architecture has a single root or parent table that translates to one or more child tables. This paradigm was adopted in the IBM Information Management System (IMS), which was

released in 1966 and quickly became popular, particularly in the banking industry. Although this method is less efficient than more recently established Database models, it is nevertheless utilized in XML systems and Geographic Information Systems (GISs).

- Relational Data Models :** E.F. Codd, an IBM researcher, first introduced Relational Data Models in 1970. They're still employed in today's relational databases, which are widely used in enterprise computing. Relational Data modeling does not necessitate a thorough knowledge of the physical aspects of the Data storage system. It reduces database complexity by explicitly joining Data segments through the use of tables.
- Entity-Relationship (ER) Data Models :** ER Data Models represent the relationships between entities in a database using formal diagrams. Data architects utilize a variety of ER modeling tools to produce visual maps that communicate database architecture goals.
- Object-oriented Data Models :** These various types of Data models became popular at the same time as object-oriented programming in the mid-1990s. The "things" in question are fictitious representations of real-world entities. Objects are organized into hierarchies of classes, each with its own set of attributes. Tables can be included in object-oriented databases, but they can also handle more sophisticated data interactions. This method is used in multimedia and hypertext databases, among other applications.
- Dimensional Data Models :** Ralph Kimball created Dimensional Data Models, which were created to speed up data retrieval for analytic purposes in a Data Warehouse. While relational and ER models focus on efficient storage, Dimensional models include redundancy to make it easier to find information for reporting and retrieval. This type of modeling is common in OLAP systems.

The Star Schema, in which data is structured into facts (measurable elements) and dimensions (reference information), and each fact is encircled by its corresponding dimensions in a star-like pattern, is a prominent Dimensional Data model. The other is the Snowflake Schema, which is similar to the star schema but incorporates additional layers of linked dimensions, increasing the complexity of the branching pattern.

## 1.14 Understand and Model Subtypes and Supertypes

In data modeling, the concepts of subtypes and supertypes are used to represent relationships between entities and to model hierarchical structures in the data. These concepts are often used in entity-relationship diagrams (ER diagrams) to visually depict the relationships between different types of entities.

### Supertypes

- A supertype is a generalization of one or more related entity types. It represents a higher-level entity that encompasses common attributes and relationships shared by its subtypes. Supertypes do not have instances on their own but serve as a template for their subtypes. In other words, a supertype defines a common set of properties and relationships that are inherited by its subtypes.
- For example, consider an entity called "Vehicle" as a supertype. It could have attributes like "Manufacturer," "Model," and "Year," as well as relationships like "Owned By" and "Registered In." The "Vehicle" supertype would be more abstract and generalized, representing common features shared by different types of vehicles, such as cars, motorcycles, and trucks.

### Subtypes

- A subtype is a specialized entity type that inherits attributes and relationships from its supertype. It represents a more specific category or subclass of the supertype, focusing on distinct characteristics that are unique to that particular subtype.
- Continuing with the "Vehicle" example, the subtypes could be "Car," "Motorcycle," and "Truck." Each of these subtypes would have its own unique attributes and relationships, such as "Number of Doors" for cars, "Number of Wheels" for motorcycles, and "Cargo Capacity" for trucks.

### Key points to understand about subtypes and supertypes

- Subtypes are mutually exclusive, meaning an instance of a particular subtype cannot belong to more than one subtype simultaneously. For instance, a vehicle cannot be both a car and a motorcycle at the same time.
- Every instance of a subtype is also an instance of its supertype. All attributes and relationships defined in the supertype are inherited by the subtypes.
- Subtypes can have additional attributes and relationships specific to their unique characteristics.
- Subtypes can participate in relationships independently. For example, a "Car" subtype may have a relationship with a "Driver" entity, while a "Truck" subtype may have a relationship with a "Shipping Company" entity.
- Subtypes can also have their own subtypes, forming a hierarchical structure.
- The use of subtypes and supertypes in data modeling helps in organizing data in a more structured and efficient manner. It allows for better representation of complex relationships and variations within the data, leading to more accurate and flexible database designs.

### Using Subtypes and Supertypes

- It is not surprising that many of the arguments that arise in data modeling are about the appropriate level of generalization, although they are not always recognized as such. We cannot easily resolve such disputes by turning to the rulebook, nor do we want to throw away interesting options too early in the modeling process.
- The ability to represent different levels of generalization requires a new diagramming convention, the box-in-box. You should be very wary about overcomplicating diagrams with too many different symbols, but this one literally adds another dimension (generalization/specialization) to our models.

### Subtypes and Supertypes as Entity Classes

Much of the confusion that surrounds the proper use of subtypes and supertypes can be cleared with a simple rule : subtypes and supertypes are entity classes.

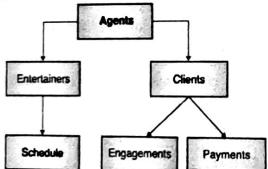
- We use the same diagramming convention (the box with rounded corners) to represent all entity classes, whether or not they are subtypes or supertypes of some other entity class(es).
- Subtypes and supertypes must be supported by definitions.
- Subtypes and supertypes can have attributes. Attributes particular to individual subtypes are allocated to those subtypes; common attributes are allocated to the supertype.

**Data Modeling and Visualization**

- Subtypes and supertypes can participate in relationships. Notice in our family tree model how neatly we have been able to capture our "mother of" and "father of" relationships by tying them to entity classes at the most appropriate level. In fact, this diagram shows most of the sorts of relationships that seem to worry modelers, in particular the relationship between an entity class and its own supertype.
- Subtypes can themselves have subtypes. We need not restrict ourselves to two levels of subtyping. In practice, we tend to represent most concepts at one, two, or three levels of generality, although four or five levels are useful from time to time.

**1.15 Understand and Model Hierarchical Data**

- A hierarchical model represents the data in a tree-like structure in which there is a single parent for each record. To maintain order there is a sort field which keeps sibling nodes into a recorded manner. These types of models are designed basically for the early mainframe database management systems, like the Information Management System (IMS) by IBM.
- This model structure allows the one-to-one and a one-to-many relationship between two/ various types of data. This structure is very helpful in describing many relationships in the real world; table of contents, any nested and sorted information.
- The hierarchical structure is used as the physical order of records in storage. One can access the records by navigating down through the data structure using pointers which are combined with sequential accessing. Therefore, the hierarchical structure is not suitable for certain database operations when a full path is not also included for each record.
- Data in this type of database is structured hierarchically and is typically developed as an inverted tree. The "root" in the structure is a single table in the database and other tables act as the branches flowing from the root. The Fig. 1.15.1 shows a typical hierarchical database structure.

**Fig. 1.15.1 : Typical Hierarchical Database Structure****Agents Database**

- In the above diagram, an agent books several entertainers, and each entertainer, in return has his/her own schedule. It is the duty of an agent to maintain several clients whose entertainment needs are to be met. A client books engagement through the agent and makes payments to the agent for his services.
- A relationship in this database model is represented by the term parent/child. A parent table can be linked with one or more child tables in this type of relationship, but a single child table can be linked with only one parent table. The tables are explicitly linked via a pointer/index or by the physical arrangement of the records within the tables.

- A user can access the data by starting at the root table and working down through the tree to the target data. the user must be familiar with the structure of the database to access the data without any complexity.

**Advantages**

- A user can retrieve data very quickly due to the presence of explicit links between the table structures.
- The referential integrity is built in and automatically enforced due to which a record in a child table must be linked to an existing record in a parent table, along with that if a record deleted in the parent table then that will cause all associated records in the child table to be deleted as well.

**Disadvantages**

- When a user needs to store a record in a child table that is currently unrelated to any record in a parent table, it gets difficulty in recording and user must record an additional entry in the parent table.
- This type of database cannot support complex relationships, and there is also a problem of redundancy, which can result in producing inaccurate information due to the inconsistent recording of data at various sites.
- Consider an example using the database diagram shown in the previous diagram. A user cannot enter a new record for the entertainer in the Entertainers table until the entertainer is assigned to a specific agent in the Agents table since a record in a child table (Entertainers) must be related to a record in the parent table (Agents).
- Therefore, this type of database suffers from the problem of redundant data. For example, if there is a many-to-many relationship between clients and entertainers; an entertainer will perform for many clients, and a client will hire many entertainers. This type of relationship in a hierarchical database cannot easily model, so developers must introduce redundant data into both the Schedule and Engagements tables.
- The Schedule table will now have client data which contains information such as client name, address, and phone number to show for whom and where each entertainer is performing. This data is redundant because it is currently stored also in the Clients table.
- The Engagements table will now contain data on entertainers which contains information such as entertainer name, phone number, and type of entertainer to indicate which entertainers are performing for a given client. This data is also redundant because it is currently stored in the Entertainers table.
- The problem with this redundancy is that it can result in producing inaccurate information because it opens the possibility of allowing a user to enter a single piece of data inconsistently.
- This problem can be solved by creating one hierarchical database specifically for entertainers and another one specifically for agents. The Entertainers database will contain only the data recorded in the Entertainers table, and the revised Agents database will contain the data recorded in Agents, Clients, Payments, and Engagements tables. there is no need of as you can define a logical child relationship between the Engagements table in the Agents database and the Entertainers table in the Entertainers database. With this relationship in place, you can retrieve a variety of information, such as a list of booked entertainers for a given client or a performance schedule for a given entertainer. The Fig. 1.15.2 describes the whole picture.

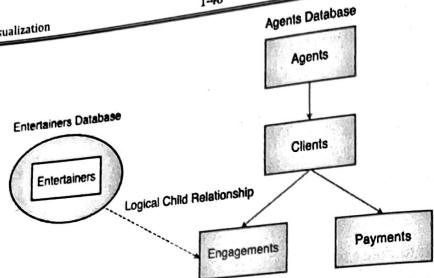


Fig. 1.15.2 : Booked entertainers for a given client

- The hierarchical database suited well to the tape storage systems which is used by mainframes in the 1970s and was very popular in organizations whose database is based on those systems. But, even though the hierarchical database provided fast and direct access to data and was useful in several circumstances, it was clear that a new database model was needed to address the growing problems of data redundancy and complex relationships among data.
- The idea behind this database model is useful for a certain type of data storage, but it is not extremely versatile and is confined to some specific uses.
- For example, where each individual person in a company may report to a given department, the department can be used as a parent record and the individual employees will represent secondary records, each of which links back to that one parent record in a hierarchical structure.

### 1.15.1 Hierarchical Database Model

- A hierarchical database model is a data model in which the data are organized into a tree-like structure. Hierarchical model is one of the oldest database models, dating from 1960's. First Hierarchical database information management system was developed jointly by North American Rockwell company and IBM.
- Hierarchy is an ordered tree and easy to understand. Data is represented in form of records. Relationship among the data is represented by records of links. A tree may be defined as a set of nodes. At the root of the tree is the single parent, the parent can have none, one or more children.

### 1.15.2 Features of Hierarchical Database Model

- One to many relationships :** It only supports one-to-many relationships. Many to many relationships are not supported.
- Problem in Deletion :** If a parent is deleted then the child automatically gets deleted.
- Hierarchy of data :** Data is represented in a hierarchical tree-like structure.
- Parent-child relationship :** Each child can have only one parent but a parent can have more than one children.
- Pointer :** Pointers are used for linking records that tell which is a parent and which child record is.

- Disk input and output is minimized :** Parent and child records are placed or stored close to each other on the storage device which minimizes the hard disk input and output.
- Fast navigation :** As parent and child are stored close to each other so access time is reduced and navigation becomes faster.
- Predefined relationship :** All relations between root, parent and child nodes are predefined in the database schema.
- Re-organization difficulty :** Hierarchy prevents the re-organization of data.
- Redundancy :** One to many relationships increases redundancy in the data which leads to the retrieval of inaccurate data.

### 1.15.3 Advantages and Disadvantages of Hierarchical Model

#### Advantages

- Promotes data sharing
- Parent/child relationship promotes conceptual simplicity and data integrity
- Database security is provided and enforced by DBMS
- Efficient with 1 : M relationships

#### Disadvantages

- Requires knowledge of physical data storage characteristics
- Navigational system requires knowledge of hierarchical path
- Changes in structure require changes in all application programs
- Implementation limitations
- No data definition
- Lack of standards

### 1.16 Understand and Model Recursive Relationships

#### Recursive Relationship in ER Diagrams

ER Diagram stands for Entity Relationship Diagram. When we draw the relationships between entities using a diagram, then it is called an entity relationship diagram. The ER diagram is just for understanding the purpose of the database administrator. We cannot use the ER diagram directly on the computer.

ER Diagrams are converted into the tabular form then they are inserted into the computer using any query language. In ER Diagrams, we use attributes, entities, and the relationships between entities. We use an oval shape to represent the entity and a diamond shape to represent the relationships between entities.

#### 1.16.1 Cardinality

- The number of the relationship between two entity sets can have, is known as Cardinality. The quantity of items in a particular set is referred to mathematically as its Cardinality. Cardinality can be used by Database Administrators to count values and tables. Cardinality in a Database typically illustrates the relationship between the data in two distinct tables by emphasizing how frequently one object occurs in comparison to another.

Data Modeling and Visualization

- Because it establishes organized linkages between different tables or entities, Cardinality is crucial in Databases. The query execution plan-a series of actions users might take to find and access data in a Database System-is significantly impacted by this. A well-organized query execution strategy can make it simpler for consumers to find the information they require fast. The Cardinality model can be used in Databases by Database Managers for a variety of purposes, but corporations often use it to evaluate customer or inventory data.
- The Mapping Cardinality, also known as the Cardinality Ratio, in a database refers to the number of entities that can be connected to another entity through a specific relation set. Although they can help with the description of relation sets with more than two entity sets, Mapping Cardinality is most beneficial when defining binary relation sets.

**Types of Cardinality in DBMS :**

- One to one cardinality
- Many to one cardinality
- One to many cardinality
- Many to many cardinality

**1.16.2 Recursive Relationship**

- If there is a relationship between two entities and they both belong to the same entity type, then this is called a recursive relationship. We can say a relationship between the same entity set is called a recursive relationship. In the ER diagram, we can use one entity that will enter and exit from itself will denote the relationship. While mapping the ER diagram into the table, we must consider the foreign key in the table for mapping the recursive relationship.
- A relationship between two entities of a similar entity type is called a recursive relationship. Here the same entity type participates more than once in a relationship type with a different role for each instance. In other words, a relationship has always been between occurrences in two different entities. However, the same entity can participate in the relationship. This is termed a recursive relationship.
- Recursive relationships are often used to represent hierarchies or networks, where an entity can be connected to other entities of the same type. For example, in an organizational chart, an employee can have a relationship with other employees who are also in a managerial position. Similarly, in a social network, a user can have a relationship with other users who are their friends.
- To represent a recursive relationship in an ER diagram, we use a self-join, which is a join between a table and itself. In other words, we create a relationship between the same entity type. The self-join involves creating two instances of the same entity and connecting them with a relationship. One instance is considered the parent, and the other instance is considered the child.
- We use cardinality constraints to specify the number of instances of the entity that can participate in the relationship. For example, in an organizational chart, an employee can have many subordinates, but each subordinate can only have one manager. This is represented as a one-to-many (1 : N) relationship between the employee entity and itself.
- Overall, recursive relationships are an important concept in ER modeling, and they allow us to represent complex relationships between entities of the same type. They are particularly useful in modeling hierarchical data structures and networks.

- A relationship between two entities of a similar entity type is called a recursive relationship. Here the same entity type participates more than once in a relationship type with a different role for each instance. In other words, a relationship has always been between occurrences in two different entities. However, the same entity can participate in the relationship. This is termed a **recursive** relationship.
- Recursive relationships are often used to represent hierarchies or networks, where an entity can be connected to other entities of the same type. For example, in an organizational chart, an employee can have a relationship with other employees who are also in a managerial position. Similarly, in a social network, a user can have a relationship with other users who are their friends.
- To represent a recursive relationship in an ER diagram, we use a self-join, which is a join between a table and itself. In other words, we create a relationship between the same entity type. The self-join involves creating two instances of the same entity and connecting them with a relationship. One instance is considered the parent, and the other instance is considered the child.
- We use cardinality constraints to specify the number of instances of the entity that can participate in the relationship. For example, in an organizational chart, an employee can have many subordinates, but each subordinate can only have one manager. This is represented as a one-to-many (1 : N) relationship between the employee entity and itself.
- Overall, recursive relationships are an important concept in ER modeling, and they allow us to represent complex relationships between entities of the same type. They are particularly useful in modeling hierarchical data structures and networks.

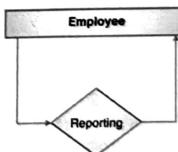
**There can be a lot of examples of recursive relationships**

For Example, if a boy marries a girl and both belong to the human entity set, then this is called the recursive relationship. So, we can indirectly say that if an entity participates more than once in a relationship for different roles, then this is called a recursive relationship. In the above example, the person entity will be used for a boy, and then, on the other hand, it will be used for a girl, and there will be a relation between them.

**Examples**

- We can understand the recursive relationships using the case study of an organization. In an organization, we have employees and a CEO. Any successful organization follows the hierarchy of employees where employees are arranged at their levels, and they report to their supervisors. We can have two roles of relationships in an organization which are: supervisor and subordinate. In an organization, each employee can have a maximum of only one supervisor, but one supervisor can supervise more than one subordinate.
- Let's suppose we have one relationship named 'reporting' between supervisors and subordinates.

So, ER diagram will look like this :



**Fig. 1.16.1 : ER diagram**

**Data Modeling and Visualization**

- The degree of the relationship is one, which means it is a recursive relationship, and a recursive relationship is also called a unary relationship. We can discuss the cardinalities of different roles in the below points :
  - The minimum cardinality of a subordinate can be zero because the CEO of a company cannot be the subordinate.
  - The maximum cardinality of a subordinate can be ONE because, at max, one supervisor is allowed to be a subordinate.
  - The minimum cardinality of a supervisor can be ZERO because subordinates of the companies cannot be supervisors of any other employees.
  - The maximum cardinality of a supervisor can be N, where one supervisor can be allowed to N employees and N can be the total working employees of the organization.
- To manage the relationships, there would be a foreign key as the manager number in each record. So, the sample table's entities would look like this :
  - Emp\_entity(Emp\_no,Emp\_FirstName, Emp\_LastName, Emp\_DOB, Emp\_NI\_Number, Manager\_no);
  - Manager no - (employee no of the employee's manager is managing no);
- We will see a lot of examples of the unary or recursive relationships, which are the following :

**Example 1**

One person can have many children, and the relationship between person and child is unary because they both belong to the same entity type as humans.

**Example 2**

One monitor can handle multiple students in class, and he is itself a student. So the relationship handling is unary or recursive relationship because the monitor and other students belong to the same entity set.

**Example 3**

When the match between team A and team B is organized, then this is also an example of a recursive relationship because both teams belong to the same entity set team, and the relation is 'match'.

**Example 4**

In a machine, one part is built by combining different small parts, so it is also one kind of recursive relationship.

## **1.17 Understand and Model Historical Data**

**What is historical data?**

- Historical data, in a broad context, is data collected about past events and circumstances pertaining to a particular subject.
- By definition, historical data includes most data generated either manually or automatically within an enterprise. Sources, among a great number of possibilities, include press releases, log files, financial reports, project and product documentation and email and other communications.

**How is historical data used and why is it important?**

- In a business context, historical data is used to make important strategic decisions about the present and future. Managers use historical data to track organizational performance over time, identify areas of improvement and make predictions about future trends.
- Businesses are collecting more data than ever and often storing it for longer, both for their own purposes and to satisfy compliance requirements.

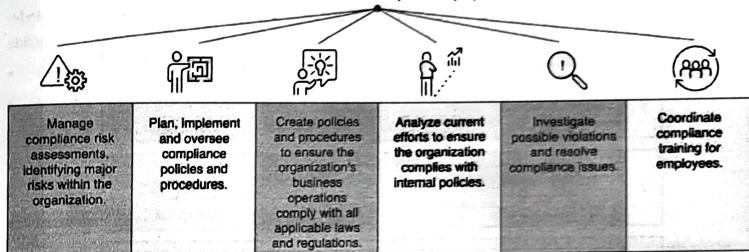
**What does a chief compliance officer do ?****A cco works with leadership and employees to:**

Fig. 1.17.1 : Historical data flow

- Historical data is one of the many tools businesses use to satisfy compliance requirements.
- Historical data can answer important questions such as the following :
  - What were our sales last quarter?
  - How many customer complaints did we receive last year?
  - How has our website traffic changed over the past six months?
- Organizations should have a plan in place for collecting, storing and managing historical data. This will ensure that historical data is available when needed and can be easily accessed and analyzed.

**How is historical data stored?**

- Historical data can be stored in a variety of ways, including databases, spreadsheets and text files. It is important to choose a storage method that is efficient and scalable so historical data can be easily accessed and analyzed.
- When choosing a storage method for historical data, it is important to consider the following questions :
  - How much data will need to be stored?
  - How often will the data need to be accessed?
  - Who will need to access the data?
  - What type of analysis will be performed on the data?
- There are many options available for storing historical data. The best option for a particular organization depends on the specific needs and requirements.

- Introduction to Data Modeling
- Storage capacities have increased significantly in recent years and cloud storage has taken much of the burden of storage administration from many enterprises.
  - However, for large organizations with strict compliance requirements and privacy controls there are also options to store historical data in a private server or data center.
- Data storage tiering hierarchy**
- How many storage tiers an organization has depends on how it classifies data.

Tier	Data category	Data description	Example storage media
0	Mission critical	<ul style="list-style-type: none"> <li>Data that support critical, high-performing workloads that cannot afford delays or disruptions in service</li> <li>Data requirements outweigh storage costs</li> </ul>	<ul style="list-style-type: none"> <li>NVMe SSDs</li> <li>RAM</li> <li>Storage-class memory (e.g. Optane)</li> </ul>
1	Hot data	<ul style="list-style-type: none"> <li>Data that is used continuously to maintain day-to-day business operations</li> <li>Data needs are balanced against storage costs</li> </ul>	<ul style="list-style-type: none"> <li>SSDs</li> <li>High-performing HDDs</li> <li>Hybrid storage systems</li> </ul>
2	Warm data	<ul style="list-style-type: none"> <li>Data that's accessed infrequently or not in constant use but might still be required on occasion</li> <li>Cost considerations are given greater priority</li> </ul>	<ul style="list-style-type: none"> <li>SATA HDDs</li> <li>Backup appliances</li> <li>Tape storage</li> <li>Cloud storage</li> </ul>
3	Cold data	<ul style="list-style-type: none"> <li>Data that is rarely accessed or updated, if at all, or stored only archival purposes</li> <li>Uses the least expensive storage</li> </ul>	<ul style="list-style-type: none"> <li>Slow spinning HDDs</li> <li>Optical discs</li> <li>Tape storage</li> <li>Archival cloud storage</li> </ul>

Fig. 1.17.2

- How organizations store data depends on how it is classified and such factors as compliance requirements.

#### How long should you retain historical data?

- The answer to this question depends on the organization and the type of data being collected.
- Some organizations are required by law or regulation to retain historical data for a certain period of time. For example, businesses in the financial services industry are subject to Financial Industry Regulatory Authority (FINRA) rules, which require the retention of business records for varying lengths of time depending on the type of record.
- Other organizations may have internal policies for how long historical data should be retained. For example a company may keep financial records for seven years or customer service records for five years.
- It is important to consult with legal and compliance teams to determine how long historical data needs to be retained. Once a decision has been made, it is important to have a plan in place for how historical data will be stored and managed over time

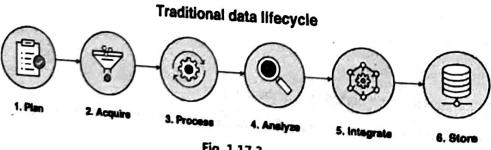


Fig. 1.17.3

- Data lifecycle management helps ensure data is not maintained without good reason or for longer than needed and is appropriately archived or disposed of as appropriate.
- Because data storage requires resources to maintain, Data Lifecycle Management (DLM) is recommended to ensure that data is not maintained without good reason or for longer than necessary and that it is properly archived or disposed of as appropriate.

#### 1.17.1 Benefits of Data Modeling

Developers, Data Architects, Business Analysts, and other stakeholders can examine and comprehend relationships among data in a Database or Data warehouse using Data Modeling. Furthermore, it has the ability to :

- Reduce software and database development errors.
- Increase enterprise-wide consistency in documentation and system architecture.
- Enhance the performance of your application and database.
- Streamline Data mapping across the organization.
- Improve communication between the development and Business Intelligence teams.
- Make database design easier and faster at the conceptual, logical, and physical levels.

#### 1.17.2 Data Modeling Tools

Today, a variety of commercial and free source Computer-Aided Software Engineering (CASE) solutions, including Data modeling, diagramming, and visualization tools, are extensively utilized. Following are a few examples :

- Erwin Data Modeler** is a Data Modelling tool that supports alternative notation approaches, including a dimensional approach, and is based on the Integration DEFinition for Information Modeling (IDEFIX) Data modeling language.
- Enterprise Architect** is a visual modeling and design tool for enterprise information systems, architectures, software applications, and databases. It is built on the foundation of object-oriented languages and standards.
- ER/Studio** is a database design software that works with several common database management systems today. Both relational and dimensional data models are supported.
- Open Source Solutions** like Open Model Sphere are examples of free Data modeling tools.

#### Review Questions

- Q. 1 A program consists of two modules. The number of errors, X, in the first module and the number of errors, Y , in the second module have the joint distribution,  $P(0, 0) = P(0, 1) = P(1, 0) = 0.2$ ,  $P(1, 1) = P(1, 2) = P(1, 3) = 0.1$ ,  $P(0, 2) = P(0, 3) = 0.05$ . Find (a) the marginal distributions of X and Y , (b) the probability of no errors in the first module, and (c) the distribution of the total number of errors in the program. Also, (d) Find out if errors in the two modules occur independently
- Q. 2 Explain in detail Positive, negative, and zero covariance with appropriate graphs

Data Modeling and Visualization

- Q. 3 Explain application of Chebyshev's inequality finance sector by using suitable example
- Q. 4 Explain any two discrete distributions (i) Geometric distribution (ii) Binomial distribution (iii) Bernoulli distribution
- Q. 5 State and explain Central limit theorem
- Q. 6 A disk has free space of 330 megabytes. Is it likely to be sufficient for 300 independent images, if each image has expected size of 1 megabyte with a standard deviation of 0.5 megabytes?
- Q. 7 With neat diagram explain any methods of graphical statistics to visualize data (i) Histograms (ii) stem-and-leaf plots  
(iii) box plots (iv) time plots (v) scatter plots.
- Q. 8 Define and explain maximum likelihood estimation
- Q. 9 What Makes a Good Data Model? Why is the Data Model Important?
- Q. 10 Explain the use of generalization and specialization in a model subtyping using Different levels of generalization
- Q. 11 Discuss various Types of Rules on Recursive Relationships

**2**

# Testing and Data Modeling

**Syllabus****Random Numbers and Simulation :** Sampling of continuous distributions, Monte Carlo methods**Hypothesis Testing :** Type I and II errors, rejection regions; Z-test, T-test, F-test, Chi-Square test, Bayesian test**Stochastic Processes and Data Modeling :** Markov process, Hidden Markov Models, Poisson Process, Gaussian

Processes, Auto-Regressive and Moving average processes, Bayesian Network, Regression, Queuing systems

## 2.1 Random Number and Simulation

- Random numbers are numbers that are generated in a way that cannot be predicted in advance. They are often used in simulations, which are computer models that mimic real-world systems.
- There are two main types of random numbers : true random numbers and pseudorandom numbers. True random numbers are generated by physical processes that are truly unpredictable, such as radioactive decay or atmospheric noise. Pseudorandom numbers are generated by algorithms that produce a sequence of numbers that appears to be random, but is actually deterministic.
- In simulations, random numbers are used to model uncertainty. For example, a simulation of a traffic jam might use random numbers to determine the arrival times of cars at an intersection. This allows the simulation to account for the fact that the arrival times of real cars are not always predictable.
- Random numbers are also used to generate different outcomes in a simulation. For example, a simulation of a game of roulette might use random numbers to determine where the ball lands. This allows the simulation to produce different results each time it is run, which makes it more realistic.
- There are a number of different ways to generate random numbers. Some of the most common methods include :
  - **Linear congruential generators** : These generators use a simple mathematical formula to generate a sequence of numbers.
  - **Mersenne twister generators** : These generators are more complex than linear congruential generators, but they produce better quality random numbers.
  - **Hardware random number generators** : These generators use physical processes to generate random numbers.
- The choice of random number generator depends on the specific application. For example, a simulation that requires very high-quality random numbers might use a hardware random number generator.
- Simulation is a powerful tool that can be used to study complex systems. By using random numbers, simulations can account for uncertainty and produce different outcomes each time they are run. This makes simulations a valuable tool for risk analysis and decision making.

Data Modeling and Visualization

- Here are some examples of how random numbers and simulation are used in real-world applications :
  - **Financial modelling** : Simulations are used to model the behavior of financial markets. This can be used to assess the risk of investments and to make trading decisions.
  - **Engineering design** : Simulations are used to design new products and systems. This can be used to test the performance of designs and to identify potential problems.
  - **Medical research** : Simulations are used to study the spread of diseases and to develop new treatments. This can be used to improve public health and to save lives.

**2.1.1 Sampling Continuous Distribution**

- Sampling of continuous distributions is the process of generating random numbers from a continuous probability distribution. This is often done in simulations, where it is necessary to generate values that follow a particular probability distribution.
- There are a number of different ways to sample from continuous distributions. One common method is to use the inverse transform sampling method. This method involves first finding the inverse of the Cumulative Distribution Function (CDF) of the desired distribution. Then, random numbers are generated from a uniform distribution and then passed through the inverse CDF to obtain values from the desired distribution.
- Another common method for sampling from continuous distributions is the rejection sampling method. This method involves first generating a random number from a proposal distribution. Then, the random number is accepted if it falls within the range of the desired distribution. Otherwise, the random number is rejected and a new one is generated.

**2.1.2 Monte Carlo Methods**

- Monte Carlo sampling methods are a class of techniques for randomly sampling a probability distribution. There are many different Monte Carlo sampling methods, but some of the most common ones for continuous distributions include :
  - **Inverse transform sampling** : This method involves first finding the inverse of the Cumulative Distribution Function (CDF) of the desired distribution. Then, random numbers are generated from a uniform distribution and then passed through the inverse CDF to obtain values from the desired distribution.
  - **Rejection sampling** : This method involves first generating a random number from a proposal distribution. Then, the random number is accepted if it falls within the range of the desired distribution. Otherwise, the random number is rejected and a new one is generated.
  - **Importance sampling** : This method involves first choosing a proposal distribution that is different from the desired distribution. Then, random numbers are generated from the proposal distribution and weighted according to the probability density function of the desired distribution.
  - **Markov chain Monte Carlo (MCMC)** : This method involves constructing a Markov chain whose stationary distribution is the desired distribution. Then, random numbers are generated from the Markov chain.
- Monte Carlo methods are a class of computational algorithms that rely on random sampling and statistical techniques to approximate numerical results. They are widely used in various fields, including physics, engineering, finance, computer science, and more, to solve problems that might be analytically intractable or complex to solve using traditional methods.

Data Modeling and Visualization

- The name "Monte Carlo" originates from the Monte Carlo Casino in Monaco, which is known for its games of chance and randomness. The fundamental idea behind Monte Carlo methods is to simulate random processes and gather statistical information from these simulations to make educated estimations or decisions.
- Here's a general overview of how Monte Carlo methods work :
  - 1. Problem Formulation** : Start with a problem that involves uncertainty or randomness, and that can be modeled using probabilistic or statistical methods.
  - 2. Random Sampling** : Generate a large number of random samples or simulations from the input space of the problem. These samples are often drawn from specific probability distributions related to the problem.
  - 3. Simulation** : Apply the model or algorithm you're investigating to each of the random samples. This could involve solving equations, performing calculations, or running simulations.
  - 4. Statistical Analysis** : Analyze the results of the simulations statistically to estimate properties of interest, such as the mean, variance, probabilities, or other relevant metrics.
  - 5. Inference and Approximation** : Use the collected statistical data to approximate solutions, make predictions, or gain insights about the original problem that would be difficult or impossible to obtain analytically.
- Examples of Monte Carlo methods in different domains :
  - **Estimating Pi** : The "Buffon's Needle" experiment is a classic example of using Monte Carlo methods to estimate the value of pi. By randomly dropping needles on a surface with parallel lines drawn on it and calculating the ratio of needles crossing the lines to the total number of needles dropped, one can approximate the value of pi.
  - **Integral Approximation** : Monte Carlo integration is a technique for numerically approximating the value of definite integrals. It involves generating random points within a bounded region and using the proportion of points falling under a curve to estimate the integral's value.
  - **Finance** : Monte Carlo methods are widely used in finance to model the behavior of financial instruments, calculate option prices, and simulate stock price movements under uncertain conditions.
  - **Physics** : In particle physics, Monte Carlo simulations are used to model the behavior of subatomic particles and predict outcomes of complex interactions.
  - **Optimization** : Monte Carlo optimization methods use random sampling to explore the solution space of optimization problems and find near-optimal solutions.
  - **Machine Learning** : In reinforcement learning, Monte Carlo methods are used to estimate the value of state-action pairs in a dynamic environment.
- Monte Carlo methods can provide accurate approximations, but their accuracy often increases with the number of simulations performed. As computational power has increased, the application of Monte Carlo methods has become more widespread and sophisticated.

- Monte Carlo methods to estimate lengths, areas, and volumes of different shapes. Let's start with each of these cases :

## 1. Estimating Length

To estimate the length of a curve or a line segment using Monte Carlo methods, you can follow these steps:

- Consider the curve or line segment you want to measure, and the bounding box that contains it.
- Generate random points within the bounding box.
- Count the number of points that fall within the curve or line segment.
- Calculate the length estimate using the formula: Length Estimate = (Number of Points Inside Curve / Total Number of Points) \* Area of Bounding Box.
- As you increase the number of random points, your length estimate should become more accurate.

## 2. Estimating Area

To estimate the area of a region using Monte Carlo methods, follow these steps:

- Consider the region you want to measure and the bounding rectangle or square that encloses it.
- Generate random points within the bounding rectangle.
- Count the number of points that fall within the region.
- Calculate the area estimate using the formula: Area Estimate = (Number of Points Inside Region / Total Number of Points) \* Area of Bounding Rectangle.
- This method works for irregularly shaped regions as well.

## 3. Estimating Volume

To estimate the volume of a three-dimensional shape using Monte Carlo methods, you can adapt the approach used for estimating area. Here's how:

- Consider the three-dimensional shape and the bounding box that encloses it.
- Generate random points within the bounding box.
- Count the number of points that fall within the shape.
- Calculate the volume estimate using the formula :  

$$\text{Volume Estimate} = (\text{Number of Points Inside Shape}/\text{Total Number of Points}) * \text{Volume of Bounding Box.}$$
- Keep in mind that the accuracy of your estimates improves as you increase the number of random points used in the simulations.
- The key idea in all of these cases is to use random sampling to approximate the ratio of the desired quantity (length, area, or volume) to the total space in which the shape is contained. By doing so, you to the total number of points generated.
- It's important to note that the accuracy of Monte Carlo estimates depends on the number of random samples used. The more samples you generate, the more accurate your estimate will be, but this may also require more computational resources.

## 2.2 Hypothesis Testing

- Hypothesis testing is a statistical method used to determine whether there is enough evidence in a sample data to draw conclusions about a population. It involves formulating two competing hypotheses, the null hypothesis ( $H_0$ ) and the alternative hypothesis ( $H_a$ ), and then collecting data to assess the evidence.
- The null hypothesis is a statement of no difference, while the alternative hypothesis is a statement of some difference. For example, the null hypothesis might be that there is no difference in the average height of men and women, while the alternative hypothesis might be that the average height of men is greater than the average height of women.

### 2.2.1 Type I and Type II Errors

- In statistics, a **Type I error** is a false positive conclusion, while a **Type II error** is a false negative conclusion.
- Making a statistical decision always involves uncertainties, so the risks of making these errors are unavoidable in hypothesis testing.
- The probability of making a Type I error is the significance level, or alpha ( $\alpha$ ), while the probability of making a Type II error is beta ( $\beta$ ). These risks can be minimized through careful planning in your study design.

Type and Type II Error		
Null hypothesis is ...	True	False
Rejected	Type I error False negative Probability = $\alpha$	Correct decision True positive Probability = $1 - \alpha$
Not Rejected	Correct decision True negative Probability = $1 - \beta$	Type II error False negative Probability = $\beta$

### 2.2.1(A) Type I Error

- A Type I error means rejecting the null hypothesis when it's actually true. It means concluding that results are **statistically significant** when, in reality, they came about purely by chance or because of unrelated factors.
- The risk of committing this error is the significance level (alpha or  $\alpha$ ) you choose. That's a value that you set at the beginning of your study to assess the statistical probability of obtaining your results (p value).
- The significance level is usually set at 0.05 or 5%. This means that your results only have a 5% chance of occurring, or less, if the null hypothesis is actually true.
- If the p value of your test is lower than the significance level, it means your results are statistically significant and consistent with the alternative hypothesis. If your p value is higher than the significance level, then your results are considered statistically non-significant.

Data Modeling and Visualization**Type I error rate**

- The null hypothesis distribution curve in Fig. 2.2.1 shows the probabilities of obtaining all possible results if the study were repeated with new samples and the null hypothesis were true in the population.
- At the tail end, the shaded area represents alpha. It's also called a critical region in statistics.
- If your results fall in the critical region of this curve, they are considered statistically significant and the null hypothesis is rejected. However, this is a false positive conclusion, because the null hypothesis is actually true in this case.

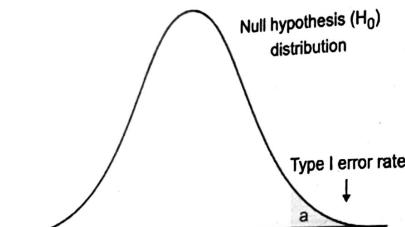
**Probability of making a Type I error**

Fig. 2.2.1 : Probability of making Type-I Error

**2.2.1(B) Type II Error**

- A Type II error means not rejecting the null hypothesis when it's actually false. This is not quite the same as "accepting" the null hypothesis, because hypothesis testing can only tell you whether to reject the null hypothesis.
- Instead, a Type II error means failing to conclude there was an effect when there actually was. In reality, your study may not have had enough statistical power to detect an effect of a certain size.
- Power is the extent to which a test can correctly detect a real effect when there is one. A power level of 80% or higher is usually considered acceptable.
- The risk of a Type II error is inversely related to the statistical power of a study. The higher the statistical power, the lower the probability of making a Type II error.

**Type II error rate**

- The alternative hypothesis distribution curve in Fig. 2.2.2 shows the probabilities of obtaining all possible results if the study were repeated with new samples and the alternative hypothesis were true in the population.
- The Type II error rate is beta ( $\beta$ ), represented by the shaded area on the left side. The remaining area under the curve represents statistical power, which is  $1 - \beta$ .
- Increasing the statistical power of your test directly decreases the risk of making a Type II error.

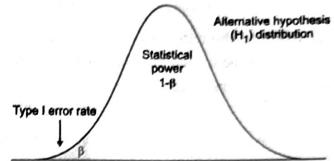
**Probability of making a Type II error**

Fig. 2.2.2 : Probability of making Type-II Error

**2.2.2 Rejection Region**

- A rejection region, also known as a critical region, is a set of values for the test statistic for which the null hypothesis is rejected. In other words, if the observed test statistic is in the rejection region, then we reject the null hypothesis and accept the alternative hypothesis.
- The rejection region is determined by the significance level of the test. The significance level is the probability of rejecting the null hypothesis when it is actually true. The most common significance levels are 0.05 and 0.01. This means that there is a 5% or 1% chance of rejecting the null hypothesis when it is actually true.
- The rejection region is typically calculated using the sampling distribution of the test statistic. The sampling distribution is the distribution of the test statistic under the assumption that the null hypothesis is true. The rejection region is typically set so that the probability of the test statistic falling into the rejection region by chance is equal to the significance level.

**2.2.3 Z Test**

A z-test is a statistical hypothesis test used to compare a sample mean to a known population mean when the population standard deviation is known. It is appropriate when the sample size is sufficiently large (typically  $n > 30$ ) or when the population follows a normal distribution. The z-test is based on the standard normal distribution, which has a mean of 0 and a standard deviation of 1.

The general process of performing a z-test involves the following steps :

- 1) Formulate the null and alternative hypotheses :** The null hypothesis ( $H_0$ ) typically states that there is no significant difference between the sample mean and the population mean. The alternative hypothesis ( $H_1$ ) asserts that there is a significant difference.
- 2) Select a significance level (alpha,  $\alpha$ ) :** The significance level determines the probability of making a Type I error (rejecting the null hypothesis when it is true). Commonly used significance levels are 0.05 (5%) and 0.01 (1%).
- 3) Calculate the test statistic (z-score) :** The z-score is computed using the formula :

$$z = \frac{(\bar{x} - \mu)}{(\sigma / \sqrt{n})}$$

where :

$\bar{x}$  is the sample mean,

$\mu$  is the population mean (the known value),

$\sigma$  is the population standard deviation (the known value),  $n$  is the sample size.

## Data Modeling and Visualization

- 4) Determine the critical value :** Based on the chosen significance level and assuming the null hypothesis is true, find the critical z-value from the standard normal distribution table.
- 5) Compare the test statistic (z-score) to the critical value :** If the absolute value of the z-score is greater than the critical value (in the rejection region), then the null hypothesis is rejected in favor of the alternative hypothesis. If the z-score falls within the non-rejection region, there is insufficient evidence to reject the null hypothesis.

The z-test is commonly used when the population standard deviation is known, and it is often used for large sample sizes or when conducting tests related to proportions.

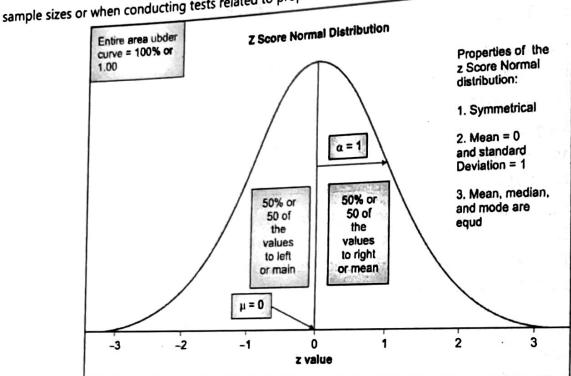


Fig. 2.2.3 : Normal Distribution of Z-Test

#### 2.2.4 t Test

- A t-test is a statistical hypothesis test used to compare the means of two groups and determine if there is a significant difference between them. It is appropriate when comparing the means of two independent samples, and it is commonly used when the sample size is small (typically  $n < 30$ ) or when the population standard deviation is unknown. The t-test is based on the t-distribution, which takes into account the variability introduced by using the sample standard deviation.
- There are two main types of t-tests:
  - Independent Samples t-test :** This test is used when comparing the means of two separate and independent groups. The groups could be from different populations, have different treatments, or be subject to different conditions.
  - Paired Samples t-test** (also known as Dependent Samples t-test) : This test is used when comparing the means of two related or paired samples. The samples are not independent, and each observation in one sample is directly related to a specific observation in the other sample.

- The general process of performing a t-test involves the following steps :
  - Formulate the null and alternative hypotheses :** The null hypothesis ( $H_0$ ) typically states that there is no significant difference between the means of the two groups, while the alternative hypothesis ( $H_a$ ) asserts that there is a significant difference.
  - Select a significance level (alpha,  $\alpha$ ) :** The significance level determines the probability of making a Type I error (rejecting the null hypothesis when it is true). Commonly used significance levels are 0.05 (5%) and 0.01 (1%).
  - Calculate the test statistic (t-value) :** The t-value is computed using the formula :

$$t = \frac{\bar{x}_1 - \bar{x}_2}{(s_{\text{pool}} * \sqrt{1/n_1 + 1/n_2})}$$

where :

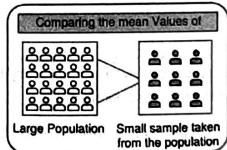
$\bar{x}_1$  and  $\bar{x}_2$  are the sample means of the two groups,

$s_{\text{pool}}$  is the pooled standard deviation, calculated based on the sample standard deviations and sample sizes,

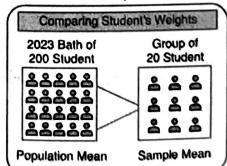
$n_1$  and  $n_2$  are the sample sizes of the two groups.

- Determine the degrees of freedom (df) :** The degrees of freedom are used to find the critical t-value from the t-distribution table.
- Compare the test statistic (t-value) to the critical value :** If the absolute value of the t-value is greater than the critical value (in the rejection region), then the null hypothesis is rejected in favor of the alternative hypothesis. If the t-value falls within the non-rejection region, there is insufficient evidence to reject the null hypothesis.
- The t-test is widely used in various fields for comparing means, such as in clinical trials, social sciences, economics, and more.

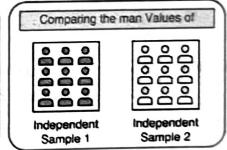
#### One-Sample t-Test



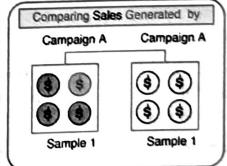
Example:



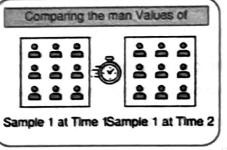
#### Two-Sample t-Test



Example:



#### Paired-Sample t-Test



Example:

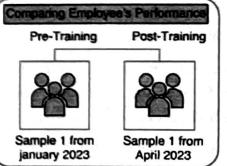


Fig. 2.2.4 : Sample T-test

### 2.2.5 Chi Square test

The chi-square test is a statistical hypothesis test used to determine if there is a significant association between two categorical variables. It is often applied to analyze data in contingency tables, where data is organized into rows and columns representing the categories of the two variables being studied.

There are different types of chi-square tests, depending on the nature of the data and the research question:

#### Chi-Square Test of Independence

This test is used when you want to determine if there is a significant association between two categorical variables, and the variables are independent of each other. For example, you might want to investigate if there is a relationship between gender (male/female) and the preference for a certain product (yes/no).

#### Chi-Square Test of Goodness-of-Fit

This test is used when you want to assess how well the observed data fits a theoretical or expected distribution. For example, you may want to know if the distribution of blood types in a population follows the expected distribution based on genetics.

The general process of performing a chi-square test involves the following steps:

##### 1) Formulate the null and alternative hypotheses

The null hypothesis ( $H_0$ ) typically states that there is no significant association between the two variables, while the alternative hypothesis ( $H_a$ ) asserts that there is a significant association.

##### 2) Select a significance level (alpha, $\alpha$ )

The significance level determines the probability of making a Type I error (rejecting the null hypothesis when it is true). Commonly used significance levels are 0.05 (5%) and 0.01 (1%).

##### 3) Create a contingency table

Organize the data into a two-dimensional table (contingency table) that shows the frequency of observations for each combination of categories of the two variables.

##### 4) Calculate the test statistic (chi-square statistic)

The chi-square statistic measures the difference between the observed frequencies in the contingency table and the expected frequencies under the assumption of independence or the expected distribution.

##### 5) Determine the degrees of freedom (df)

The degrees of freedom depend on the size of the contingency table and are used to find the critical chi-square value from the chi-square distribution table.

##### 6) Compare the test statistic (chi-square value) to the critical value

If the chi-square value exceeds the critical value (in the rejection region), then the null hypothesis is rejected in favor of the alternative hypothesis. If the chi-square value falls within the non-rejection region, there is insufficient evidence to reject the null hypothesis.

The chi-square test is commonly used in various fields, including social sciences, biology, medicine, and market research, to analyze categorical data and determine if there are significant relationships between variables.

### Example

	Observed			Total
	rep.	dem.	other	
male	26	13	5	44
female	20	29	7	56
	46	42	12	100

	Expected			
	rep.	dem.	other	
male	20.24	18.48	5.28	
female	25.76	23.52	6.72	

#### Chi-square test

$H_0$  • Political party

• independent upon gender

$H_1$  • Political party

• dependent upon gender



Fig. 2.2.5 : Chi-Square Test

### 2.2.6 F Test

- The F-test is a statistical hypothesis test used to compare the variances of two or more groups or populations. It is commonly used in analysis of variance (ANOVA) and regression analysis to assess whether the variability between the groups is significantly different from the variability within the groups.
- The F-test is based on the F-distribution, which is a probability distribution that arises when comparing the ratios of two independent chi-square distributions. The F-distribution has two parameters, degrees of freedom for the numerator ( $df_1$ ) and degrees of freedom for the denominator ( $df_2$ ).
- There are different types of F-tests, depending on the context in which it is used :
  - One-Way ANOVA F-test :** This test is used when comparing the means of two or more groups to determine if there is a significant difference between their population variances. The one-way ANOVA compares the variance among the groups (explained variance) to the variance within the groups (unexplained variance).
  - Two-Way ANOVA F-test :** In two-way ANOVA, the F-test is used to examine the interaction effect between two independent factors and their influence on the dependent variable.
  - F-test in Regression Analysis :** In regression analysis, the F-test is used to assess the overall significance of a linear regression model. It compares the variation explained by the regression model to the residual variation (unexplained variation).
- The general process of performing an F-test involves the following steps :
  - Formulate the null and alternative hypotheses :** The null hypothesis ( $H_0$ ) typically states that there is no significant difference between the variances, while the alternative hypothesis ( $H_a$ ) asserts that there is a significant difference.
  - Select a significance level (alpha,  $\alpha$ ) :** The significance level determines the probability of making a Type I error (rejecting the null hypothesis when it is true). Commonly used significance levels are 0.05 (5%) and 0.01 (1%).

- Data Modeling and Visualization**
- 3) **Calculate the test statistic (F-value)** : The F-value is computed by dividing the variance among the groups (or model variance) by the variance within the groups (or residual variance).
  - 4) **Determine the degrees of freedom ( $df_1$  and  $df_2$ )** : The degrees of freedom depend on the specific test being conducted and are used to find the critical F-value from the F-distribution table.
  - 5) **Compare the test statistic (F-value) to the critical value** : If the F-value exceeds the critical value (in the rejection region), then the null hypothesis is rejected in favor of the alternative hypothesis. If the F-value falls within the non-rejection region, there is insufficient evidence to reject the null hypothesis.
  - The F-test is a powerful tool for comparing variances and assessing the significance of relationships in various statistical analyses.

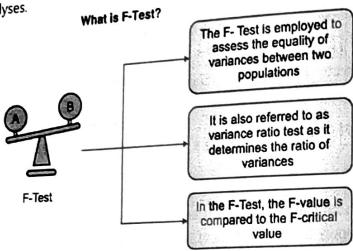


Fig. 2.2.6 : F-test

## 2.2.7 Bayesian Test

- In Bayesian statistics, hypothesis testing is performed in a different way compared to classical or frequentist statistics. Bayesian hypothesis testing involves updating prior beliefs about the parameters of interest using observed data to obtain a posterior distribution, which represents the updated beliefs after incorporating the data.
- The main steps in Bayesian hypothesis testing are as follows :

  - 1) **Prior Distribution** : Before any data is collected, the analyst starts with a prior distribution, which represents their beliefs about the parameters of interest based on previous knowledge, expert opinions, or other relevant information. The prior distribution reflects what the analyst believes about the parameters before observing any data.
  - 2) **Likelihood Function** : The likelihood function is a statistical model that represents the probability of observing the data given the parameters of interest. It describes the chance of getting the observed data under different parameter values.
  - 3) **Bayes' Theorem** : Bayes' theorem is used to update the prior beliefs to the posterior beliefs after observing the data. It mathematically combines the prior distribution and the likelihood function to obtain the posterior distribution.
  - Posterior  $\propto$  Prior  $\times$  Likelihood
  - 4) **Posterior Distribution** : The posterior distribution represents the updated beliefs about the parameters after incorporating the observed data. It combines the prior information and the information contained in the data to provide a complete probability distribution for the parameters.

- 5) **Bayesian Decision Rule** : To make decisions based on the posterior distribution, one can use various approaches, such as choosing the parameter value with the highest probability (maximum a posteriori estimation) or computing credible intervals to summarize the uncertainty about the parameter estimates.
- Bayesian hypothesis testing offers several advantages, including the ability to incorporate prior knowledge, the ability to handle small sample sizes, and the straightforward interpretation of results in terms of probabilities. However, it requires specifying appropriate prior distributions, which can sometimes be subjective and challenging, especially when little prior information is available.
- Overall, Bayesian hypothesis testing provides a different perspective on hypothesis testing, allowing for a more flexible and probabilistic approach to drawing conclusions about parameters of interest. It is increasingly used in various fields, especially when dealing with complex models and limited data.

## 2.3 Stochastic Processes and Data Modeling

- Stochastic processes are a type of mathematical object that describes the evolution of a system over time. They are used to model systems that are subject to random variation, such as stock prices, traffic flow, and weather patterns.
- Data modeling is the process of creating a mathematical model of a system from data. This model can then be used to make predictions about the future behavior of the system.
- Stochastic processes and data modeling are closely related. Stochastic processes can be used to generate data, and data modeling can be used to fit stochastic processes to data.

### 2.3.1 Markov Process

- A Markov process is a stochastic process where the probability of transitioning to a particular state only depends on the current state. This means that the past history of the process does not affect the future. Markov processes are often used to model systems that exhibit memoryless behavior, such as the arrival of customers at a store.
- Formally, a Markov process is defined as a set of states and a transition matrix that specifies the probabilities of moving from one state to another in discrete time steps. The properties of a Markov process can be described as follows :
  - o **State Space** : The set of all possible states that the process can be in is called the state space. Each state represents a distinct situation or condition.
  - o **Markov Property** : The Markov property implies that the future state of the process is independent of the past states, given the present state. In other words, the probability of transitioning to a particular state depends solely on the current state and not on the sequence of states that led to the current state.
  - o **Transition Matrix** : The transition matrix, often denoted by  $P$ , is a square matrix where each element  $P(i, j)$  represents the probability of moving from state  $i$  to state  $j$  in one time step. The elements of each row sum up to 1, as the process must move to one of the possible states.
  - o **Homogeneity** : The probabilities of transitioning from one state to another remain constant over time. This property is often referred to as the time-homogeneity of the Markov process.

## Data Modeling and Visualization

- Markov processes find applications in various fields, such as physics, biology, economics, finance, and computer science. They are particularly useful for modeling systems with discrete states and random transitions, where memorylessness is a reasonable assumption. Examples of real-life applications of Markov processes include modeling weather patterns, stock price movements, queueing systems and natural language processing tasks.

## 2.3.2 Hidden Markov Model (HMM)

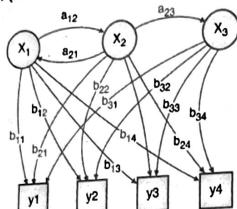


Fig. 2.3.1 : Probabilistic parameters of a hidden Markov model

- A hidden Markov model (HMM) is a type of Markov process where the current state is not directly observable. Instead, the state is hidden and can only be inferred from the observed data. HMMs are often used to model systems where the underlying state is not directly measurable, such as the speech of a person or the behavior of a financial market.
- It is widely used in various fields, including speech recognition, natural language processing, bioinformatics, finance, and more, where there is an interest in modeling sequential data with underlying hidden structure.
- The basic concept of a Hidden Markov Model involves two main components :
  - Hidden States :** The system is assumed to have a set of hidden states, which represent the underlying unobservable processes. These states form a Markov chain, meaning that the probability of transitioning from one hidden state to another depends only on the current state and not on the past history. Hidden states are not directly observable.
  - Observable Symbols :** Each hidden state emits observable symbols or observations according to certain probability distributions. These observations provide information about the underlying hidden state. The relationship between the hidden states and the observable symbols is modeled using emission probabilities.
- Formally, an HMM is defined by the following components :
  - State Space :** The set of all possible hidden states.
  - Transition Probabilities :** The probabilities of transitioning from one hidden state to another. These probabilities are represented in a transition matrix.
  - Emission Probabilities :** The probabilities of observing a particular symbol or observation given a hidden state. These probabilities are represented in an emission matrix.

- The process of using an HMM typically involves two main tasks :

- Learning :** Given a sequence of observations, the goal is to estimate the parameters of the HMM (transition probabilities and emission probabilities) that best fit the observed data. This is often done using algorithms like the Expectation-Maximization (EM) algorithm or the Baum-Welch algorithm.
  - Inference :** Given the learned HMM parameters, the task is to infer the most likely sequence of hidden states that generated the observed data. This is achieved using algorithms like the Viterbi algorithm or the Forward-Backward algorithm.
- HMMs are powerful tools for modeling and analyzing sequential data with hidden structures. They allow for efficient representation of complex systems, and their ability to deal with unobserved states makes them well-suited for a wide range of applications where there is uncertainty or missing information.

## Example

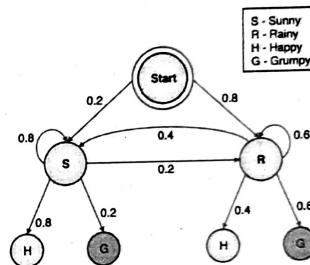


Fig. 2.3.2 : HMM

## 2.3.3 Poisson Process

- A Poisson process is a stochastic process where the number of events occurring in a given time interval is a random variable that follows a Poisson distribution. Poisson processes are often used to model the occurrence of events over time, such as the number of phone calls received in a call center or the number of customers arriving at a store.
- A Poisson process is a stochastic process that models the occurrence of events over time or space. It is named after the French mathematician Siméon Denis Poisson, who studied this process in the early 19th century. The Poisson process is widely used in various fields, including telecommunications, queuing theory, biology, and reliability analysis, where the occurrences of events are random and independent.
- The key characteristics of a Poisson process are as follows :
  - Time or Space :** The Poisson process can be defined either in continuous time (time-based Poisson process) or in discrete space (space-based Poisson process). In the continuous-time case, it models events occurring in continuous time intervals, while in the discrete-space case, it models events occurring at specific spatial locations.
  - Independent Increments :** The occurrences of events are assumed to be independent of each other. In other words, the probability of an event occurring in a particular interval or location is not influenced by whether an event occurred in any other interval or location.

Data Modeling and Visualization

- 3) Stationary :** The rate of event occurrences is constant over time or space. This means that the average rate of events per unit of time or space remains the same throughout the process.
- 4) Homogeneity :** The process is time-homogeneous, meaning that the probability of an event occurring in a given interval or location is the same for all time or space.
- Formally, the Poisson process is characterized by a single parameter, denoted by  $\lambda$  (lambda), which represents the average rate of event occurrences per unit of time or space. The probability of observing  $k$  events in a time interval or at a specific location is given by the Poisson distribution :
- $$P(X = k) = \frac{e^{(-\lambda)} \cdot \lambda^k}{k!}$$
- where  $X$  is the random variable representing the number of events,  $e$  is the base of the natural logarithm (approximately 2.71828), and  $k!$  denotes the factorial of  $k$ .
- The Poisson process is often used to model rare events or occurrences that happen at a constant rate. It is widely used in queuing theory to analyze the behavior of waiting lines, in telecommunications to study the arrival of packets, and in biology to model the distribution of mutations or the occurrence of certain phenomena.



Fig. 2.3.3 : Poisson process

**2.3.4 Gaussian Process**

- A Gaussian process is a stochastic process where the values of the process are normally distributed. Gaussian processes are often used to model uncertainty in data.
- Auto-regressive and moving average processes : Auto-regressive (AR) and moving average (MA) processes are stochastic processes that are used to model the dependence of the current value of a process on its past values. AR processes are models of the past values of the process, while MA processes are models of the errors in the past values of the process.
- In a Gaussian process, a function is represented as a collection of random variables, and any finite subset of these random variables follows a joint Gaussian distribution. In other words, a Gaussian process defines a distribution over functions, rather than a single function. This distribution is determined by two main components :
  - Mean Function :** The mean function represents the average behavior of the underlying function and is often assumed to be zero (although it can be non-zero in some cases). It captures the overall trend of the data.
  - Covariance Function (Kernel) :** The covariance function, also known as the kernel function, characterizes the relationship between data points and measures the similarity or correlation between function values at different points. It quantifies how the function values change as the input points change. Different choices of the covariance function lead to different Gaussian process models, allowing for flexibility in modeling various patterns and behaviors.

The key properties of a Gaussian process are :

**1) Flexibility**

Gaussian processes can model a wide range of functions, from smooth to highly non-linear, depending on the choice of the kernel function.

**2) Uncertainty Estimation**

In addition to predicting the most likely function, Gaussian processes provide a measure of uncertainty in the predictions, which is crucial in decision-making and for quantifying confidence intervals.

**3) Bayesian Framework**

Gaussian processes are inherently Bayesian, meaning that prior beliefs can be incorporated into the model, and uncertainty is propagated through to the predictions.

**4) Non-parametric**

- Gaussian processes are considered non-parametric models because the number of parameters increases with the size of the dataset, allowing them to adapt to the complexity of the data.
- Gaussian processes have applications in various fields, including regression, time series analysis, spatial modeling, optimization, and Bayesian optimization, among others. They are particularly useful when data is limited, noisy, or when the underlying function is not well-defined.
- Although Gaussian processes are computationally intensive for large datasets due to their non-parametric nature, they remain a popular choice in many applications due to their flexibility, interpretability, and uncertainty-aware predictions. Various approximation techniques, such as sparse GPs or kernel approximation methods, are employed to scale Gaussian processes to larger datasets.

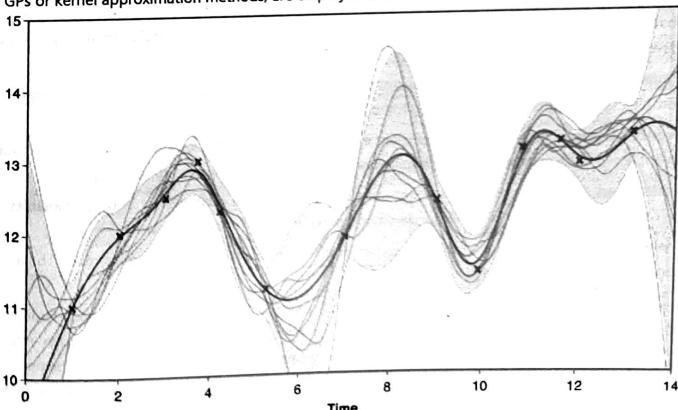


Fig. 2.3.4 : Gaussian process

**Data Modeling and Visualization****2.3.5 Bayesian Network**

- A Bayesian network is a probabilistic graphical model that represents the joint probability distribution of a set of variables. Bayesian networks are often used to model uncertainty in data and to make predictions about the future.
- It provides a powerful framework for modeling uncertainty and making predictions or inferences based on available evidence. Bayesian networks are widely used in artificial intelligence, machine learning, decision analysis, and expert systems.
- The key components of a Bayesian network are :
  - Nodes** : Each node in the Bayesian network represents a random variable. These variables can be observable (e.g., sensor measurements) or unobservable (e.g., latent variables or hidden states).
  - Directed Edges (Arrows)** : Directed edges (arrows) between nodes represent probabilistic dependencies. If there is an arrow from Node A to Node B, it indicates that Node B is conditionally dependent on Node A. In other words, the value of Node B is influenced by the value of Node A.
  - Conditional Probability Tables (CPTs)** : Each node in the Bayesian network has an associated conditional probability table (CPT), which specifies the conditional probabilities of the node given its parents (nodes with incoming arrows). These CPTs capture the probabilistic relationships between nodes in the network.
- The Bayesian network encodes the joint probability distribution of all the variables in the system using the chain rule of probability. This allows for efficient inference and reasoning about the relationships between variables and their probabilities.

The main advantages of Bayesian networks are :

- Probabilistic Reasoning** : Bayesian networks allow for explicit modeling of uncertainty and provide a formal framework for reasoning under uncertainty.
- Modular Representation** : The graphical structure of Bayesian networks allows for a modular representation of complex systems, making it easier to understand and update models.
- Inference** : Bayesian networks enable efficient inference and prediction, such as computing the probability of an event given evidence (posterior probability) or the most likely explanation for observed data.
- Causality and Explanation** : Bayesian networks can model causal relationships between variables, making them useful for causal reasoning and understanding the impact of interventions.
- Learning** : Bayesian networks can be learned from data, allowing for automatic model building from observed relationships.

Applications of Bayesian networks include medical diagnosis, fault diagnosis, risk assessment, natural language processing, image recognition, and many other areas where uncertainty and complex relationships need to be modeled and analyzed. Bayesian networks have become a valuable tool in decision support systems and expert systems, where they help to combine domain knowledge with data-driven reasoning to make informed decisions.

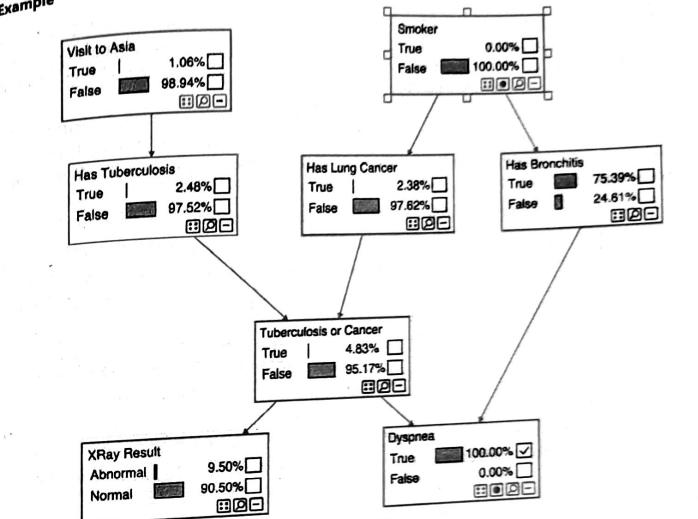
**Data Modeling and Visualization****Example**

Fig. 2.3.5

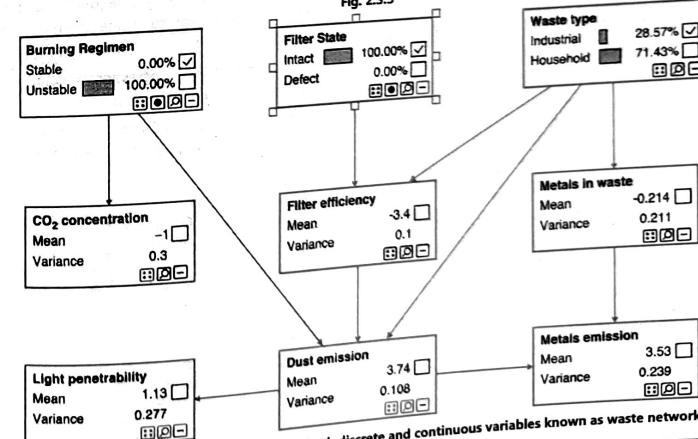


Fig. 2.3.6 : A Simple Bayesian network with both discrete and continuous variables known as waste network

**2.3.6 Regression**

- Regression is a statistical method used to model the relationship between a dependent variable (also known as the outcome or response variable) and one or more independent variables (also known as predictor variables or features). The goal of regression analysis is to understand how changes in the independent variables are associated with changes in the dependent variable.
- The fundamental idea behind regression is to fit a mathematical function (the regression model) to the observed data in such a way that it provides the best possible approximation of the relationship between the variables. This allows us to make predictions about the dependent variable based on the values of the independent variables.
- There are several types of regression models, with the most common ones being :
  - Simple Linear Regression :** This involves modeling the relationship between one dependent variable and one independent variable. It assumes that the relationship can be approximated by a straight line equation ( $y = mx + b$ ), where  $y$  is the dependent variable,  $x$  is the independent variable,  $m$  is the slope, and  $b$  is the intercept.
  - Multiple Linear Regression :** This extends the simple linear regression to model the relationship between one dependent variable and two or more independent variables. The equation becomes  $y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$ , where  $b_0$  is the intercept,  $b_1, b_2, \dots, b_n$  are the coefficients, and  $x_1, x_2, \dots, x_n$  are the independent variables.
  - Polynomial Regression :** In this type of regression, the relationship between the variables is approximated by a polynomial equation of a higher degree than one. It is useful when the relationship is nonlinear.
  - Logistic Regression :** This is used when the dependent variable is binary or categorical. Logistic regression estimates the probability of the dependent variable being in a particular category based on the values of the independent variables.
  - Ridge Regression and Lasso Regression :** These are variations of linear regression that include regularization techniques to prevent overfitting and improve the model's generalization.
- Regression analysis is widely used in various fields, including economics, finance, social sciences, engineering, and machine learning. It helps researchers and analysts to understand the underlying relationships between variables, make predictions, and draw conclusions based on the observed data.

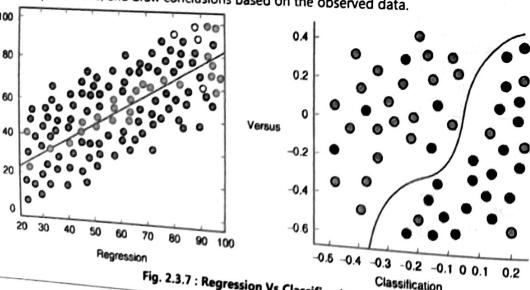


Fig. 2.3.7 : Regression Vs Classification

**2.3.7 Queuing Systems**

- A queuing system, also known as a queuing model or waiting line system, is a mathematical model used to study and analyze the behavior of waiting lines or queues. Queuing systems are prevalent in various real-world scenarios, such as traffic management, customer service centers, computer systems, telecommunications, transportation, and manufacturing.
- A queuing system consists of the following components :
  - A queuing system, also known as a queuing system, is a mathematical model that is used to analyze and study the behavior of entities (such as customers, data packets, or tasks) as they wait in line (queue) for service. Queuing systems are used to understand and optimize a wide range of real-world scenarios where entities arrive at a service point or facility and wait their turn to be served.
- Key components of a queuing system include :
  - Arrival Process :** This refers to how entities arrive at the system. Arrivals can follow various patterns, such as random arrival times, fixed intervals, or scheduled arrivals.
  - Service Process :** The service process describes how entities are served once they reach the front of the queue. Service times can also follow different patterns, such as exponential, uniform, or other distributions.
  - Queue Discipline :** This determines the order in which entities are served from the queue. Common queue disciplines include first-come-first-served (FCFS), last-come-first-served (LCFS), priority-based, etc.
  - Queue Capacity :** This specifies the maximum number of entities that the queue can accommodate at any given time. If the queue is full, arriving entities may be turned away or managed through some overflow mechanism.
  - Queue Length :** The current number of entities in the queue waiting to be served.
  - Service Facilities :** These are the resources responsible for serving the entities. Service facilities can be a single server or multiple servers, each with its own characteristics.
  - Performance Metrics :** Queuing systems are analyzed based on various performance metrics, including average waiting time, average queue length, utilization of the service facility, and more.
  - Queueing Models :** Queueing systems are often represented and analyzed using mathematical models. Different models are used based on the specifics of the system, such as the number of servers, the arrival and service time distributions, and the queue discipline.
- Examples of queuing systems can be found in various fields :
  - Telecommunications :** Call centers, data networks, and internet routers often use queuing models to optimize resource allocation and minimize waiting times.
  - Retail :** Checkout lines in supermarkets and other retail stores can be modeled using queuing systems to improve customer flow and reduce wait times.
  - Manufacturing :** Production lines with limited resources and processes that require sequential steps can benefit from queuing analysis to optimize efficiency.

- 4. Healthcare :** Hospitals and clinics use queuing models to manage patient flow, appointment scheduling, and resource allocation.
- 5. Transportation :**
- o Traffic flow, public transportation systems, and airport security lines can be analyzed using queuing techniques to reduce congestion and improve overall efficiency.
  - o Queuing theory, the mathematical foundation of queuing systems, provides insights into the behavior of these systems and helps decision-makers make informed choices to optimize performance and resource allocation.
  - o The study of queuing systems involves the application of queuing theory, which is a branch of operations research and applied probability. Queuing theory uses mathematical models to describe and analyze the behavior of waiting lines under different conditions and assumptions. Various analytical methods, simulation techniques, and numerical methods are employed to solve queuing models and obtain relevant performance metrics.
  - o Queuing theory has wide-ranging applications and is crucial for decision-making in designing and managing systems with waiting lines, ensuring resource allocation, and providing efficient customer service in various industries and domains.

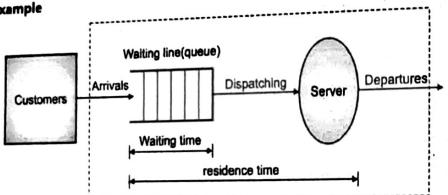
**Queuing System Example**

Fig. 2.3.8 : Queuing System

- Little's Law is a fundamental theorem in queuing theory that provides a relationship between the average number of entities in a queuing system, the average arrival rate of entities, and the average time those entities spend in the system. It's named after John D.C. Little, who introduced the law in 1961.
- Mathematically, Little's Law can be expressed as :

Average Number of Entities in the System = Average Arrival Rate \* Average Time in the System

In symbols :

$$L = \lambda * W$$

Where :

L is the average number of entities (customers, packets, tasks, etc.) in the queuing system.

$\lambda$  (lambda) is the average arrival rate of entities to the system.

W is the average time that an entity spends in the system.

- This law applies to a wide range of queuing systems, as long as the system is stable (the arrival rate is less than or equal to the service rate) and the entities follow the First-Come-First-Served (FCFS) queue discipline.
- Little's Law is intuitive and provides a useful tool for analyzing and understanding queuing systems. It essentially states that the average number of entities in the system is equal to the product of the average arrival rate and the average time each entity spends in the system.

**Review Questions**

- Q. 1 Estimate lengths, areas, and volumes by Monte Carlo methods
- Q. 2 Discuss level of significance of Type I and Type II errors (Z-test about a population mean). The number of concurrent users for some internet service provider has always averaged 5000 with a standard deviation of 800. After an equipment upgrade, the average number of users at 100 randomly selected moments of time is 5200. Does it indicate, at a 5% level of significance, that the mean number of concurrent users has increased? Assume that the standard deviation of the number of concurrent users has not changed.
- Q. 3 Compare two variances of F-tests.
- Q. 4 Stochastic process  $X(t)$  is Markov, under what two conditions future will be same as present and past. Define Poisson process. Poisson process is a suitable stochastic model in rare events. Justify?
- Q. 5 How does a queuing system work? What happens with a job when it goes through a queuing system?
- Q. 6 What is Queuing system and an illustration to the Little's Law with neat graph.



3

## **Basics of Data Visualization**

## Syllabus

**Syllabus** Computational Statistics and Data Visualization, Types of Data Visualization, Presentation and Exploratory Graphics, Graphics and Computing, Statistical Historiography, Scientific Design Choices in Data Visualization, Higher-dimensional Displays and Special Structures, Static Graphics : Complete Plots, Customization, Extensibility, Other Issues : 3-D Plots, Speed, Output Formats, Data Handling.

## **Other Issues : 3-D Plots, Speed, Output Formats, Data I/O**

3.1 Computational Statistics and Data Visualization

### 3.1.1 Types of Data Visualization

Data visualization comes in various forms, each suited to display specific types of data and convey different types of insights. Here are some common types of data visualization:

1. **Bar Charts** : Bar charts represent data using rectangular bars of varying lengths, with the length of each bar proportional to the value it represents. They are useful for comparing discrete categories or showing changes over time.
  2. **Line Charts** : Line charts use lines to connect data points, making them ideal for showing trends and patterns over continuous intervals or time.
  3. **Pie Charts** : Pie charts display data as slices of a circle, representing proportions or percentages of a whole. They are suitable for illustrating part-to-whole relationships.
  4. **Scatter Plots** : Scatter plots use individual data points to represent two variables, making them valuable for showing relationships, correlations, or clusters within the data.
  5. **Area Charts** : Area charts are similar to line charts but fill the area between the line and the x-axis, often used to show cumulative totals or stacked data.
  6. **Heatmaps** : Heatmaps use color intensity to represent the magnitude of values in a two-dimensional matrix. They are effective in displaying patterns or correlations in large datasets.
  7. **Bubble Charts** : Bubble charts are similar to scatter plots, but they use bubbles of different sizes to represent a third variable. They are useful for displaying three dimensions of data simultaneously.
  8. **Histograms** : Histograms group data into bins and display the frequency or distribution of values within each bin. They are useful for understanding the data's underlying distribution.
  9. **Box Plots (Box-and-Whisker Plots)** : Box plots provide a visual summary of the data's distribution, displaying the median, quartiles and outliers.
  10. **Choropleth Maps** : Choropleth maps use colors or patterns to represent data by geographic regions, making them ideal for showing spatial distributions or regional variations.

11. **Tree Maps** : Tree maps represent hierarchical data as nested rectangles, with the size of each rectangle proportional to the value it represents.
  12. **Network Graphs** : Network graphs display relationships between nodes and edges, helpful in visualizing connections, networks, or social interactions.
  13. **Word Clouds** : Word clouds visually represent the frequency of words in a text, with more frequently occurring words displayed in larger fonts.

3.1.2 Computational Statistics

Computational Statistics is a dynamic and interdisciplinary field that merges the principles of statistics, mathematics, and computer science to address complex and challenging problems in data analysis and statistical modeling. It emerged as a response to the exponential growth of data and the need to process and analyze vast datasets efficiently.

## **Key Aspects of Computational Statistics**

- Algorithm Development** : Computational Statistics involves the design and implementation of numerical algorithms to perform statistical computations, simulations, and optimizations. These algorithms are crafted to handle large-scale data, complex models, and computationally intensive tasks.
  - Simulation Techniques** : Monte Carlo simulations and other stochastic methods play a crucial role in Computational Statistics. They enable researchers to model and analyze intricate systems under uncertainty, providing valuable insights into real-world processes.
  - Bayesian Inference** : The use of computational methods is instrumental in Bayesian statistics. It allows practitioners to perform probabilistic inference and tackle complex Bayesian models, where traditional analytical solutions may not be feasible.
  - High-Performance Computing** : With the advent of big data, Computational Statistics leverages high-performance computing techniques, such as parallel processing and distributed computing, to process and analyze massive datasets efficiently.
  - Machine Learning Integration** : Computational Statistics often intersects with machine learning, incorporating advanced techniques to build predictive models, identify patterns, and gain deeper insights from data.
  - Resampling Methods** : Techniques like bootstrapping and cross-validation are widely employed in Computational Statistics to assess model performance, estimate uncertainties, and make inferences about population parameters.
  - Optimization Techniques** : Numerical optimization methods are used to estimate parameters in statistical models and find optimal solutions for various problems, including maximum likelihood estimation.

### **3.1.3 Applications of Computational Statistics**

- Data Science** : Computational Statistics is at the core of modern data science, enabling professionals to extract valuable information and patterns from vast and diverse datasets.
  - Bioinformatics** : In bioinformatics, Computational Statistics plays a pivotal role in analyzing genetic data, protein structures, and biological networks, facilitating discoveries in genomics and drug development.

## Data Modeling and Visualization

3. **Financial Modeling** : Computational Statistics is essential in modeling financial markets, risk assessment, and option pricing, aiding financial analysts and institutions in making informed decisions.
4. **Environmental Science** : Computational Statistics is used to analyze environmental data, climate models, and ecological systems, contributing to our understanding of environmental changes and conservation efforts.

In conclusion, Computational Statistics has become an indispensable tool in contemporary data-driven research and decision-making. Its integration of statistical methodologies and computational power empowers researchers, analysts, and data scientists to unravel complex problems and extract valuable insights from vast and intricate datasets. As technology continues to advance, Computational Statistics is expected to remain at the forefront of modern data analysis, making significant contributions across various domains.

**3.2 Presentation and Exploratory Graphics**

- Presentation graphics in visualization refer to the visual representations of data that are designed and optimized for communicating information to an audience. These graphics are specifically created for presentations, reports, and other forms of communication, aiming to convey insights clearly, concisely, and persuasively. The primary goal of presentation graphics is to enhance understanding and retention of information, making complex data more accessible to a broader audience.
- Characteristics and considerations of presentation graphics in visualization include :
  1. **Clarity and Simplicity** : Presentation graphics prioritize clarity and simplicity. They avoid unnecessary clutter and distractions, focusing on the key message and data patterns.
  2. **Visual Appeal** : Well-designed presentation graphics use colors, fonts, and layout to create visually appealing representations that capture the audience's attention.
  3. **Storytelling** : Presentation graphics are often part of a larger narrative or story. They help guide the audience through the data insights, supporting the message being conveyed.
  4. **Use of Annotations** : Annotations, labels, and captions are frequently employed to provide context and explanation for the visual elements, ensuring that the audience interprets the data accurately.
  5. **Interactivity (in some cases)** : While presentation graphics are often static, interactive elements may be incorporated to engage the audience and allow them to explore the data further.
  6. **Data Integrity** : Despite simplification, presentation graphics should maintain data accuracy and integrity, avoiding misleading or incorrect interpretations.
  7. **Audience Consideration** : Presentation graphics are tailored to the specific audience to ensure the information is comprehensible to them. This may involve adjusting the level of detail, using appropriate language, or focusing on relevant aspects.
- Common types of presentation graphics used in data visualization include :
  - o **Bar Charts and Column Charts** : Ideal for comparing data across categories.
  - o **Line Charts** : Effective in showing trends and changes over time.
  - o **Pie Charts** : Useful for illustrating proportions or percentages of a whole.

## Data Modeling and Visualization

- o **Area Charts** : Suitable for displaying cumulative data or stacked information.
- o **Scatter Plots** : Helpful in visualizing relationships and correlations between variables.
- o **Infographics** : Combine various visual elements to present complex information in an engaging manner.
- In conclusion, presentation graphics in visualization play a crucial role in transforming raw data into compelling and informative visual narratives. By employing design principles and focusing on clarity and simplicity, these graphics empower presenters to communicate data-driven insights effectively, fostering better understanding and decision-making by their audience.

**Note on Exploratory Graphics**

- Exploratory Graphics is an essential aspect of data visualization that aims to understand and gain insights from data through visual exploration. It involves the creation of graphical representations of data without a specific preconceived goal or hypothesis, allowing analysts and researchers to discover patterns, trends, and relationships that might not be immediately apparent from the raw data.
- In the initial stages of data analysis, exploratory graphics serve as a powerful tool for data exploration and hypothesis generation. By visualizing the data in various ways, analysts can identify potential outliers, spot interesting patterns, and formulate questions that lead to more in-depth investigations.
- The main features of Exploratory Graphics include flexibility and interactivity. Visualizations are often dynamic, enabling users to interact with the data, zoom in or out, filter, and change the visualization parameters on-the-fly. This interactivity empowers analysts to explore different aspects of the data and delve deeper into areas of interest.
- Exploratory graphics also emphasize the use of simple and intuitive visual representations to convey complex information effectively. Common types of exploratory graphics include scatter plots, histograms, box plots, density plots, and heatmaps.
- In summary, Exploratory Graphics is a crucial step in the data analysis process, providing a visual lens through which analysts can understand data patterns, generate hypotheses, and refine their understanding of the underlying data structure. By embracing creativity, curiosity, and interactivity, exploratory graphics support data-driven decision-making and lay the groundwork for more focused and rigorous data analysis techniques.

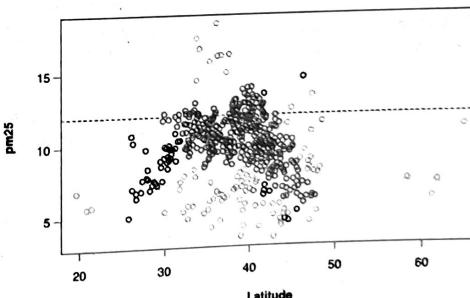


Fig. 3.2.1 : Lognormal Distribution

**Data Modeling and Visualization**

Load lesson dependencies

`library(swirl)``#install_from_swirl("Exploratory Data Analysis")``library(ggplot2)``library(jpeg)`

- The lognormal distribution graph indicates that the probability of body fat falling in the range of 20 to 24% is 0.1864. This information can help the researcher determine how many candidates they will need to assess to obtain a sufficient sample size.
- Random variables are a fundamental concept in statistics and probability theory. Discrete random variables take on countable, specific values, while continuous random variables assume uncountably infinite values. Understanding the properties of both types is crucial in many statistical applications.
- Exploratory Graphics is a fundamental aspect of data visualization that aims to understand and gain insights from data through visual exploration. It involves the creation of graphical representations of data without a specific preconceived goal or hypothesis, allowing analysts and researchers to discover patterns, trends, and relationships that might not be immediately apparent from the raw data.
- The primary objective of exploratory graphics is to provide an initial, visual overview of the data, enabling analysts to identify interesting features, outliers, clusters, and potential relationships between variables. By visualizing the data in various ways, analysts can formulate questions and hypotheses that lead to more in-depth investigations and refined data analysis.
- Characteristics of Exploratory Graphics in visualization include :
  - Flexibility and Interactivity :** Exploratory graphics often employ dynamic and interactive visualizations that allow users to manipulate the data representation in real-time. This interactivity empowers analysts to explore different aspects of the data and drill down into specific subsets or time periods for deeper insights.
  - Rapid Data Exploration :** With exploratory graphics, analysts can quickly visualize data from different angles and gain an intuitive understanding of its distribution, shape, and patterns. This agility is particularly valuable when dealing with large datasets or complex multivariate data.
  - Hypothesis Generation :** By visually examining the data without predefined expectations, exploratory graphics can lead to the generation of hypotheses and research directions that might have been overlooked otherwise.
  - Identifying Outliers and Anomalies :** Exploratory graphics can effectively highlight outliers and anomalies in the data, prompting further investigation into potential data errors or interesting phenomena.
  - Variable Relationships :** Through scatter plots and other visualizations, exploratory graphics help identify potential relationships and correlations between variables, which can guide the formulation of statistical models.
  - Common visualization techniques used in exploratory graphics include scatter plots, histograms, box plots, density plots, heatmaps, parallel coordinate plots, and interactive dashboard visualizations. The choice of visualization type depends on the data's nature and the insights the analyst aims to discover.

- In summary, Exploratory Graphics in visualization serves as a powerful and essential tool for data exploration and hypothesis generation. By fostering creativity, curiosity, and interactivity, exploratory graphics enable analysts to uncover valuable insights, lay the groundwork for further analysis, and make data-driven decisions with confidence.

**3.3 Graphics and Computing****3.3.1 Graphics in Data Visualization**

- Graphics play a central role in data visualization, acting as the visual bridge between raw data and human understanding. As a key component of data visualization, graphics are responsible for transforming complex datasets into intuitive and visually appealing representations that convey meaningful insights effectively. Graphics are essential in simplifying information, highlighting patterns, and supporting data-driven decision-making.
- The primary purpose of graphics in data visualization is to facilitate data exploration, analysis, and communication. They enable analysts, researchers, and decision-makers to grasp the underlying patterns, trends, and relationships within the data, making it easier to derive valuable insights and actionable conclusions.
- Graphics in data visualization offer various benefits :
  - Visual Perception :** Graphics leverage the human visual system's strengths to process and comprehend large amounts of data more efficiently than textual or numerical representations.
  - Communication :** Well-designed graphics enhance communication by conveying complex information in a concise and accessible manner, enabling both technical and non-technical audiences to understand the data.
  - Exploration and Discovery :** Graphics facilitate data exploration, allowing users to interact with the visualizations to uncover hidden patterns, outliers, and correlations.
  - Context and Comparison :** Graphics provide a clear context for data, enabling easy comparison between different data points, categories, or time periods.
  - Storytelling :** Graphics are instrumental in crafting data-driven narratives, presenting data insights as compelling stories that resonate with the audience.
- Data visualization employs a wide range of graphic types, each suited to represent specific types of data and insights. Common types of graphics include bar charts, line charts, scatter plots, pie charts, heatmaps, area charts, histograms, box plots, and more. The choice of graphic type depends on the data's nature, the analysis goals, and the target audience.
- However, while graphics are powerful tools for data visualization, their effectiveness relies on proper design and consideration of data accuracy. Poorly designed or misleading graphics can lead to misinterpretations and incorrect conclusions. Therefore, it is crucial to adhere to best practices in data visualization, ensuring that graphics are truthful, clear, and aesthetically pleasing.

**Data Modeling and Visualization**

- In conclusion, graphics are the backbone of data visualization, allowing data analysts and decision-makers to transform complex data into meaningful insights. By using graphics effectively, data visualization enhances understanding, improves decision-making, and unlocks the full potential of data as a valuable resource for solving real-world challenges.
- Graphic speed and power are crucial considerations in data visualization, particularly when dealing with large datasets or complex visualizations. Both aspects impact the user experience, the ability to explore data efficiently, and the overall effectiveness of the visualization in conveying insights.

**1. Speed**

- Graphic speed refers to how quickly a visualization can be rendered and updated, especially in response to user interactions or changes in the underlying data.
- Fast rendering is essential for creating responsive and interactive visualizations, allowing users to explore data in real-time and make quick comparisons or analyses.
- Slow graphic rendering can lead to frustration and hinder the exploration process, especially when dealing with dynamic or interactive visualizations.

**2. Power**

- Graphic power relates to the capacity of a visualization system to handle and display large datasets or complex visual representations.
- Powerful visualization systems can efficiently process and display vast amounts of data without sacrificing performance or visual quality.
- The ability to handle big data effectively enables analysts to explore intricate patterns and gain deeper insights from extensive datasets.

**Strategies to Improve Graphic Speed and Power :**

- Data Aggregation :** For large datasets, data aggregation can significantly improve speed and reduce computational overhead. Aggregating data points into appropriate bins or summary statistics before visualization reduces the number of elements rendered on the screen.
- Simplification Techniques :** In cases where high levels of detail are not essential, simplification techniques can be applied to reduce the complexity of visual elements without compromising the overall message.
- Progressive Loading :** For web-based or interactive visualizations, progressive loading can be implemented to show a preliminary representation quickly and gradually add more detail as the user interacts with the visualization.
- Use of GPU (Graphics Processing Unit) :** Leveraging the computational power of GPUs can enhance graphic speed, especially for computationally intensive visualizations like 3D or real-time simulations.
- Parallel Processing :** In scenarios where computations can be parallelized, using parallel processing techniques can expedite data manipulation and rendering, thereby improving overall speed.
- Data Filtering and Culling :** Implementing data filtering and culling techniques can remove unnecessary workload.

**Balancing Speed and Power**

Achieving a balance between speed and power is essential. While fast rendering is critical for smooth interactions, sacrificing visual quality or analytical depth may hinder meaningful data exploration. Optimal performance often depends on the specific use case, data size, and user requirements. Striking the right balance ensures that the visualization remains informative, responsive, and user-friendly.

- In conclusion, graphic speed and power are vital considerations in data visualization to deliver efficient, interactive, and informative visual representations. By employing strategies that improve rendering performance and handle large datasets effectively, data visualization tools can empower users to explore data rapidly, derive valuable insights, and make data-driven decisions with confidence.
- Covariance formula is a statistical formula, used to evaluate the relationship between two variables. It is one of the statistical measurements to know the relationship between the variance between the two variables. Let us say X and Y are any two variables, whose relationship has to be calculated. Thus the covariance of these two variables is denoted by  $\text{Cov}(X,Y)$ .

**3.3.2 Computing in Data Visualization**

- Computing plays a pivotal role in data visualization, enabling the creation of insightful and interactive visual representations of complex datasets. As data sets continue to grow in size and complexity, computational techniques are becoming increasingly essential to process, analyze, and render data efficiently.
- By leveraging the power of modern computing technologies, data visualization tools can handle large volumes of data, implement sophisticated algorithms, and provide users with dynamic and real-time visualizations.
- One of the primary tasks of computing in data visualization is data processing and preparation. Raw data is often messy, incomplete, or stored in various formats, making it challenging to visualize directly. Computational methods are employed to clean, transform, and preprocess the data, ensuring its readiness for visualization.
- This may involve data cleaning to handle missing values or outliers, data integration to merge disparate datasets, and data aggregation to summarize information for efficient visualization.
- Moreover, computing in data visualization facilitates the generation of various visual elements, such as charts, graphs, maps, and interactive components. For instance, when creating a scatter plot, the computer calculates the coordinates of each data point, determines appropriate scales for the axes, and renders the visualization based on the provided specifications. The computational process involves mapping data values to visual properties, such as color, size, or position, and coordinating the elements to create an informative and visually appealing display.
- Interactivity is a significant aspect of modern data visualization, allowing users to interact with visualizations and explore data from different angles. Computational techniques enable real-time responsiveness, enabling users to zoom, filter, or manipulate the visualization on the fly. For example, in a geographic map visualization, users can zoom in and out, pan across regions, and apply filters to focus on specific data subsets. These interactions require computational power to process user inputs and dynamically update the visualization accordingly.

Data Modeling and Visualization

- Furthermore, advanced computing techniques, such as machine learning and data mining algorithms, are increasingly integrated into data visualization tools. Machine learning models can be utilized to identify patterns, clusters, or anomalies in the data, providing valuable insights for visualization. Additionally, data mining techniques can be employed to discover meaningful associations and correlations, supporting the creation of informative visualizations that reveal hidden relationships within the data.
- High-performance computing is another critical aspect of data visualization, especially when dealing with massive datasets or computationally intensive visualizations. Parallel processing and distributed computing are employed to speed up data processing, rendering, and analysis, ensuring that visualizations remain responsive and scalable to handle big data.
- While computing advancements have significantly enhanced data visualization capabilities, it is essential to consider the potential challenges. As the complexity of visualizations increases, computational demands also grow, and resource constraints may arise. Careful optimization and utilization of computing resources are necessary to maintain fast rendering and responsiveness without compromising visual quality or accuracy.
- In conclusion, computing in data visualization is the backbone of modern visualization tools, enabling data processing, visual rendering, interactivity, and analysis. By harnessing the power of computational techniques, data visualization empowers users to explore and communicate complex data effectively, unlocking valuable insights and facilitating data-driven decision-making across various domains.

**3.3.3 Computational Power and Rendering in Visualization**

- Computational power and rendering are two critical components in data visualization that significantly impact the creation, interaction, and display of visual representations of data. As data volumes and complexity continue to grow, the role of computational power becomes increasingly essential in handling large datasets and implementing sophisticated visualization techniques. Meanwhile, rendering speed and efficiency are crucial for delivering responsive and interactive visualizations, enhancing the user experience and enabling real-time exploration of data.
- Computational power in data visualization refers to the computational resources, such as processing speed, memory, and parallelization capabilities, available to perform data manipulations, statistical calculations, and rendering tasks. With the advent of high-performance computing and modern GPUs (Graphics Processing Units), data visualization tools can leverage parallel processing to tackle computationally intensive tasks efficiently.
- This enhanced computational power enables handling vast datasets, supporting complex analytical operations, and generating high-quality visualizations even for real-time and interactive scenarios.
- Data visualization often involves preprocessing and transforming raw data into a suitable format for visualization. Computational power is instrumental in these data preparation steps, which may include cleaning the data, aggregating it, or performing feature engineering to extract meaningful insights.
- Additionally, advanced statistical algorithms and machine learning models can be integrated into data visualization tools to analyze data and uncover underlying patterns or trends. The computational power allows these algorithms to be applied effectively, enhancing the depth and accuracy of the visual insights presented to users.

Data Modeling and Visualization

- Furthermore, the rendering process is a critical aspect of data visualization, determining how visual elements are displayed on the screen. Rendering speed directly influences the responsiveness and interactivity of visualizations. High computational power allows for quick data rendering and dynamic updates, ensuring smooth interactions between users and visualizations. For instance, in interactive visualizations like zoomable maps or real-time charts, rendering speed is crucial in providing users with a seamless experience as they interact and explore the data.
- To achieve efficient rendering, optimization techniques are employed, such as data aggregation, caching, and level of detail (LOD) rendering. These techniques help reduce the computational load by prioritizing essential data points and simplifying the visual representation while maintaining the overall integrity of the visualization. For example, in a scatter plot with a large number of data points, aggregating or sampling the data can improve rendering speed without sacrificing data fidelity.
- Despite advancements in computational power, certain challenges remain in data visualization, particularly when dealing with extremely large datasets or complex visualizations. Balancing the need for high-performance rendering with the desire for comprehensive data exploration requires careful optimization and utilization of computational resources.
- In conclusion, computational power and rendering are critical components in data visualization, enabling the creation of sophisticated visualizations, real-time interactivity, and efficient handling of large datasets. By harnessing the capabilities of modern computing technologies, data visualization tools can empower users to gain deeper insights from their data, make data-driven decisions, and effectively communicate complex information in a visually engaging manner. As computing technology continues to evolve, data visualization is expected to reach new heights in performance, accuracy, and usability, advancing its applications across various fields and industries.
- Computation plays a fundamental role in data visualization, encompassing a range of processes and algorithms that transform raw data into meaningful visual representations. It involves data processing, statistical calculations, and rendering tasks, which collectively enable the creation of insightful and interactive visualizations that aid in understanding complex datasets. Computation enhances the capabilities of data visualization tools and empowers users to explore data from various perspectives, identify patterns, and make data-driven decisions effectively.

**3.3.4 Features of Computation in Data Visualization**

- Data Preprocessing :** Computation in data visualization involves data preprocessing to clean, transform, and prepare the raw data for visualization. This step ensures that the data is in a suitable format for analysis and visualization, addressing issues such as missing values, outliers, and data inconsistencies.
- Aggregation and Summarization :** Computation enables data aggregation and summarization to reduce data complexity and facilitate the visualization of large datasets. Aggregating data into meaningful groups or summary statistics allows users to gain insights from the overall trends without overwhelming visual clutter.
- Statistical Analysis :** Computation in data visualization encompasses various statistical calculations, such as mean, median, standard deviation, and correlation. These calculations provide valuable statistical insights that can be represented visually to aid in data exploration and understanding.

- Data Modeling and Visualization**

3-11

  - 4. Mapping Data to Visual Properties :** Computation is employed to map data values to visual properties such as position, size, color, and shape. This mapping allows for the creation of data-driven visual elements that represent the underlying data accurately.
  - 5. Interactivity :** Computation enables interactive data visualization, allowing users to interact with visualizations dynamically. Users can filter, zoom, pan, and explore data in real-time, gaining deeper insights through exploration and manipulation.

### 3.3.5 Applications of Computation in Data Visualization

- ### 3.3.5 Applications of Computation in Data Visualization

  - Business Intelligence**: Computation in data visualization is widely used in business intelligence applications to analyze and present business data in the form of dashboards, charts, and graphs. This facilitates data-driven decision-making and supports strategic planning.
  - Scientific Research**: Computation is essential in scientific research to visualize complex data from various disciplines, including biology, physics, chemistry, and environmental science. It helps researchers explore data patterns, identify trends, and communicate findings effectively.
  - Exploratory Data Analysis**: Computation aids in exploratory data analysis, allowing data scientists and analysts to interactively explore data, identify outliers, detect patterns, and formulate hypotheses.
  - Geospatial Visualization**: Computation is used in geospatial visualization to process and render geographic data, such as maps, satellite imagery, and spatial analyses. It enables users to visualize geographical patterns and relationships.
  - Data Journalism**: Computation plays a role in data journalism by supporting the creation of data-driven visualizations that enhance storytelling and communicate complex information to a broader audience.
  - Data-driven Art and Design**: Computation is utilized in data visualization for artistic and creative purposes, transforming data into aesthetically pleasing visual representations, such as generative art and data sculptures.

In summary, computation is the backbone of data visualization, providing the necessary tools and techniques to process, analyze, and represent data visually. Its features enable data exploration, pattern recognition, and interactivity, making it an indispensable tool for gaining insights from data and effectively communicating complex information. The applications of computation in data visualization are vast and span across numerous fields, driving data-driven decision-making, scientific research, and innovative data communication approaches.

Then,

### 3.4 Statistical Historiography

- Statistical historiography in data visualization refers to the study and exploration of the historical development and evolution of statistical graphics and visual representations of data. It delves into the origins of various visualization techniques, their transformations over time, and the influential figures who contributed to their advancements. This area of inquiry sheds light on the rich history of data visualization, illustrating how human ingenuity and technological progress have shaped the field, leading to the sophisticated visualizations we have today.

- The roots of statistical historiography in data visualization can be traced back to the early works of pioneers like William Playfair, who is often credited with the invention of the line graph, bar chart, and pie chart in the late 18th century. His groundbreaking visualizations laid the foundation for subsequent developments in data visualization and inspired other statisticians and scientists to explore the potential of graphical representations in conveying information effectively.
  - The 19th and early 20th centuries witnessed the contributions of eminent statisticians and data visualization practitioners such as Florence Nightingale, John Snow, and Charles Minard. Florence Nightingale's pioneering work in using coxcomb charts to illustrate mortality statistics during the Crimean War is a prominent example of the power of data visualization in influencing public policy and decision-making. John Snow's famous cholera map, depicting the distribution of cholera cases in London, provided a compelling visualization that helped identify the source of the outbreak.
  - The advent of computing in the mid-20th century revolutionized data visualization. The development of computer graphics and visualization tools enabled the creation of more complex and dynamic visualizations. The work of statisticians like John Tukey and Edward Tufte in the latter half of the 20th century further advanced the field, emphasizing the importance of effective visual communication and the need to avoid misleading graphical representations.
  - With the proliferation of the internet and advancements in technology, data visualization has become more accessible and widely used in various domains. Interactive and web-based visualizations allow users to engage with data, leading to the emergence of data journalism and data-driven storytelling. Moreover, the integration of data visualization with artificial intelligence and machine learning has opened new frontiers in data exploration and pattern recognition.
  - The historical context of data visualization helps modern practitioners understand the evolution of visualization techniques and provides insights into the choices and challenges faced by earlier statisticians. By studying historical visualizations, data visualization practitioners can learn from past successes and failures, inspiring the development of novel and impactful visualizations.
  - Statistical historiography in data visualization is an interdisciplinary endeavor, drawing on contributions from statisticians, computer scientists, designers, and domain experts. By examining the evolution of data visualization techniques over time, researchers gain a deeper appreciation of the transformative power of data visualization in scientific discovery, decision-making, and communication.
  - In conclusion, statistical historiography in data visualization is a valuable field of study that explores the historical context and development of visual representations of data. It highlights the contributions of visionaries who shaped the field and underscores the enduring significance of data visualization in interpreting complex information and facilitating data-driven insights. By understanding the historical roots of data visualization, practitioners can build on past achievements, continuously refining and innovating visualizations to meet the evolving needs of the data-driven world.
  - Statistical historiography in data visualization refers to the study and examination of the historical development and evolution of statistical graphics and visual representations of data. It involves analyzing the origins of various visualization techniques, tracing their transformations over time, and understanding the contributions of influential figures who have shaped the field. This area of research sheds light on the rich history of data visualization, revealing how human creativity and technological advancements have influenced the field, leading to the sophisticated visualizations we have today.

### 3.4.1 Features of Statistical Historiography

- 1. Historical Context :** Statistical historiography in data visualization provides a historical context for understanding the evolution of data visualization techniques. It explores the origins of key visual representations, such as bar charts, line graphs, scatter plots, and pie charts, dating back to the 18th and 19th centuries.
- 2. Contributions of Pioneers :** The study delves into the contributions of pioneering statisticians and data visualization practitioners who laid the groundwork for modern data visualization. Figures like William Playfair, Florence Nightingale, John Snow, and Charles Minard are among those whose innovative work has had a lasting impact on the field.
- 3. Evolution of Techniques :** Statistical historiography traces the evolution of visualization techniques from static and hand-drawn visualizations to interactive and computer-generated graphics. It highlights the role of technological advancements in shaping the complexity and interactivity of visualizations.
- 4. Influence on Decision-Making :** The historical study of data visualization reveals instances where visualizations have influenced public policy and decision-making. Visualizations like John Snow's cholera map and Florence Nightingale's coxcomb charts are notable examples of how data visualization has been a catalyst for change.
- 5. Impact on Communication :** Statistical historiography explores how visual communication has evolved over time, emphasizing the importance of effective data presentation and avoiding misleading graphical representations. The work of statisticians like John Tukey and Edward Tufte in promoting good visualization practices is a significant aspect of this study.

### 3.4.2 Applications of Statistical Historiography

- 1. Informing Best Practices :** By understanding the historical context and lessons learned from past visualizations, practitioners can develop better visualization practices. Insights gained from historical examples can inform the design and presentation of modern visualizations, ensuring they are accurate, informative, and effective.
- 2. Inspiring Innovation :** The examination of historical visualizations can serve as a source of inspiration for contemporary data visualization practitioners. Learning from the creative solutions and approaches of past visionaries can inspire new ideas and novel visualization techniques.
- 3. Supporting Research and Education :** Statistical historiography provides a valuable resource for researchers, educators, and students in the field of data visualization. Historical examples can be used to illustrate concepts and demonstrate the evolution of visualization techniques in educational settings.
- 4. Enhancing Data Journalism :** Understanding the history of data visualization can enrich data journalism practices. Journalists can draw from historical examples to create engaging and impactful data-driven storytelling and communication.
- 5. Advancing Interdisciplinary Collaboration :** The interdisciplinary nature of statistical historiography in data visualization encourages collaboration between statisticians, computer scientists, designers, and domain experts. By studying the historical development of visualizations, researchers from diverse fields can gain insights that foster innovation and cooperation.

In conclusion, statistical historiography in data visualization is an important field of study that examines the historical context and evolution of visual representations of data. It highlights the contributions of influential figures and explores the impact of data visualization on decision-making, communication, and scientific discovery. By drawing from historical examples, practitioners can develop better visualization practices and find inspiration for innovative visualizations that continue to push the boundaries of data exploration and understanding.

## 3.5 Scientific Design Choices in Data Visualization

### What are Design Choices in Data Visualization?

- Scientific design choices in data visualization refer to the deliberate and well-informed decisions made by data visualization practitioners to ensure the accuracy, clarity, and effectiveness of visual representations of data. As data visualization serves as a means of communication between data and its audience, making scientifically sound design choices is crucial to convey meaningful insights and support data-driven decision-making.
- One of the primary scientific design choices in data visualization involves selecting the appropriate type of visualization for the data at hand. Different types of data, such as categorical, numerical, or temporal, call for specific visualization techniques. For example, bar charts are ideal for comparing discrete categories, while line charts are well-suited for visualizing trends over time. Making the right choice ensures that the visualization effectively represents the underlying data structure and relationships.
- Another critical design choice is determining the visual encodings used to represent data attributes. Visual encodings map data values to visual properties, such as position, size, color, and shape. Scientifically choosing appropriate encodings is essential to avoid misleading interpretations and accurately convey quantitative information. For example, using a color scale that is perceptually uniform ensures that viewers can make reliable comparisons between different color-coded data points.
- Design choices related to labels, annotations, and titles are also essential in data visualization. Scientifically crafting informative labels and annotations adds context and clarity to the visualizations, guiding the audience in interpreting the data accurately. Providing titles and descriptions that convey the visualization's purpose and the data it represents aids in communicating the key message effectively.
- Color choices in data visualization demand careful consideration. Scientifically using colors that are visually distinguishable and have clear associations with data attributes helps viewers perceive patterns and relationships. Properly chosen color palettes aid in enhancing the visualization's aesthetics while avoiding visual clutter and confusion.
- Scale and axes design are fundamental elements in data visualization. Scientifically setting appropriate scales for axes ensures that data is presented in proportion, without distorting the information. Axes labels and ticks are chosen thoughtfully to facilitate accurate reading and interpretation of data values.
- Consideration of data density and data-ink ratio is an important scientific design choice. Eliminating non-essential elements, often referred to as chartjunk, reduces visual noise and enhances the clarity of the visualization. Maximizing the data-ink ratio, i.e., the proportion of ink dedicated to displaying data, helps to emphasize the crucial data elements and minimize distractions.

### 3.4.1 Features of Statistical Historiography

- Historical Context :** Statistical historiography in data visualization provides a historical context for understanding the evolution of data visualization techniques. It explores the origins of key visual representations, such as bar charts, line graphs, scatter plots, and pie charts, dating back to the 18th and 19th centuries.
- Contributions of Pioneers :** The study delves into the contributions of pioneering statisticians and data visualization practitioners who laid the groundwork for modern data visualization. Figures like William Playfair, Florence Nightingale, John Snow, and Charles Minard are among those whose innovative work has had a lasting impact on the field.
- Evolution of Techniques :** Statistical historiography traces the evolution of visualization techniques from static and hand-drawn visualizations to interactive and computer-generated graphics. It highlights the role of technological advancements in shaping the complexity and interactivity of visualizations.
- Influence on Decision-Making :** The historical study of data visualization reveals instances where visualizations have influenced public policy and decision-making. Visualizations like John Snow's cholera map and Florence Nightingale's coxcomb charts are notable examples of how data visualization has been a catalyst for change.
- Impact on Communication :** Statistical historiography explores how visual communication has evolved over time, emphasizing the importance of effective data presentation and avoiding misleading graphical representations. The work of statisticians like John Tukey and Edward Tufte in promoting good visualization practices is a significant aspect of this study.

### 3.4.2 Applications of Statistical Historiography

- Informing Best Practices :** By understanding the historical context and lessons learned from past visualizations, practitioners can develop better visualization practices. Insights gained from historical examples can inform the design and presentation of modern visualizations, ensuring they are accurate, informative, and effective.
- Inspiring Innovation :** The examination of historical visualizations can serve as a source of inspiration for contemporary data visualization practitioners. Learning from the creative solutions and approaches of past visionaries can inspire new ideas and novel visualization techniques.
- Supporting Research and Education :** Statistical historiography provides a valuable resource for researchers, educators, and students in the field of data visualization. Historical examples can be used to illustrate concepts and demonstrate the evolution of visualization techniques in educational settings.
- Enhancing Data Journalism :** Understanding the history of data visualization can enrich data journalism practices. Journalists can draw from historical examples to create engaging and impactful data-driven storytelling and communication.
- Advancing Interdisciplinary Collaboration :** The interdisciplinary nature of statistical historiography in data visualization encourages collaboration between statisticians, computer scientists, designers, and domain experts. By studying the historical development of visualizations, researchers from diverse fields can gain insights that foster innovation and cooperation.

In conclusion, statistical historiography in data visualization is an important field of study that examines the historical context and evolution of visual representations of data. It highlights the contributions of influential figures and explores the impact of data visualization on decision-making, communication, and scientific discovery. By drawing from historical examples, practitioners can develop better visualization practices and find inspiration for innovative visualizations that continue to push the boundaries of data exploration and understanding.

### 3.5 Scientific Design Choices in Data Visualization

#### What are Design Choices in Data Visualization?

- Scientific design choices in data visualization refer to the deliberate and well-informed decisions made by data visualization practitioners to ensure the accuracy, clarity, and effectiveness of visual representations of data. As data visualization serves as a means of communication between data and its audience, making scientifically sound design choices is crucial to convey meaningful insights and support data-driven decision-making.
- One of the primary scientific design choices in data visualization involves selecting the appropriate type of visualization for the data at hand. Different types of data, such as categorical, numerical, or temporal, call for specific visualization techniques. For example, bar charts are ideal for comparing discrete categories, while line charts are well-suited for visualizing trends over time. Making the right choice ensures that the visualization effectively represents the underlying data structure and relationships.
- Another critical design choice is determining the visual encodings used to represent data attributes. Visual encodings map data values to visual properties, such as position, size, color, and shape. Scientifically choosing appropriate encodings is essential to avoid misleading interpretations and accurately convey quantitative information. For example, using a color scale that is perceptually uniform ensures that viewers can make reliable comparisons between different color-coded data points.
- Design choices related to labels, annotations, and titles are also essential in data visualization. Scientifically crafting informative labels and annotations adds context and clarity to the visualizations, guiding the audience in interpreting the data accurately. Providing titles and descriptions that convey the visualization's purpose and the data it represents aids in communicating the key message effectively.
- Color choices in data visualization demand careful consideration. Scientifically using colors that are visually distinguishable and have clear associations with data attributes helps viewers perceive patterns and relationships. Properly chosen color palettes aid in enhancing the visualization's aesthetics while avoiding visual clutter and confusion.
- Scale and axes design are fundamental elements in data visualization. Scientifically setting appropriate scales for axes ensures that data is presented in proportion, without distorting the information. Axes labels and ticks are chosen thoughtfully to facilitate accurate reading and interpretation of data values.
- Consideration of data density and data-ink ratio is an important scientific design choice. Eliminating non-essential elements, often referred to as chartjunk, reduces visual noise and enhances the clarity of the visualization. Maximizing the data-ink ratio, i.e., the proportion of ink dedicated to displaying data, helps to emphasize the crucial data elements and minimize distractions.

- Basics of Data Visualization**
- Interactivity is another scientific design choice that can significantly enhance data visualization. Providing interactive elements, such as tooltips, zooming, or filtering, enables users to explore data in-depth and gain insights through active engagement with the visualizations. However, interactivity should be employed judiciously to avoid overwhelming the audience and maintain focus on the essential information.
  - In conclusion, scientific design choices in data visualization are essential for creating clear, accurate, and effective visual representations of data. By thoughtfully selecting appropriate visualization types, encodings, colors, labels, and interactivity, data visualization practitioners ensure that the visualizations align with the data's characteristics and effectively convey insights. Scientifically sound design choices enhance the communicative power of data visualization, making it a valuable tool for data exploration, analysis, and decision-making in various domains.

#### Which are the Design choices in Data visualization?

Design choices in data visualization encompass a wide range of decisions that data visualization practitioners make to create clear, informative, and visually appealing visual representations of data. These design choices include :

1. **Visualization Type** : Selecting the appropriate type of visualization based on the data's nature and the insights to be conveyed. Common types include bar charts, line charts, scatter plots, pie charts, heatmaps, and more.
2. **Visual Encodings** : Choosing visual properties (e.g., position, size, color, shape) to map data attributes accurately. Each data attribute should be represented using an appropriate visual encoding to facilitate understanding.
3. **Color Palette** : Deciding on a color palette that is visually distinguishable and meaningful for the data being visualized. Proper color choices aid in highlighting patterns and relationships within the data.
4. **Labels and Annotations** : Adding informative labels, titles, and annotations to provide context and help the audience understand the data and its implications.
5. **Scale and Axes** : Setting appropriate scales for the axes to ensure that the data is accurately represented and easily interpretable. Well-designed axes labels and ticks aid in data comprehension.
6. **Data Density and Simplification** : Considering the data density and employing techniques to simplify the visualization when dealing with large datasets to avoid visual clutter and improve readability.
7. **Chartjunk Removal** : Eliminating unnecessary graphical elements (chartjunk) that do not contribute to conveying the data effectively, improving the data-ink ratio.
8. **Data-Ink Ratio** : Maximizing the data-ink ratio by dedicating a higher proportion of graphical elements to display actual data, reducing distractions and enhancing the clarity of the visualization.
9. **Consistency and Formatting** : Ensuring consistency in the use of visual elements, fonts, and formatting across different components of the visualization.
10. **Interaction** : Incorporating interactive elements to create a cohesive and professional appearance.
11. **Accessibility** : Designing visualizations with consideration for accessibility, making them inclusive and usable for all users, including those with disabilities.

12. **Responsiveness** : Ensuring that the visualization adapts to different screen sizes and devices to maintain usability and readability across various platforms.
  13. **User-Centric Design** : Prioritizing the needs and expectations of the audience, keeping the target users in mind when making design choices to optimize the user experience.
  14. **Storytelling** : Incorporating elements that enable data-driven storytelling, using visualizations to convey a narrative and engage the audience effectively.
  15. **Aesthetics** : Focusing on creating aesthetically pleasing visualizations that attract and engage the audience, enhancing the overall impact of the data presentation.
- By carefully considering and implementing these design choices, data visualization practitioners can create meaningful and impactful visualizations that effectively communicate insights and support data-driven decision-making.

### 3.6 Higher-Dimensional Displays and Special Structures

- Higher-dimensional displays in data visualization refer to techniques and methods designed to represent and visualize data with more than three dimensions. While traditional 2D and 3D visualizations are well-established and widely used, higher-dimensional data poses unique challenges due to human perceptual limitations. Higher-dimensional displays address this issue by employing various strategies to project and present multi-dimensional data in ways that are comprehensible and insightful.
- One common approach in higher-dimensional displays is dimensionality reduction, where complex data is transformed into a lower-dimensional space while preserving essential characteristics. Techniques such as Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE) are used to project high-dimensional data into 2D or 3D space, allowing for visualization and exploration. Dimensionality reduction is particularly useful when dealing with datasets with numerous variables or attributes, enabling analysts to gain an overview of the data's structure and patterns.
- Parallel coordinate plots are another powerful tool for higher-dimensional displays. In this visualization, each data point is represented as a polyline that connects coordinates along different axes, with each axis corresponding to a different variable. This method enables the simultaneous visualization of multiple dimensions and can reveal relationships and patterns that might be challenging to discern in lower-dimensional views.
- Scatterplot matrices are also effective for displaying higher-dimensional data. In scatterplot matrices, multiple scatter plots are arranged in a grid, with each plot representing the relationship between a pair of variables. This approach facilitates the comparison of all variable combinations, providing a comprehensive view of the data's multivariate distribution and correlation structure.
- Another technique for higher-dimensional displays is data aggregation or binning. In this approach, high-dimensional data is grouped or summarized into lower-dimensional arrays or distributions. Aggregating data reduces the complexity of visualizations and can help identify trends or patterns in the data that might not be evident in the raw multi-dimensional representation.

- Data Modeling and Visualization**
- Glyph-based visualizations are commonly used in higher-dimensional displays to encode multiple dimensions using graphical symbols or glyphs. For example, color, size, shape, and texture of glyphs can be utilized to represent different variables, providing an intuitive and compact representation of multi-dimensional data.
  - Advances in interactive visualization tools have greatly enhanced higher-dimensional displays. Interactive techniques, such as brushing and linking, enable users to select subsets of data in one view and see corresponding data highlighted in other views. This interactivity allows analysts to drill down into specific dimensions or subsets of the data, gaining deeper insights from complex multi-dimensional datasets.
  - Machine learning algorithms, such as clustering and classification, can also aid in higher-dimensional displays by automatically organizing and labeling data based on patterns in the high-dimensional space. This can help reveal inherent groupings and structures within the data, supporting data exploration and understanding.
  - Despite the utility of higher-dimensional displays, they come with challenges related to visual clutter and the curse of dimensionality. Visualizing data with many dimensions can result in overcrowded visualizations and make it difficult to interpret individual data points.
  - As the number of dimensions increases, the data becomes sparser in the visual space, leading to potential loss of information and patterns. Addressing these challenges requires careful consideration of appropriate visualization techniques, interactivity, and the use of dimensionality reduction methods to aid in meaningful data exploration.
  - In conclusion, higher-dimensional displays in data visualization are vital for understanding complex datasets with multiple attributes. By employing techniques such as dimensionality reduction, parallel coordinate plots, scatterplot matrices, and data aggregation, analysts can gain valuable insights from multi-dimensional data.
  - Interactivity, glyph-based visualizations, and machine learning algorithms further enhance the capabilities of higher-dimensional displays, providing powerful tools for data exploration, pattern recognition, and knowledge discovery across various domains. However, it is essential to address challenges related to visual clutter and the curse of dimensionality to create informative and visually appealing higher-dimensional visualizations.

#### Examples of Higher-dimensional Displays in Data Visualization

##### 1. t-SNE Visualization

t-SNE is a dimensionality reduction technique commonly used for visualizing high-dimensional data in 2D or 3D space. Let's consider an example using the popular Iris dataset, which has four dimensions (sepal length, sepal width, petal length, and petal width) for each sample :

```
'''python
```

```
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.manifold import TSNE

# Load the Iris dataset
data = load_iris()
```

#### Data Modeling and Visualization

```
X, y = data.data, data.target
```

```
# Perform t-SNE to reduce data to 2D
```

```
tsne = TSNE(n_components=2, random_state=42)
```

```
X_tsne = tsne.fit_transform(X)
```

```
# Plot the 2D t-SNE visualization
```

```
plt.scatter(X_tsne[:, 0], X_tsne[:, 1], c=y, cmap='viridis')
```

```
plt.colorbar()
```

```
plt.xlabel('t-SNE Component 1')
```

```
plt.ylabel('t-SNE Component 2')
```

```
plt.title('t-SNE Visualization of Iris Dataset')
```

```
plt.show()
```

```
'''
```

##### 2. Parallel Coordinate Plot

Parallel coordinate plots are used to visualize multi-dimensional data. Let's consider an example using the famous "wine" dataset, which has 13 dimensions (features) for each wine sample :

```
'''python
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
from pandas.plotting import parallel_coordinates
```

```
# Load the Wine dataset
```

```
url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data'
```

```
df = pd.read_csv(url, header=None)
```

```
df.columns = ['Class', 'Alcohol', 'Malic Acid', 'Ash', 'Alcalinity', 'Magnesium', 'Total Phenols',
```

```
'Flavanoids', 'Nonflavanoid Phenols', 'Proanthocyanins', 'Color Intensity', 'Hue',
```

```
'OD280/OD315', 'Proline']
```

```
# Normalize the data for parallel coordinate plot
```

```
df_normalized = (df - df.min()) / (df.max() - df.min())
```

```
# Plot parallel coordinate plot
```

```
plt.figure(figsize=(10, 6))
```

```
parallel_coordinates(df_normalized, 'Class', colormap='viridis')
```

```
plt.title('Parallel Coordinate Plot of Wine Dataset')
```

```
plt.show()
```

```
'''
```

**Data Modeling and Visualization****3. Glyph-based Visualization**

Glyph-based visualizations use graphical symbols to represent multiple dimensions. Let's consider an example using the "Seaborn" library to create a scatter plot matrix, where each scatter plot shows the relationship between two attributes of the Iris dataset:

```
'''python
Import seaborn as sns
from sklearn.datasets import load_iris
```

```
# Load the Iris dataset
data = load_iris()
df = pd.DataFrame(data.data, columns=data.feature_names)
df['target'] = data.target
```

```
# Create a scatter plot matrix (Pairplot)
sns.pairplot(df, hue='target', palette='viridis')
plt.suptitle('Scatter Plot Matrix of Iris Dataset', y=1.02)
plt.show()
'''
```

Formulas and codes for all possible higher-dimensional displays are beyond the scope of a single response. The examples provided demonstrate some commonly used higher-dimensional display techniques. There are many other visualization libraries and techniques available, such as multidimensional scaling (MDS), parallel sets, and hyperbolic trees, among others, which can be explored to visualize higher-dimensional data effectively. Each technique may require different libraries, algorithms, and data preprocessing steps, depending on the specific visualization task and data characteristics.

**3.6.1 Special Structure in Data Visualization**

- Special structures in data visualization refer to unique patterns, configurations, or relationships within the data that require specialized visualization techniques to effectively represent and understand. These structures often arise in specific types of data or domains, and standard visualization methods may not be sufficient to capture their complexity or reveal their underlying insights.
- By employing specialized visualization techniques, data visualization practitioners can uncover hidden patterns and gain deeper insights into the data, leading to better decision-making and understanding within their respective fields.
- One example of special structures is time series data. Time series data represents observations collected over a sequence of time intervals, such as daily stock prices, monthly temperature records, or hourly website traffic. Time series visualizations often involve line charts or area plots to depict trends and seasonal patterns over time. Techniques like seasonality decomposition, where the data is separated into trend, seasonal, and residual components, allow for better understanding and forecasting of temporal patterns.

- Network data is another type of data with special structures that require dedicated visualization techniques. Network visualizations represent relationships between entities as nodes (vertices) connected by edges (links). Social networks, transportation networks, and biological networks are examples of domains where network data arises. Visualization methods like node-link diagrams, matrix plots, and force-directed layouts are commonly used to reveal the structure and connections within networks.
- Geospatial data also has special structures that demand specific visualization approaches. Geospatial visualizations display data on maps, allowing analysts to understand spatial patterns and correlations. Choropleth maps, heatmaps, and point-based maps are some common techniques used to represent geospatial data. Geospatial visualizations find applications in fields such as urban planning, environmental monitoring, and epidemiology.
- Data with hierarchical structures presents another challenge for visualization. Hierarchical data consists of nested categories or levels, like organizational hierarchies or file directory structures. Treemaps, sunburst charts, and dendograms are visualization techniques used to represent hierarchical data, enabling users to explore data at different levels of granularity.
- Parallel coordinates plots are employed for visualizing multivariate data, where each axis represents a different variable. Parallel coordinates are especially useful when dealing with high-dimensional datasets, as they allow for the simultaneous visualization of multiple dimensions and facilitate the identification of patterns and relationships.
- Additionally, data with spatial or temporal uncertainty requires specialized visualization techniques to represent the uncertainty intervals. Uncertainty visualization methods like error bars, confidence bands, and probabilistic graphical models help to communicate the reliability and variability of data points, aiding in decision-making and risk assessment.
- Another important aspect of special structures in data visualization is the use of interactive techniques. Interactive visualization allows users to dynamically explore the data, filter and drill down into specific segments, and gain insights from different perspectives. Techniques like brushing and linking, zooming, and filtering are often applied to enhance user interactions with complex datasets.
- In summary, special structures in data visualization encompass a diverse range of patterns and relationships that require dedicated visualization techniques to reveal their unique insights. Time series data, network data, geospatial data, hierarchical data, and multivariate data are some examples of data types with special structures.
- Specialized visualization methods enable analysts to gain deeper insights, uncover hidden patterns, and support decision-making across various domains. Interactive techniques play a vital role in exploring complex data, empowering users to engage with the data visualization and extract meaningful information effectively. As data continues to grow in complexity and diversity, the development of specialized visualization techniques will remain crucial in facilitating data exploration and understanding.
- Probability Plotting :** A method of finding parameter values where the data is plotted on special plotting paper and parameters are derived from the visual plot.

### 3.7 Static Graphics

- Static graphics in data visualization refer to visual representations of data that are fixed and non-interactive. Unlike dynamic or interactive visualizations, static graphics are pre-designed and do not allow users to manipulate or explore the data in real-time. While dynamic visualizations offer greater interactivity and flexibility, static graphics play a critical role in data communication, presentation, and dissemination.
- One of the primary advantages of static graphics is their simplicity and ease of understanding. Since these visualizations are pre-designed and fixed, they can be carefully crafted to highlight specific patterns, trends, or insights in the data. Static graphics are particularly useful for communicating a clear and concise message to a broad audience, making them ideal for data-driven reports, infographics, and academic publications.
- Static graphics are often preferred when the data story is straightforward, and there is no need for user interaction or exploration. For instance, when presenting simple comparisons, distributions, or trends, a static bar chart, line chart, or pie chart can effectively convey the information without overwhelming the audience with unnecessary complexity.
- Another advantage of static graphics is their reproducibility and consistency. Once a static visualization is designed, it can be easily replicated and reused in various contexts without changing the original message. This reproducibility ensures that the data insights are consistently communicated to different audiences or across different platforms.
- Static graphics are commonly used for data visualization in printed materials, such as books, magazines, and reports, where interactivity is not possible. They are also used in presentations and slideshows to convey data insights to an audience during conferences or meetings.
- Designing static graphics requires careful attention to aesthetics and visual clarity. Since there is no opportunity for user interaction, the visual elements and annotations must be well-crafted to facilitate easy comprehension and understanding. Effective use of color, size, and layout is essential to ensure that the data points are distinct and easily identifiable.
- In contrast to dynamic visualizations, which may require more computational resources and processing time, static graphics are computationally efficient. They can be generated and rendered quickly, making them suitable for large-scale data sets or real-time data updates.
- However, static graphics have their limitations. They may not be suitable for complex or high-dimensional data sets, as the fixed nature of the visualization may not capture all the intricacies and patterns in the data. In such cases, interactive or dynamic visualizations can offer more in-depth exploration and analysis opportunities.
- In conclusion, static graphics in data visualization play a crucial role in presenting data insights in a concise and visually appealing manner. Their simplicity, reproducibility, and ability to communicate a clear message make them valuable tools for data communication and reporting. While static graphics may not offer the interactivity and exploration capabilities of dynamic visualizations, they remain a fundamental component of data visualization, catering to specific communication needs and contexts where user interaction is not required.

### 3.7.1 Complete Plots in Data Visualization

- Complete plots in data visualization refer to visual representations that encompass all necessary elements to present data comprehensively and convey meaningful insights effectively. These plots are self-contained and provide a holistic view of the data, enabling users to grasp key patterns, relationships, and trends at a glance. Complete plots are often used in exploratory data analysis, presentations, and reports where a comprehensive understanding of the data is essential.
- One of the most common examples of a complete plot is the box plot (box-and-whisker plot). It displays the distribution of a dataset's values, including the median, quartiles, and any outliers, providing a clear overview of the data's spread and central tendency in a single visualization.
- Scatter plots with regression lines are another example of complete plots. They not only show the relationship between two variables but also include a regression line that indicates the overall trend and direction of the relationship.
- Complete plots often include informative labels, titles, and legends to provide context and aid in data interpretation. Clear axis labels and units are essential to understanding the scale and measurement of the data.
- Another feature of complete plots is the use of appropriate color schemes and visual elements to highlight key data points or distinguish different categories or groups. Color can be used to represent different data classes or highlight specific data segments.
- Annotations are commonly used in complete plots to provide additional information or insights. Annotations might include callouts for particular data points or text boxes to explain notable observations.
- Complete plots often have a balanced combination of simplicity and visual complexity. They are designed to convey the essential information without overwhelming the audience with excessive visual elements or unnecessary complexity.
- In summary, complete plots in data visualization are visual representations that encompass all necessary elements to present data comprehensively and effectively. They include features such as box plots, scatter plots with regression lines, informative labels and legends, appropriate color schemes, annotations, and a balance between simplicity and visual complexity. Complete plots are valuable tools in exploratory data analysis, presentations, and reports, as they provide a comprehensive understanding of the data at a glance, enabling better decision-making and data-driven insights.

#### Implementation of Complete Plots

The implementation of complete plots in data visualization can be achieved using various programming languages and libraries that support data visualization. I will demonstrate the implementation of two complete plots : the box plot and the scatter plot with a regression line using Python and the popular data visualization library, Matplotlib.

##### 1. Box Plot

The following Python code demonstrates how to create a box plot using Matplotlib :

```
"python
import matplotlib.pyplot as plt"
```

Data Modeling and Visualization

import numpy as np

```
# Sample data for the box plot
data = [np.random.normal(0, 1, 100), np.random.normal(2, 1, 100), np.random.normal(-2, 1, 100)]
```

```
# Create the box plot
plt.boxplot(data, labels=['Group 1', 'Group 2', 'Group 3'])
plt.xlabel('Groups')
plt.ylabel('Values')
plt.title('Box Plot Example')
plt.grid(True)
plt.show()
```

...  
In this example, we generate three random datasets and create a box plot to visualize their distributions. The box plot shows the quartiles, median, and potential outliers for each group.

**2. Scatter Plot with Regression Line**

The following Python code demonstrates how to create a scatter plot with a regression line using Matplotlib:

```
'''python
import matplotlib.pyplot as plt
import numpy as np
# Sample data for the scatter plot
x = np.random.rand(100)
y = 2 * x + np.random.normal(0, 0.2, 100)

# Create the scatter plot
plt.scatter(x, y, label='Data Points', color='b')
plt.plot(x, 2 * x, label='Regression Line', color='r')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.title('Scatter Plot with Regression Line')
plt.legend()
plt.grid(True)
plt.show()
```

...  
In this example, we generate random data points for the x and y variables and create a scatter plot to visualize their relationship. We also add a red regression line that represents the linear relationship between the two variables.

These examples illustrate the implementation of complete plots in data visualization using Matplotlib in Python. Depending on the specific data and requirements, other data visualization libraries like Seaborn, Plotly, or ggplot2 in R can also be used to create complete plots with similar functionalities. The choice of library may vary based on the programming language, preferred visual aesthetics, and additional features required for the visualization.

**3.7.2 Customization**

Customization in data visualization refers to the ability to modify and personalize various aspects of a visual representation to suit specific needs, preferences, or requirements. Customization empowers data visualization practitioners to tailor the appearance and functionality of visualizations, making them more effective in conveying insights and supporting data-driven decision-making. There are several aspects of a visualization that can be customized:

- 1. Visual Elements :** Practitioners can customize the visual elements of a data visualization, such as colors, shapes, sizes, and styles, to make the visualization visually appealing and suitable for the target audience. Color palettes can be chosen to represent specific data categories or highlight patterns effectively. Customizing shapes and sizes can help distinguish data points or emphasize important elements.
- 2. Axes and Scales :** Customizing axes and scales allows for better data representation and interpretation. Practitioners can set specific axis ranges, intervals, and tick marks to ensure the data is presented accurately and meaningfully. Logarithmic scales, for instance, can be applied to handle data with a wide range of values.
- 3. Labels and Annotations :** Customizing labels and annotations is essential for providing context and explaining the data visualization. Practitioners can add informative titles, axis labels, and data point labels to aid interpretation. Annotations can be used to highlight specific data points or events of interest.
- 4. Layout and Composition :** Customizing the layout and composition of a visualization involves arranging multiple visual elements to create a coherent and informative presentation. Practitioners can adjust the positioning, alignment, and spacing of different components to enhance the overall aesthetics and readability.
- 5. Interactivity :** Customizing interactivity enables users to interact with the visualization, allowing them to explore and analyze the data from different perspectives. Interactive elements such as tooltips, zooming, panning, and filtering can be added to engage users and facilitate data exploration.
- 6. Theme and Style :** Customizing the overall theme and style of a visualization can align it with branding guidelines or the visual identity of a project. Consistent themes and styles create a cohesive look and feel across multiple visualizations, enhancing their presentation and professionalism.
- 7. Data Aggregation :** Customizing data aggregation methods can be useful when dealing with large datasets. Practitioners can choose appropriate aggregation techniques like binning, averaging, or summing to simplify the visualization and reduce visual clutter.
- 8. Animation :** Customizing animation effects can be applied to demonstrate changes over time or to reveal specific data patterns dynamically. Animated visualizations are engaging and can effectively convey temporal trends or data transitions.

Customization in data visualization allows practitioners to cater to specific audiences, domain requirements, and communication goals. It fosters creativity and flexibility in data presentation, enabling practitioners to create visualizations that are unique, impactful, and tailored to the specific characteristics of the data and the intended audience. By leveraging customization options, data visualization becomes a powerful tool for effective data communication, storytelling, and data-driven insights.

### 3.7.3 Extensibility

Extensibility in data visualization and modeling refers to the capability of data visualization tools and modeling frameworks to be easily extended, customized, and adapted to handle new data types, visualization types, algorithms, or functionalities. It enables data scientists, researchers, and developers to build on existing tools and models, adding new features or modifying existing ones to meet specific needs and address unique challenges. Extensibility plays a crucial role in empowering users to tailor data visualization and modeling solutions to their specific requirements, making them more versatile, efficient, and effective.

#### In Data Visualization

- Custom Visualizations** : Extensible data visualization tools allow users to create custom visualizations beyond the standard chart types. Users can develop new visualization types that suit the specific characteristics of their data and the insights they want to communicate.
- Integration of External Libraries** : Extensibility enables seamless integration with external libraries and APIs, empowering users to leverage additional functionalities and visual elements available in other tools or frameworks.
- Data Connectors** : Extensible visualization tools can be extended to support connections with various data sources, databases, or APIs, allowing users to easily access and visualize data from multiple platforms.
- Interactivity** : Extensible visualization tools often provide APIs and libraries that enable developers to add interactive features, such as tooltips, zooming, or filtering, to enhance user engagement and data exploration.
- Theming and Styling** : Extensibility allows users to define and apply custom themes and styles to data visualizations, aligning them with branding guidelines or the visual identity of a project.

#### In Data Modeling

- Custom Algorithms** : Extensible modeling frameworks allow data scientists to implement and integrate custom algorithms and models to address specific data analysis or machine learning tasks.
- Feature Engineering** : Extensibility enables users to extend feature extraction and engineering capabilities by adding new data preprocessing techniques or transforming existing features to better suit the modeling requirements.
- Incorporating External Libraries** : Extensible modeling frameworks facilitate the integration of external libraries and modules, enabling users to leverage state-of-the-art machine learning algorithms or domain-specific tools.
- Hyperparameter Tuning** : Extensibility allows for custom hyperparameter tuning algorithms or strategies to optimize model performance effectively.

**5. Support for Domain-specific Data Types** : Some domains may require specialized data types not commonly supported in standard modeling frameworks. Extensibility allows users to handle such data types, ensuring the modeling framework is more inclusive and applicable across diverse fields.

The concept of extensibility is essential in the rapidly evolving fields of data visualization and modeling. It fosters collaboration and knowledge sharing, as researchers and developers can contribute new capabilities and improvements to existing tools and frameworks.

Additionally, extensibility encourages the development of open-source projects and community-driven initiatives, making advanced data visualization and modeling techniques accessible to a wider audience. By embracing extensibility, data visualization tools, and modeling frameworks become more versatile and adaptable, enabling users to push the boundaries of data analysis and derive more meaningful insights from complex datasets.

## 3.8 Other Issues

### 3.8.1 What are 3D-Plots ?

3D plots in data visualization and modeling refer to visual representations and modeling techniques that involve three-dimensional data points and relationships. These plots are used when data contains three dimensions (variables) and require visualization or analysis in a 3D space. 3D plots are particularly useful when dealing with data that has spatial, temporal, or volumetric attributes, as they allow for a more comprehensive understanding of the data's structure and patterns.

#### In Data Visualization

- 3D Scatter Plot** : A 3D scatter plot represents data points in a three-dimensional space, with each point defined by its three attributes. This type of plot is useful for visualizing the distribution and relationships between three variables.
- 3D Surface Plot** : A 3D surface plot represents data as a continuous surface over a 3D space. It is commonly used to visualize functions or datasets with continuous variables, showing how the data changes across the three dimensions.
- 3D Line Plot** : A 3D line plot shows lines connecting data points in a 3D space. It is often used to visualize time series data or spatial trajectories with three dimensions.
- 3D Bar Plot** : A 3D bar plot represents data using bars or columns in a 3D space. It is suitable for comparing data across three dimensions, similar to a 2D bar chart but with an additional dimension.

#### In Data Modeling

- 3D Feature Space** : In machine learning, 3D feature space refers to the space formed by three feature variables or attributes. Visualizing data in a 3D feature space can provide insights into the relationships between features and their impact on the target variable.
- 3D Decision Boundaries** : In classification tasks, 3D decision boundaries are used to separate different classes or categories based on the values of three features. Visualizing decision boundaries in a 3D space can help understand the model's classification performance.

**Data Modeling and Visualization**

- 3. 3D Model Evaluation :** In model evaluation, 3D plots can be used to visualize model performance across three dimensions, such as hyperparameter tuning. This visualization helps identify optimal model settings and understand the model's sensitivity to different parameters.

It is essential to note that while 3D plots can provide valuable insights, they also have limitations. Visualizing data in three dimensions can be challenging for human perception, as our eyes are more accustomed to 2D representations. 3D plots can lead to visual clutter and distortions if not carefully designed. Additionally, 3D plots may not be suitable for very high-dimensional data, as it becomes challenging to interpret and comprehend visualizations in higher-dimensional spaces.

In summary, 3D plots in data visualization and modeling provide a valuable tool for visualizing and analyzing three-dimensional data. They enable users to explore spatial, temporal, or volumetric relationships and gain a deeper understanding of the data's characteristics. However, it is crucial to use 3D plots judiciously and consider the complexity and interpretability of the visualizations, particularly when dealing with high-dimensional data.

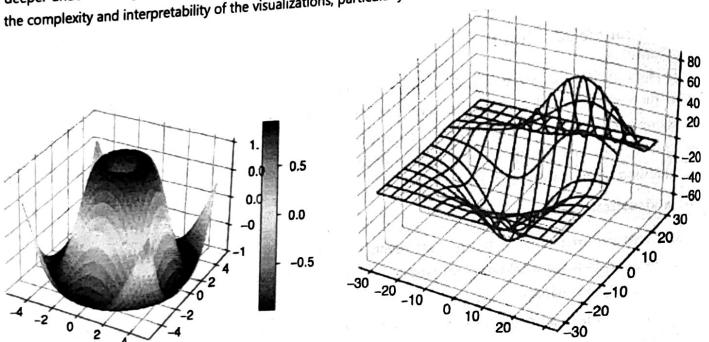


Fig. 3.8.1 : 3D Model Evaluation

### 3.8.2 Output Formats

Output formats in data visualization and modeling refer to the various ways in which visualizations and modeling results can be presented, shared, or saved for different purposes and audiences. The choice of output format depends on factors such as the complexity of the data, the target audience, the level of interactivity required, and the intended use of the visualizations or model results. There are several common output formats used in data visualization and modeling:

- 1. Static Images :** Static images are one of the most straightforward output formats for data visualizations. They are commonly used for presenting visualizations in reports, presentations, and publications. Formats like PNG, JPEG, or SVG (Scalable Vector Graphics) are commonly used for saving static images of visualizations. These formats ensure that the visualizations are easily shareable and can be embedded in various documents or web pages.

- 2. Interactive Web-based Visualizations :** Web-based visualizations offer a high degree of interactivity, allowing users to explore and interact with the data in real-time. JavaScript libraries and frameworks like D3.js, Plotly, and Bokeh are often used to create interactive visualizations that can be embedded in websites or web applications. Interactive visualizations are valuable for data exploration, as they enable users to interactively filter, zoom, and explore the data.
- 3. Interactive Dashboards :** Dashboards provide a collection of interactive visualizations and insights, offering a comprehensive view of the data. They are commonly used for real-time monitoring, business intelligence, and decision-making. Tools like Tableau, Power BI, and Grafana enable users to create interactive dashboards that can be accessed and shared with stakeholders.
- 4. Animated Visualizations :** Animated visualizations are output formats that use animation to show changes over time or data transitions. They are particularly useful for visualizing time series data, simulations, and dynamic processes. Animated GIFs or video formats are often used to save and share animated visualizations.
- 5. Data Files :** In data modeling, the output format can include data files containing model predictions, feature importance scores, or other model outputs. Common file formats like CSV (Comma-Separated Values) or JSON (JavaScript Object Notation) are used to save data for further analysis or integration with other tools.
- 6. Model Serialization :** In machine learning, model serialization is used to save trained models for later use or deployment. Serialized models can be saved in formats like Pickle (Python-specific), ONNX (Open Neural Network Exchange), or PMML (Predictive Model Markup Language) to ensure compatibility with different environments.
- 7. 3D Printing :** In certain domains like engineering and medical visualization, 3D printing can be used as an output format to create physical representations of 3D models or data visualizations. This allows for a tangible and hands-on exploration of complex data structures.

The choice of output format depends on the specific requirements of the data visualization or modeling task. Static images are suitable for simple visualizations in reports, while interactive web-based visualizations and dashboards are ideal for dynamic exploration. Animated visualizations are effective for showing temporal patterns, and data files and model serialization are essential for sharing data and model results with other tools or platforms. By considering the appropriate output format, data visualization and modeling outputs can be effectively communicated and utilized for decision-making, research, and data-driven insights.

### 3.8.3 Data Handling

- Data handling in data visualization is a crucial process that involves managing, processing, and preparing the data before visualizing it. It is the foundation upon which successful data visualizations are built, as the quality, accuracy, and structure of the data directly impact the effectiveness and interpretability of the visual representations. Proper data handling ensures that the visualizations are meaningful, insightful, and aligned with the objectives of the data analysis.
- The first step in data handling is data acquisition. This involves obtaining the data from various sources, such as databases, files, APIs, or web scraping. The data may come in different formats, including CSV, Excel, JSON, or relational databases. Data acquisition requires careful consideration of the data's structure and content to ensure that it can be effectively processed for visualization.

- After data acquisition, the next step is data cleaning. Data cleaning involves identifying and rectifying errors, missing values, duplicates, and inconsistencies in the dataset. Cleaning the data is essential for ensuring that the visualizations accurately represent the underlying information and do not mislead the audience. Data cleaning may involve imputing missing values, removing outliers, or resolving inconsistencies in data entries.
- Data transformation is another critical aspect of data handling. Data may need to be transformed or reshaped to fit the requirements of the chosen visualization technique. This could involve aggregating data, converting data types, or normalizing values. Data transformation ensures that the data is in a suitable format for visualization, enabling the representation of patterns and relationships effectively.
- For many visualization tasks, data may come from multiple sources or be stored in separate datasets. Data integration is the process of combining relevant information from different datasets to create a unified dataset for visualization. Integration may involve merging, joining, or appending datasets based on common identifiers.
- Handling categorical data is another aspect of data handling. Categorical data needs to be converted into numeric values for visualization purposes. Common techniques for handling categorical data include one-hot encoding and label encoding, which transform categorical variables into numerical representations.
- Time series data requires special handling to account for temporal aspects. Data preprocessing for time series involves parsing and transforming timestamps, handling missing values over time, and aligning data for temporal comparisons. Time series data is often visualized in line charts, bar charts, or other time-based visualizations.
- For large datasets, data reduction techniques may be applied to simplify the data without losing important insights. Sampling or summarization methods can be used to reduce data size and improve visualization performance, particularly when dealing with data that cannot be visualized directly due to its volume.
- Data preprocessing is a crucial step in data handling, ensuring that the data is suitable for the chosen visualization method and modeling algorithms. Preprocessing tasks may include scaling, standardization, and feature engineering to enhance the data's suitability for analysis and visualization.
- Data filtering allows users to focus on specific subsets of data that are relevant to the visualization task. Filtering can be based on specific criteria, time periods, or other user-defined constraints. Data filtering enables users to explore and present targeted insights without overwhelming visualizations with excessive information.
- Data exploration is an iterative process throughout data handling, where analysts visually inspect the data to gain insights, identify potential issues, and make informed decisions about data transformation and visualization choices. Data exploration allows users to understand the data's characteristics, identify interesting patterns, and uncover relationships that may inform the design of visualizations.
- In conclusion, data handling is a critical phase in data visualization that involves acquiring, cleaning, transforming, and preparing the data for visualization. Proper data handling ensures that the data is accurate, structured, and aligned with the objectives of the visualization task. It lays the foundation for effective data visualization, enabling users to create meaningful and insightful visual representations of complex information. By effectively handling the data, data visualization becomes a powerful tool for data exploration, analysis, and communication, supporting data-driven decision-making and enhancing our understanding of the underlying data patterns and trends.

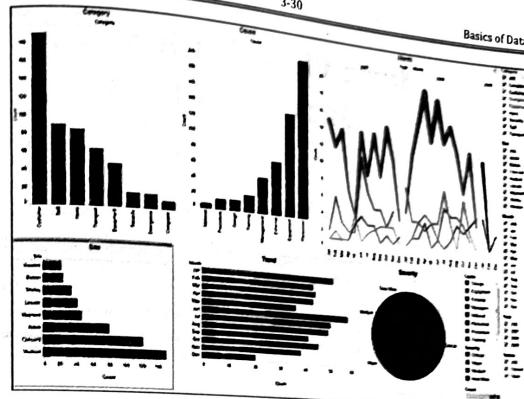


Fig. 3.8.2 : Handling categorical data

#### Review Questions

- Q. 1 What are the types of data visualization ?
- Q. 2 What is data handling?
- Q. 3 Explain key aspects of computational statistics.
- Q. 4 Explain various strategies to improve graphic speed and power.
- Q. 5 Discuss examples of higher-dimensional displays in data visualization t-SNE visualization.



DATA MODELING AND VISUALIZATION

DATA HANDLING

COMPUTATIONAL STATISTICS

DATA EXPLORATION

DATA VISUALIZATION

# 4

# Data Visualization and Data Wrangling

## Syllabus

- Data Wrangling :** Hierarchical Indexing, Combining and Merging Data Sets Reshaping and Pivoting. Data Visualization matplotlib : Basics of matplotlib, plotting with pandas and seaborn, other python visualization tools  
**Data Visualization Through their Graph Representations :** Data and Graphs Graph Layout Techniques, Force-directed Techniques Multidimensional Scaling, The Pulling Under Constraints Model, Bipartite Graphs.

## 4.1 Data Visualization

Data visualization is the graphical representation of information and data to convey insights, patterns, and trends more effectively. It involves using visual elements like charts, graphs, and maps to present data in a way that is easy to understand and interpret. Data visualization is an essential part of the data analysis process as it helps to:

- 1. Communicate Insights :** Visualizations make complex data more accessible and comprehensible to a wide audience. They allow analysts and decision-makers to communicate their findings and insights in a clear and concise manner.
- 2. Identify Patterns and Trends :** By visualizing data, patterns, trends, and correlations that might be difficult to spot in raw numbers can become apparent, helping to gain a deeper understanding of the data.
- 3. Discover Anomalies and Outliers :** Visualizations can highlight unusual data points or outliers, which may be significant for further investigation.
- 4. Explore Relationships :** Graphs and scatter plots can display relationships between variables, making it easier to explore the connections between different data points.
- 5. Support Decision Making :** Visualizations aid decision-making processes by presenting data in a format that enables users to grasp information quickly and make informed choices.

Popular types of data visualizations include :

- 1. Line Charts :** Used to display trends and changes over time.
- 2. Bar Charts :** Used to compare categorical data.
- 3. Pie Charts :** Used to show proportions or percentages of a whole.
- 4. Scatter Plots :** Used to examine the relationship between two variables.
- 5. Histograms :** Used to display the distribution of numerical data.

- 6. Heatmaps :** Used to visualize patterns and correlations in large datasets.
- 7. Geographic Maps :** Used to represent data spatially.
- 8. Box Plots :** Used to show the distribution of data and identify outliers.
- 9. Area Charts :** Used to visualize accumulated data over time.
- 10. Bubble Charts :** Used to display three dimensions of data through the size of the bubbles.

Data visualization is employed across various fields, including business, finance, marketing, healthcare, social sciences, engineering, and more. With the abundance of data available today, effective data visualization is crucial for extracting valuable insights and making data-driven decisions. It is achieved using various data visualization tools and libraries, such as Matplotlib, Seaborn, Plotly, D3.js, Tableau, and many others.

## 4.2 Data Wrangling

Data wrangling, also known as data munging or data preprocessing, refers to the process of cleaning, transforming, and preparing raw data into a format suitable for analysis. It is an essential step in the data analysis pipeline because raw data often comes in messy or unstructured formats and may contain errors, missing values, or inconsistencies. Data wrangling aims to make the data usable and ready for further analysis or modeling.

The data wrangling process typically involves the following steps :

- 1. Data Collection :** Gathering the raw data from various sources such as databases, files, APIs, web scraping, or surveys.
- 2. Data Cleaning :** Identifying and correcting errors, inconsistencies, and inaccuracies in the data. This step involves handling missing data, correcting typos, standardizing data formats, and removing duplicates.
- 3. Data Transformation :** Converting the data into a suitable format for analysis. This may involve reshaping data, aggregating information, and creating new variables or features that provide valuable insights.
- 4. Data Integration :** Combining data from different sources and integrating them into a unified dataset for analysis.
- 5. Data Reduction :** Reducing the data's size while preserving its essential characteristics. This may include sampling, filtering, or summarizing the data.
- 6. Data Enrichment :** Incorporating external data or additional information to enhance the dataset's quality and depth.
- 7. Handling Outliers :** Identifying and dealing with outliers that may skew the analysis or modeling results.
- 8. Data Formatting :** Ensuring that the data is in the appropriate data types (e.g., numeric, categorical, datetime) for analysis and modeling tools.
- 9. Handling Imbalanced Data :** Addressing issues with imbalanced datasets, which occur when one class or category is significantly underrepresented compared to others.
- 10. Feature Engineering :** Creating new features that may provide better insights or improve model performance.

Data wrangling often involves using programming languages like Python or R, along with various libraries and tools that facilitate data manipulation, such as Pandas, NumPy, and Scikit-learn in Python, and dplyr, tidyR, and ggplot2 in R. Additionally, data wrangling may require using SQL for database querying and manipulation.

Effective data wrangling is critical in ensuring that the data used for analysis is accurate, consistent, and relevant, enabling better decision-making and insights.

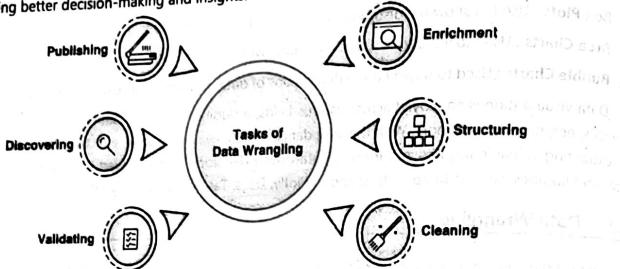


Fig. 4.2.1 : Data Wrangling

#### 4.2.1 Hierarchical Indexing

- Hierarchical indexing, also known as multi-level indexing, is a powerful feature in some data manipulation tools, such as Pandas in Python. It allows you to have multiple index levels in a Pandas DataFrame or Series. This type of indexing enables you to work with higher-dimensional data in a tabular, two-dimensional structure.
- With hierarchical indexing, you can think of your data as being organized in multiple levels of rows and columns, providing a way to represent and work with data that has complex hierarchical relationships or multiple dimensions.
- In a Pandas DataFrame, hierarchical indexing is usually achieved by creating a MultiIndex, which is an index that consists of more than one level. Each level of the index corresponds to a different dimension or category of data. For example, in a sales dataset, you might have a MultiIndex with levels for "Region," "Year," and "Month."
- Here's an example of creating a hierarchical index in Pandas :

```

import pandas as pd

# Sample data
data = {
    'Sales': [100, 150, 200, 150, 100, 200],
    'Profit': [20, 30, 40, 25, 35, 45]
}

# Create a DataFrame with hierarchical indexing
index = pd.MultiIndex.from_tuples([('North', 2021, 'Jan'), ('North', 2021, 'Feb'), ('South', 2021, 'Jan'), ('South', 2021, 'Feb')], names=['Region', 'Year', 'Month'])

df = pd.DataFrame(data, index=index)

print(df)

```

A screenshot of a Jupyter Notebook cell displaying a DataFrame named "df". The DataFrame has three columns: "Sales" (Index 0), "Profit" (Index 1), and unnamed (Index 2). The index is a MultiIndex with levels "Region", "Year", and "Month". The data is as follows:

	Sales	Profit	unnamed
Region	Year	Month	
North	2021	Jan	100
	2021	Feb	150
	2021	Mar	200
South	2021	Jan	120
	2021	Feb	180
	2021	Mar	220

Once you have created a DataFrame with hierarchical indexing, you can use it to perform various data manipulations efficiently. For example, you can use the .loc accessor to select data based on specific levels of the index or use .groupby to perform group-based operations on different levels of the index.

Hierarchical indexing is particularly useful when dealing with data with nested or hierarchical structures, time series data with multiple dimensions, or data that naturally falls into multiple categories. It allows for more sophisticated data analysis and exploration while retaining the convenience of tabular data structures.

A screenshot of a Jupyter Notebook cell displaying a DataFrame with hierarchical indexing. The index has three levels: "region", "state", and "individuals". The data is as follows:

region	state	individuals	family_members	state_pop
0 East	South Central	Alabama	2570.0	864.8
1	Pacific	Alaska	1434.0	582.9
2	Mountain	Arizona	7259.0	2656.9
3 West	South Central	Arkansas	2288.0	432.0
4	Pacific	California	18908.0	28964.0

#### 4.2.2 Combining and Merging Data Sets Reshaping and Pivoting

Combining and merging data sets, as well as reshaping and pivoting data, are important data manipulation techniques used to transform and organize data for analysis. These operations are commonly performed in data wrangling tasks to prepare data for further analysis or modeling. They are often used in conjunction with tools like Pandas in Python.

##### 1. Combining and Merging Data Sets

- Concatenation** : Concatenation is the process of combining two or more data sets along a particular axis (either rows or columns). In Pandas, the pd.concat() function is used to concatenate data frames. It allows you to stack data frames on top of each other (along rows) or side by side (along columns).
- Merging** : Merging is used to combine data sets based on common columns or indices. Pandas provides the pd.merge() function to perform different types of database-style merges, such as inner join, outer join, left join, and right join. Merging allows you to combine information from different data sets based on shared keys.

##### 2. Reshaping Data

- Pivoting** : Pivoting is a process of reshaping data by converting rows into columns. In Pandas, you can use the pd.pivot\_table() or the pd.pivot() function to create a pivot table from a DataFrame.
- Melting** : Melting is the opposite of pivoting, where columns are converted into rows. The pd.melt() function in Pandas is used to melt a DataFrame, transforming it from a wide format to a long format.

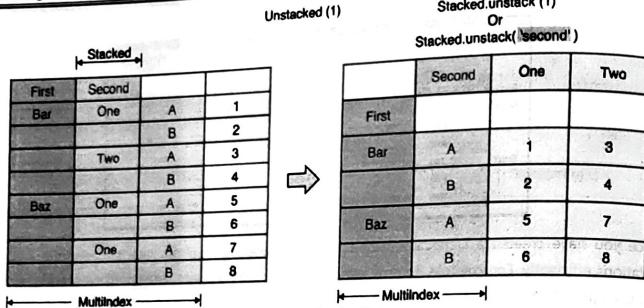


Fig. 4.2.2 : Reshaping Data

## 4.3 Data Visualization Matplotlib

- matplotlib is a popular Python library used for data visualization. It provides a wide range of functionalities for creating static, interactive, and publication-quality plots and charts. With matplotlib, you can visualize data in various formats, including line plots, scatter plots, bar plots, histograms, pie charts, and more.
- matplotlib is the most popular Python library for producing plots and other 2D data visualizations. It was originally created by John D. Hunter (JDH) and is now maintained by a large team of developers. It is well-suited for creating plots suitable for publication. It integrates well with IPython (see below), thus providing a comfortable interactive environment for plotting and exploring data. The plots are also interactive; you can zoom in on a section of the plot and pan around the plot using the toolbar in the plot window.

**Colors, Markers and Line Styles** : matplotlib's main plot function accepts arrays of X and Y coordinates and optionally a string abbreviation indicating color and line style. For example, to plot x versus y with green dashes, you would execute :

```
ax.plot(x, y, 'g-')
```

This way of specifying both color and linestyle in a string is provided as a convenience; in practice if you were creating plots programmatically you might prefer not to have to munge strings together to create plots with the desired style. The same plot could also have been expressed more explicitly as :

```
ax.plot(x, y, linestyle='--', color='g')
```

There are a number of color abbreviations provided for commonly-used colors, but any color on the spectrum can be used by specifying its RGB value (for example, '#CECECE'). You can see the full set of data points. Since matplotlib creates a continuous line plot, interpolating between points, it can occasionally marker type and line style (Fig. 4.3.1):

```
In [ ]: plt.plot(randn(30).cumsum(), 'ko--')
```

This could also have been written more explicitly as :

```
plot(randn(30).cumsum(), color='k', linestyle='dashed', marker='o')
```

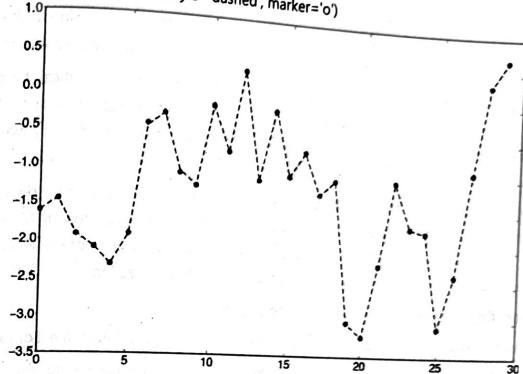


Fig. 4.3.1 : Line plot with markers example

### Matplotlib Configuration

- matplotlib comes configured with color schemes and defaults that are geared primarily toward preparing figures for publication. Fortunately, nearly all of the default behavior can be customized via an extensive set of global parameters governing figure size, subplot spacing, colors, font sizes, grid styles, and so on. There are two main ways to interact with the matplotlib configuration system. The first is programmatically from Python using the rc method. For example, to set the global default figure size to be 10 x 10, you could enter :

```
plt.rc('figure', figsize=(10, 10))
```

- The first argument to rc is the component you wish to customize, such as 'figure', 'axes', 'xtick', 'ytick', 'grid', 'legend' or many others. After that can follow a sequence of keyword arguments indicating the new parameters. An easy way to write down the options in your program is as a dictionary :

```
font_options = {'family': 'monospace', 'weight': 'bold', 'size': 'small'} plt.rc('font', **font_options)
```

- For more extensive customization and to see a list of all the options, matplotlib comes with a configuration file `matplotlibrc` in the `matplotlib/mpl-data` directory. If you can copy this file and place it in your home directory titled `.matplotlibrc`, it will be loaded each time you use matplotlib.

### 4.3.2 Plotting with Pandas and Seaborn

- matplotlib is actually a fairly low-level tool. You assemble a plot from its base components : the data display (the type of plot: line, bar, box, scatter, contour, etc.), legend, title, tick labels, and other annotations. Part of the reason for this is that in many cases the data needed to make a complete plot is spread across many objects. In pandas we have row labels, column labels, and possibly grouping information. This means that many kinds of fully-formed plots that would ordinarily require a lot of matplotlib code can be expressed in one or two concise statements. Therefore, pandas has an increasing number of high-level plotting methods for creating standard visualizations that take advantage of how data is organized in DataFrame object.

- The Python community has adopted a number of naming conventions for commonly used modules: import numpy as np import pandas as pd import matplotlib.pyplot as plt This means that when you see np.arange, this is a reference to the arange function in NumPy. This is done as it's considered bad practice in Python software development to import everything (from numpy import \*) from a large package like NumPy.
- While NumPy by itself does not provide very much high-level data analytical functionality, having an understanding of NumPy arrays and array-oriented computing will help you use tools like pandas much more effectively. If you are new to Python and just looking to get your hands dirty working with data using pandas.
- Both Pandas and Seaborn are powerful libraries that work seamlessly together for data visualization. Pandas provides a convenient interface to handle data, and Seaborn is built on top of Matplotlib, enhancing its capabilities and providing higher-level abstractions for creating visually appealing and informative plots. Here's an overview of how to use Pandas and Seaborn for data visualization :

**1. Importing Libraries :** First, import the necessary libraries :

**2. Loading Data with Pandas :** Use Pandas to load your data into a DataFrame. For example :

**3. Basic Pandas Plotting :** Pandas provides built-in functions for basic data visualization. You can use methods like plot(), plot.bar(), plot.scatter(), plot.hist(), etc., on the DataFrame :

**4. Seaborn for Enhanced Visualization :** Seaborn enhances the visual aesthetics of Matplotlib plots and offers high-level functions for more complex visualizations. You can use Seaborn's functions on Pandas DataFrames to create more informative plots with less code :

**5. Additional Customization :** Seaborn provides additional customization options for plots, such as adjusting color palettes, adding labels, setting plot size, and more.

- Remember that Seaborn will automatically improve the visual appearance of your plots, but you can further customize them using Matplotlib functions if needed.
- By combining the data manipulation capabilities of Pandas and the enhanced visualization features of Seaborn, you can efficiently explore and communicate insights from your data in an elegant and effective manner.

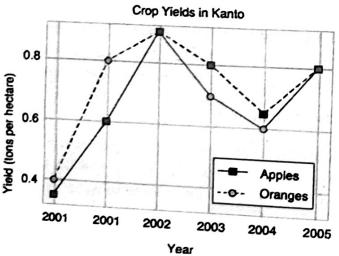


Fig. 4.3.2

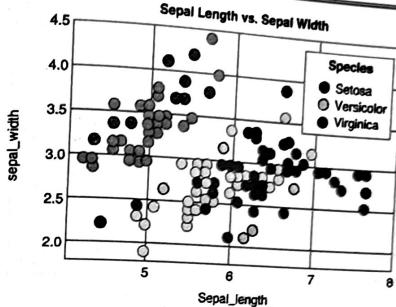


Fig. 4.3.3

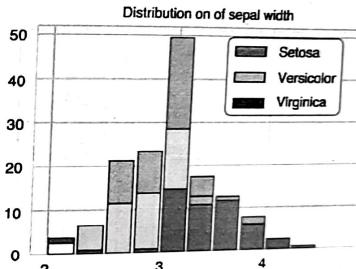


Fig. 4.3.4

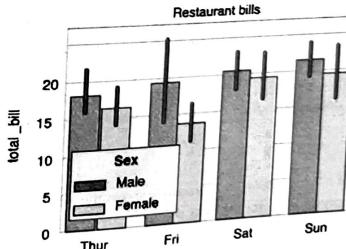


Fig. 4.3.5

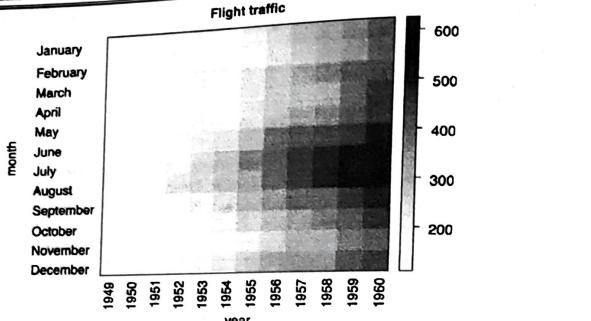


Fig. 4.3.6

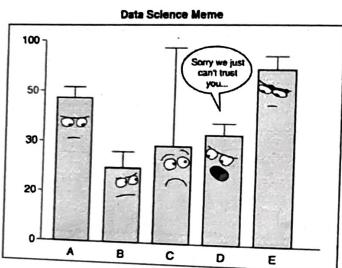


Fig. 4.3.7

- To get started with pandas, you will need to get comfortable with its two workhorse data structures: Series and DataFrame. While they are not a universal solution for every problem, they provide a solid, easy-to-use basis for most applications.
- Series A Series is a one-dimensional array-like object containing an array of data (of any NumPy data type) and an associated array of data labels, called its index. The simplest Series is formed from only an array of data:  
`obj = Series([4, 7, -5, 3])`
- The string representation of a Series displayed interactively shows the index on the left and the values on the right. Since we did not specify an index for the data, a default one consisting of the integers 0 through N - 1 (where N is the length of the data) is created. You can get the array representation and index object of the Series via its values and index attributes, respectively.  
`In [6]: obj.values`

`n [7]: obj.index`

Often it will be desirable to create a Series with an index identifying each data point:

`In [8]: obj2 = Series([4, 7, -5, 3], index=['d', 'b', 'a', 'c'])``In [9]: obj2``obj2.index`

Compared with a regular NumPy array, you can use values in the index when selecting single values or a set of values : `In [11]: obj2['a']`

### 4.3.3 Python Visualization Tool Ecosystem

- As is common with open source, there are a plethora of options for creating graphics in Python (too many to list). In addition to open source, there are numerous commercial libraries with Python bindings. In this chapter and throughout the book, I have been primarily concerned with matplotlib as it is the most widely used plotting tool in Python. While it's an important part of the scientific Python ecosystem, matplotlib has plenty of shortcomings when it comes to the creation and display of statistical graphics.
- MATLAB users will likely find matplotlib familiar, while R users (especially users of the excellent ggplot2 and trelis packages) may be somewhat disappointed (at least as of this writing). It is possible to make beautiful plots for display on the web in matplotlib, but doing so often requires significant effort as the library is designed for the printed page. Aesthetics aside, it is sufficient for most needs.
- In pandas, I, along with the other developers, have sought to build a convenient user interface that makes it easier to make most kinds of plots commonplace in data analysis. There are a number of other visualization tools in wide use. I list a few of them here and encourage you to explore the ecosystem.

#### Chaco

Chaco (<http://code.enthought.com/chaco/>), developed by Enthought, is a plotting tool-kit suitable both for static plotting and interactive visualizations. It is especially well suited for expressing complex visualizations with data interrelationships. Compared with matplotlib, Chaco has much better support for interacting with plot elements and rendering is very fast, making it a good choice for building interactive GUI applications.

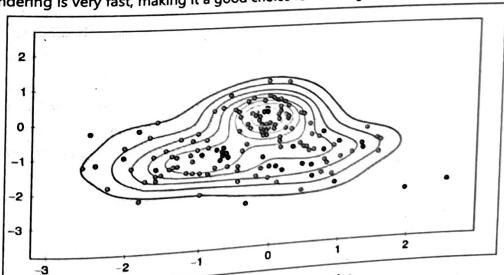


Fig. 4.3.8 : A Chaco example plot

**mayavi**

The mayavi project, developed by Prabhu Ramachandran, Gaël Varoquaux, and others, is a 3D graphics toolkit built on the open source C++ graphics library VTK. mayavi, like matplotlib, integrates with IPython so that it is easy to use interactively. The plots can be panned, rotated, and zoomed using the mouse and keyboard. While I don't show any mayavi-using code here, there is plenty of documentation and examples available online. In many cases, I believe it is a good alternative to a technology like WebGL, though the graphics are harder to share in interactive form.

**Other Packages**

Of course, there are numerous other visualization libraries and applications available in Python: PyQwt, Vusz, gnuplot-py, biggles, and others. I have seen PyQwt put to good use in GUI applications built using the Qt application framework using PyQt. While many of these libraries continue to be under active development (some of them are part of much larger applications), I have noted in the last few years a general trend toward web-based technologies and away from desktop graphics. I'll say a few more words about this in the next section.

**Future of Visualization Tools**

Visualizations built on web technologies (that is, JavaScript-based) appear to be the inevitable future. Doubtlessly you have used many different kinds of static or interactive visualizations built in Flash or JavaScript over the years. New toolkits (such as d3.js and its numerous off-shoot projects) for building such displays are appearing all the time. In contrast, development in non web-based visualization has slowed significantly in recent years. This holds true of Python as well as other data analysis and statistical computing environments like R. The development challenge, then, will be building tighter integration between data analysis and preparation tools, such as pandas, and the web browser. I am hopeful that this will become a fruitful point of collaboration between Python and non-Python users as well.

#### **4.4 Data Visualization Through their Graph Representations**

Apart from Matplotlib, Pandas, and Seaborn, there are several other Python libraries and tools for data visualization. Each of these tools has its own strengths and focuses on specific types of visualizations or interactivity. Here are some popular ones :

- Plotly** : Plotly is a powerful library for interactive and web-based visualizations. It supports a wide range of chart types, including line plots, scatter plots, bar plots, heatmaps, 3D plots, and more. Plotly can produce interactive visualizations that can be embedded in web applications or Jupyter Notebooks.
- Bokeh** : Bokeh is another library for creating interactive visualizations in Python. It focuses on providing minimal coding, making it useful for building interactive dashboards and data applications.
- Altair** : Altair is a declarative statistical visualization library that allows users to build visualizations by specifying the data transformation and visual encoding rules. It is built on top of Vega-Lite, which is a high-level grammar for visual analysis.
- Holoviews** : Holoviews is a high-level library that simplifies the creation of interactive visualizations. It allows users to create complex visualizations by composing data objects with concise Python code.

- Geopandas** : Geopandas is an extension of Pandas, specifically designed for working with geospatial data. It allows for plotting geospatial data on maps and creating choropleth maps.
- Folium** : Folium is a Python library for creating interactive maps using Leaflets. It allows for visualizing data on interactive maps with various map styles and markers.
- Networkx** : Networkx is a library for working with complex networks and graphs. It provides tools to create, analyze, and visualize graph data structures.
- Yellowbrick** : Yellowbrick is a visualization library specifically designed for machine learning. It offers tools to visualize model evaluation, feature analysis, and hyper parameter tuning.
- D3.js** : Although not a Python library, D3.js is a popular JavaScript library for creating interactive and data-driven visualizations on the web. Python developers can use D3.js through Python wrappers like mpld3. Each of these libraries has its unique features and use cases, so choosing the right visualization tool depends on the specific requirements of your data analysis and the type of visualizations you want to create.

#### **4.4.1 Data and Graphs Graph Layout Techniques**

In the context of data visualization and graph theory, graph layout techniques refer to the methods used to determine the positions of nodes (vertices) and edges in a graph, so it can be visually represented on a 2D or 3D plane. Proper graph layout is crucial for creating clear and meaningful visualizations of complex networks and relationships. Different graph layout algorithms aim to arrange nodes and edges in ways that maximize clarity, minimize overlaps, and highlight meaningful patterns or structures within the data.

**Here are some common graph layout techniques**

- Force-Directed Layout** : Force-directed layout is a widely used technique for arranging nodes in a graph. It simulates physical forces between nodes and edges, where nodes repel each other, and edges act like springs pulling connected nodes together. This algorithm seeks to find a balance where the forces are in equilibrium, resulting in a visually pleasing and distributed layout.
- Circular Layout** : In a circular layout, nodes are positioned in a circular pattern, often based on their order or specific attributes. It is suitable for visualizing cyclic graphs or situations where a hierarchical or radial structure is evident.
- Hierarchical Layout** : Hierarchical layout arranges nodes in levels or layers, showing hierarchical relationships. This approach is commonly used for visualizing trees or organizational structures.
- Tree Layout** : Similar to hierarchical layout, the tree layout is used to represent hierarchical relationships, where nodes are placed in a tree-like structure. It is ideal for visualizing decision trees and other hierarchical data.
- Layered Layout (Layered Drawings)** : Layered layout techniques divide the graph into layers, with nodes in each layer aligned horizontally or vertically. This approach is particularly useful for visualizing directed acyclic graphs (DAGs) and flowcharts.
- Kamada-Kawai Layout** : Kamada-Kawai is a graph layout algorithm that considers the distance between nodes in the graph as a cost function to minimize. It aims to position nodes in a way that preserves the graph's overall topology and minimizes edge crossings.

7. **Fruchterman-Reingold Layout** : This is another force-directed layout algorithm that places nodes based on a spring-repulsion model. It attempts to find an optimal layout that balances attractive forces between connected nodes and repulsive forces between all nodes.
8. **Spectral Layout** : Spectral layout techniques use the eigenvalues and eigenvectors of the graph's adjacency or Laplacian matrix to derive node positions. They can reveal the underlying structure of the graph and are useful for clustering and community detection.
9. **Grid Layout** : Grid layout arranges nodes on a regular grid pattern. It is suitable for organizing small and simple graphs, ensuring a consistent and structured representation.

Choosing the appropriate graph layout technique depends on the characteristics of the data, the type of graph being visualized, and the insights you want to highlight. Some graph layout algorithms are available in popular network visualization libraries, such as Networkx (Python) and Cytoscape (Java/JavaScript), while others can be implemented using custom code or specialized graph visualization tools.

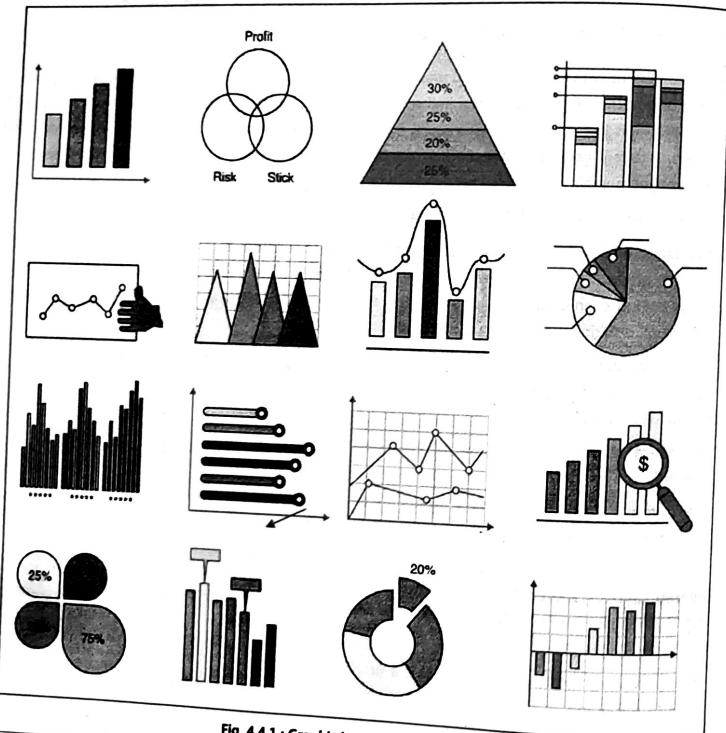


Fig. 4.4.1 : Graphic layout techniques

#### 4.4.2 Force-Directed Techniques Multidimensional Scaling

Force-directed techniques and Multidimensional Scaling (MDS) are two different approaches used in graph visualization and dimensionality reduction, both serving distinct purposes.

1. **Force-Directed Techniques** : Force-directed techniques are used to layout graphs in a visually appealing way. These techniques simulate physical forces between nodes and edges to determine their positions. Nodes repel each other, while edges act like springs pulling connected nodes together. The objective is to find a balanced configuration where the forces are in equilibrium, resulting in a layout that minimizes edge crossings and node overlaps. This approach aims to capture the global structure of the graph, highlighting clusters and patterns.
2. **Multidimensional Scaling (MDS)** : Multidimensional Scaling is a dimensionality reduction technique used for visualizing high-dimensional data in a lower-dimensional space. It is primarily used in data analysis and visualization, not specifically for graph layout. MDS attempts to preserve the pairwise distances (or dissimilarities) between data points in the lower-dimensional space as closely as possible to their original high-dimensional distances.

MDS can be used with various distance metrics, such as Euclidean distance or correlation distance, and is particularly useful when trying to understand the underlying structure or similarity relationships in data. In the context of graph visualization, MDS can be applied to analyze and visualize the similarity between nodes based on their connections or attributes.

##### Difference between Force-directed Techniques and MDS

The main difference between force-directed techniques and MDS lies in their objectives and application domains :

- Force-directed techniques are used for graph layout to create visually appealing representations of graphs, highlighting the overall structure and relationships between nodes and edges.

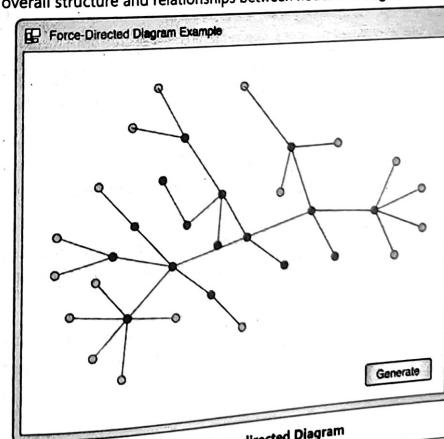


Fig. 4.4.2 : Force-directed Diagram

- Multidimensional Scaling, on the other hand, is primarily used for dimensionality reduction and data visualization. It is applied to analyze the underlying similarity or dissimilarity relationships among data points in a high-dimensional space by projecting them into a lower-dimensional space.
- While both force-directed techniques and MDS are used in visualization, they have distinct applications and goals. Force-directed techniques are specific to graph layout, whereas MDS is more general and applicable to a broader range of data visualization tasks, especially when trying to understand the structure and relationships in high-dimensional data.

#### 4.4.3 The Pulling under Constraints Model

- The Pulling Under Constraints (PUC) model is a specific approach used in force-directed graph layout algorithms. In force-directed graph layout, the goal is to find an arrangement of nodes and edges that represents the graph in a visually pleasing way, with reduced edge crossings and node overlaps. The PUC model introduces additional constraints on the positioning of nodes and edges to achieve more controlled and meaningful graph layouts.
- In traditional force-directed graph layout algorithms, nodes are treated as charged particles that repel each other, while edges act as springs pulling connected nodes together. These forces are balanced to find an equilibrium configuration. However, in complex graphs, this approach can result in chaotic or undesirable layouts.
- The PUC model addresses this issue by introducing constraints on the positions of nodes and edges based on specific graph characteristics or user-defined rules. These constraints aim to impose a level of structure and control on the graph layout, leading to more informative and organized visualizations.

**The key features of the Pulling Under Constraints model include**

- Fixed Positions** : The PUC model allows some nodes to have fixed positions, meaning their positions are pre-defined and not affected by the forces during the layout process. This feature is useful when you want to anchor specific nodes in predetermined positions, like in hierarchical or tree structures.
- Inertial Constraints** : Inertia is added to the forces acting on nodes, simulating their resistance to rapid changes in position. This ensures that nodes move smoothly, which helps in stabilizing the layout and reducing oscillations during the layout process.
- Edge Bundling** : PUC supports edge bundling, where multiple edges between the same pair of nodes are grouped together into a single curved path. Edge bundling helps in reducing visual clutter caused by dense graph structures.
- Edge Routing Constraints** : Constraints are applied to edges to guide their routes and prevent overlapping with other edges or nodes. This ensures that the edges follow meaningful paths and do not obscure important information.
- Hierarchical Layout** : PUC allows for hierarchical layout, where nodes are organized into levels or layers, and edges are mostly routed vertically or horizontally between levels, following the hierarchical structure of the graph.

The Pulling Under Constraints model is particularly useful for visualizing complex graphs, such as social networks, biological networks, or large-scale datasets, where traditional force-directed algorithms may lead to cluttered and hard-to-interpret layouts. By incorporating additional constraints and rules, the PUC model helps create more structured and informative graph visualizations. It has been implemented in several graph visualization libraries and software tools, supporting researchers and analysts in gaining insights from complex network data.

#### 4.4.4 Bipartite Graphs

A bipartite graph, also known as a bigraph, is a type of graph in graph theory that consists of two sets of nodes, such that all edges in the graph connect nodes from one set to nodes from the other set. In other words, there are no edges that connect nodes within the same set. Formally, a graph  $G(V, E)$  is bipartite if its vertex set  $V$  can be partitioned into two disjoint sets  $V_1$  and  $V_2$ , such that every edge in  $E$  connects a node from  $V_1$  to a node from  $V_2$ .

In a bipartite graph, there are no cycles of odd length. This property is known as the "bipartite property." As a result, bipartite graphs are planar, and they do not contain any odd cycles (e.g., triangles) in their structure.

**Bipartite graphs find numerous applications in various real-world scenarios such as :**

- Matching Problems** : Bipartite graphs are often used to model matching problems, where the goal is to find pairs of elements from the two sets that satisfy certain criteria. For example, matching students with courses, employees with tasks, or advertisers with target customers.
- Recommendation Systems** : In recommendation systems, bipartite graphs can be used to represent user-item interactions. Each set corresponds to users and items, and edges represent interactions (e.g., user ratings, product purchases).
- Social Network Analysis** : Bipartite graphs can be used to analyze social networks, where one set represents users and the other set represents groups or communities. Edges represent memberships or affiliations.
- Transportation Networks** : In transportation systems, bipartite graphs can be used to model the flow of goods between suppliers and consumers or transportation routes between origin and destination points.
- Resource Allocation** : Bipartite graphs can be used to represent resource allocation problems, such as assigning tasks to workers or allocating funds to projects.
  - To visually represent a bipartite graph, nodes from one set are typically placed on the left side of the graph, and nodes from the other set are placed on the right side. Edges connect nodes between the two sets.
  - Bipartite graphs can be analyzed and manipulated using graph algorithms and techniques, and they provide valuable insights into the relationships and interactions between entities in various domains.

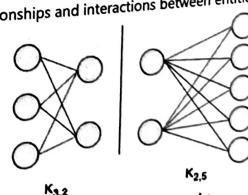


Fig. 4.4.3 : Bipartite graphs

# 5

# Data Aggregation and Analysis

## Syllabus

**Data Aggregation and Group operations :** Group by Mechanics, Data aggregation, General split-apply-combine, Pivot tables and cross tabulation 67 Time Series

**Data Analysis :** Date and Time Data Types and Tools, Time series Basics, date Ranges, Frequencies and Shifting, Time Zone Handling, Periods and Periods Arithmetic, Resampling and Frequency conversion, Moving Window Functions

## 5.1 Introduction : Data Aggregation

- Data aggregation is any process in which information is gathered and expressed in a summary form, for purposes such as statistical analysis.
- A common aggregation purpose is to get more information about particular groups based on specific variables such as age, profession or income.
- The information about such groups can then be used for Web site personalization to choose content and advertising likely to appeal to an individual belonging to one or more groups for which data has been collected. For example, a site that sells music CDs might advertise certain CDs based on the age of the user and the data aggregate for their age group.
- Online Analytic Processing (OLAP) is a simple type of data aggregation in which the marketer uses an online reporting mechanism to process the information.
- A common aggregation purpose is to get more information about particular groups based on specific variables such as age, profession, or income. The information about such groups can then be used for Web site personalization to choose content and advertising likely to appeal to an individual belonging to one or more groups for which data has been collected. For example, a site that sells music CDs might advertise certain CDs based on the age of the user and the data aggregate for their age group. Online Analytic Processing (OLAP) is a simple type of data aggregation in which the marketer uses an online reporting mechanism to process the information.

## 5.2 Data Aggregation and Group Operations

- Data aggregation and group operations are essential techniques in data analysis and are particularly useful when dealing with large datasets. These methods allow analysts to summarize and gain insights from the data efficiently.

- Categorizing a dataset and applying a function to each group, whether an aggregation or transformation, is often a critical component of a data analysis workflow. After loading, merging, and preparing a dataset, you may need to compute group statistics or possibly pivot tables for reporting or visualization purposes. pandas provides a flexible group by interface, enabling you to slice, dice, and summarize datasets in a natural way.
- One reason for the popularity of relational databases and SQL (which stands for "structured query language") is the ease with which data can be joined, filtered, transformed, and aggregated. However, query languages like SQL are somewhat constrained in the kinds of group operations that can be performed. As you will see, with the expressiveness of Python and pandas, we can perform quite complex group operations by utilizing any function that accepts a pandas object or NumPy array. In this chapter, you will learn how to :
  - Split a pandas object into pieces using one or more keys (in the form of functions, arrays, or Data Frame column names).
  - Calculate group summary statistics, like count, mean, or standard deviation, or a user-defined function.
  - Apply within-group transformations or other manipulations, like normalization, linear regression, rank, or subset selection.
  - Compute pivot tables and cross-tabulations.
  - Perform quantile analysis and other statistical group analyses.

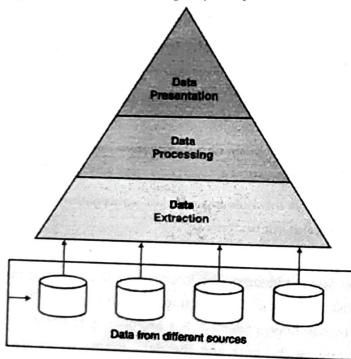


Fig. 5.2.1 : Data Aggregation

### 5.2.1 Group by Mechanics

Group by Mechanics in data modeling and visualization is a fundamental technique used to organize and summarize data based on specific criteria. This method plays a crucial role in data analysis, enabling analysts to gain insights and make data-driven decisions. The "Group by" operation involves dividing a dataset into distinct groups based on the values of one or more key variables and then performing calculations or aggregations within each group. It is widely supported in various data manipulation and visualization tools, including spreadsheet software like Microsoft Excel and data analysis libraries like pandas in Python.

### 1. Data Organization

Before applying the "Group by" operation, the data needs to be structured in a tabular format, often represented as a data frame or a table with rows and columns. Each row in the dataset represents a single observation, and each column represents a specific attribute or variable associated with that observation.

### 2. Identifying Key Variables

To use the "Group by" operation, one or more key variables need to be selected to define the groups. These variables should have categorical or discrete values, such as product categories, geographical regions, or customer IDs. The "Group by" operation will group the data based on the unique values of these key variables.

### 3. Grouping Process

Once the key variables are chosen, the "Group by" operation will group the data according to the unique values of these variables. All rows with the same value(s) for the key variable(s) will be combined into a separate group. As a result, the dataset is partitioned into multiple subsets, each representing a distinct group.

### 4. Aggregation and Calculation

- After the data is grouped, analysts can perform various aggregation functions and calculations on each group independently. Common aggregation functions include sum, mean, median, count, standard deviation, and more. These calculations provide insights into the characteristics and patterns within each group.
- For example, consider a sales dataset with columns for product categories, sales dates, and sales amounts. By grouping the data based on the product categories, you can obtain separate groups for each category. You can then apply the sum function to calculate the total sales for each product category, providing a concise summary of the sales data for different products.

### 5. Multi-Level Grouping

- Group by Mechanics also supports multi-level grouping, where you can use multiple key variables to create a hierarchical grouping structure. This allows for deeper analysis by drilling down into subgroups within larger groups.
- For instance, you can group sales data by both product categories and sales regions, providing insights into how different products perform in various regions.

### 6. Visualization

Once the data is grouped and aggregated, visualizing the results can be highly effective in communicating insights. Bar charts, line charts, pie charts, and other visualizations can help compare and contrast the data across different groups. Visualizations make it easier to identify trends, patterns, and outliers within the dataset.

### 7. Interpretation and Decision-Making

The insights gained from Group by Mechanics enable data-driven decision-making. Analysts can identify top-performing categories, regions with the highest sales, or trends in customer behavior, which can guide marketing strategies, inventory management, and resource allocation.

## 5.2.2 Data Aggregation

- Data Aggregation is any process whereby data is gathered and expressed in a summary form. When data is aggregated, atomic data rows - typically gathered from multiple sources - are replaced with totals or summary statistics. Groups of observed aggregates are replaced with summary statistics based on those observations. Aggregate data is typically found in a data warehouse, as it can provide answers to analytical questions and also dramatically reduce the time to query large sets of data.
- Data aggregation can enable analysts to access and examine large amounts of data in a reasonable time frame. A row of aggregate data can represent hundreds, thousands or even more atomic data records. When the data is aggregated, it can be queried quickly instead of requiring all of the processing cycles to access each underlying atomic data row and aggregate it in real time when it is queried or accessed.
- As the amount of data stored by organizations continues to expand, the most important and frequently accessed data can benefit from aggregation, making it feasible to access efficiently.

### What does data aggregation do?

Data aggregators summarize data from multiple sources. They provide capabilities for multiple aggregate measurements, such as sum, average and counting.

### Examples of aggregate data :

- Voter turnout by state or county. Individual voter records are not presented, just the vote totals by candidate for the specific region.
- Average age of customer by product. Each individual customer is not identified, but for each product, the average age of the customer is saved.
- Number of customers by country. Instead of examining each customer, a count of the customers in each country is presented.
- Data aggregation can also result in a similar effect to data anonymization - as individual data elements with personally identifiable details are combined and replaced with a summary representing a group as a whole. An example of this is creating a summary that shows the aggregate average salary for employees by department, rather than browsing through individual employee records with salary data.
- Aggregate data does not need to be numeric. You can, for example, count the number of any non-numeric data element.
- Before aggregating, it is crucial that the atomic data is analyzed for accuracy and that there is enough data for the aggregation to be useful. For example, counting votes when only 5% of results are available is not likely to produce a relevant aggregate for prediction.

### How do data aggregators work?

- Data aggregators work by combining atomic data from multiple sources, processing the data for new insights and presenting the aggregate data in a summary view. Furthermore, data aggregators usually provide the ability to track data lineage and can trace back to the underlying atomic data that was aggregated.

1. **Collection :** First, data aggregation tools may extract data from multiple sources, storing it in large databases as atomic data. The data may be extracted from Internet of Things (IoT) sources, such as the following:
  - social media communications;
  - news headlines;
  - personal data and browsing history from IoT devices; and
  - call centers, podcasts, etc. (through speech recognition).
2. **Processing :** Once the data is extracted, it is processed. The data aggregator will identify the atomic data that is to be aggregated. The data aggregator may apply predictive analytics, artificial intelligence (AI) or machine learning algorithms to the collected data for new insights. The aggregator then applies the specified statistical functions to aggregate the data.
3. **Presentation :** Users can present the aggregated data in a summarized format that itself provides new data. The statistical results are comprehensive and high quality.
- Data aggregation may be performed manually or through the use of data aggregators. However, data aggregation is often performed on a large-scale basis, which makes manual aggregation less feasible. Furthermore, manual aggregation risks accidental omission of crucial data sources and patterns.

## 5.2.2(A) Uses for Data Aggregation

- Data aggregation can be helpful for many disciplines, such as finance and business strategy decisions, product planning, product and service pricing, operations optimization and marketing strategy creation. Users may be data analysts, data scientists, data warehouse administrators and subject matter experts.
- Aggregated data is commonly used for statistical analysis to obtain information about particular groups based on specific demographic or behavioral variables, such as age, profession, education level or income.
- For business analysis purposes, data can be aggregated into summaries that help leaders make well-informed decisions. User data can be aggregated from multiple sources, such as social media communications, browsing history from IoT devices and other personal data, to give companies critical insights into consumers.

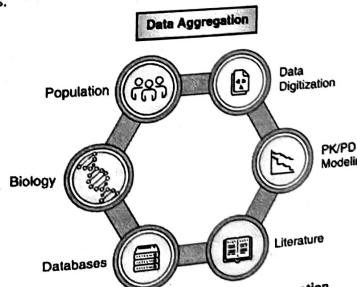


Fig. 5.2.2 : Applications of Data Aggregation

- Data aggregation is a crucial technique in data modeling and visualization, used to summarize and condense large datasets into more manageable and insightful representations. This process involves applying mathematical or statistical operations to combine data within each group, leading to concise and meaningful results. Data aggregation simplifies complex data, allowing analysts to identify patterns, trends, and overall characteristics, which are vital for data-driven decision-making. Here's a comprehensive explanation of data aggregation :

## 1. Aggregation Functions

Aggregation functions are mathematical or statistical operations that consolidate data within each group. Common aggregation functions include sum, mean (average), median, count, standard deviation, minimum, and maximum. These functions help create summary statistics, enabling users to understand the distribution and characteristics of the data effectively.

## 2. Data Organization

Before performing data aggregation, the data should be structured in a tabular format, typically represented as a data frame or a database table. Each row represents an individual observation, while each column contains attributes or variables associated with those observations.

## 3. Group by Mechanics

Data aggregation often goes hand in hand with the "Group by" operation. By grouping the data based on one or more key variables, the dataset is divided into distinct groups. Aggregation functions are then applied to each group independently, resulting in summary statistics for each group.

For instance, consider a sales dataset with columns for product categories, sales dates, and sales amounts. After grouping the data by product categories, you can apply the sum function to calculate the total sales for each category, providing a concise overview of sales performance across different products.

## 4. Multi-Level Aggregation

Similar to multi-level grouping, data aggregation also supports multi-level aggregation. In this case, you can apply aggregation functions to multiple key variables, creating a hierarchical summary. This enables deeper insights by examining subgroups within larger groups.

Continuing with the previous example, you could perform multi-level aggregation by both product categories and sales regions. This would provide a comprehensive view of total sales for each product category across different regions.

## 5. Time-based Aggregation

In time series data, data aggregation is crucial for summarizing and analyzing trends over time. Time-based aggregation involves grouping data into specific time intervals, such as days, weeks, or months, and applying aggregation functions to calculate relevant statistics for each interval.

Time-based aggregation can help identify seasonality, trends, and patterns in time series data. For instance, you can calculate the average daily sales for each month or the total monthly revenue for a specific product.

## 6. Visualization

After performing data aggregation, visualizing the results can be highly effective in presenting the summarized information. Bar charts, line charts, area charts, and other visualizations can provide a clear representation of aggregated data, making it easier to compare and analyze trends across different groups.

## 7. Decision-Making and Analysis

Data aggregation plays a pivotal role in data-driven decision-making. By summarizing large datasets into meaningful insights, analysts can make informed choices, identify growth opportunities, optimize business processes, and address potential issues.

### 5.2.3 Split Apply Combine

In Python we do this by using GroupBy and it involves one or more of the three steps of the Split-Apply-Combine strategy. Let us start by defining each of the three steps :

Fig. 5.2.3 shows the Split-Apply-Combine using an aggregation function.

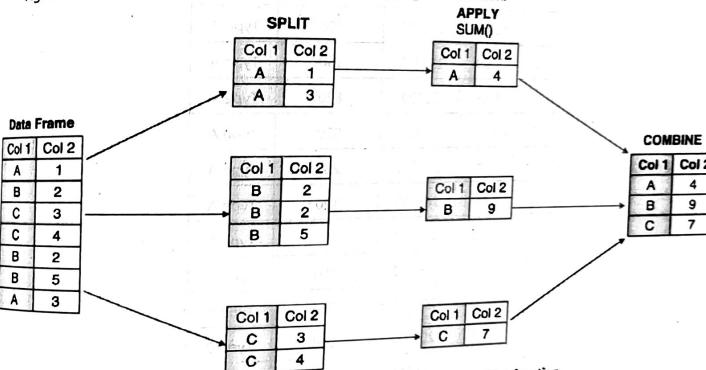


Fig. 5.2.3 : The Split-Apply-Combine using an aggregation function.

- Split** : Split the data into groups based on some criteria thereby creating a Group by object. (We can use the column or a combination of columns to split the data into groups)
- Apply** : Apply a function to each group independently. (Aggregate, Transform, or Filter the data in this step)
- Combine** : Combine the results into a data structure (Pandas Series, Pandas Data Frame) Dataset

- To go a bit deeper, let's create a fictitious data to serve as an example. Have a thorough look at the data frame (`data_sales`) given below, because it will be used throughout this article.
- To access the code used in this article please visit this link: [TechKnowledge Publications](#)

## Data Modeling and Visualization

Import Libraries and create a small dataset to work on.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Create an Example Data-set in the form of dictionary having key value pairs.

```
sales_dict={'colour':['Yellow','Black','Blue','Red','Yellow','Black','Blue',
                     'Red','Yellow','Black','Blue','Red','Red','Blue','Red'],
            'sales':[100000,150000,200000,90000,120000,
                     100000,150000,200000,90000,120000,140000,130000,400000,350000],
            'transactions':[100,150,200,90,130,70,250,150,100,130,80,90,300,150,170,230,280],
            'product':[('type A','type A','type A','type A','type A','type A','type A',
                       'type A','type A','type A','type A','type A','type A','type A',
                       'type A','type B','type B','type B','type B','type B','type B')]
```

```
data_sales=pd.DataFrame(sales_dict)
```

```
data_sales
```

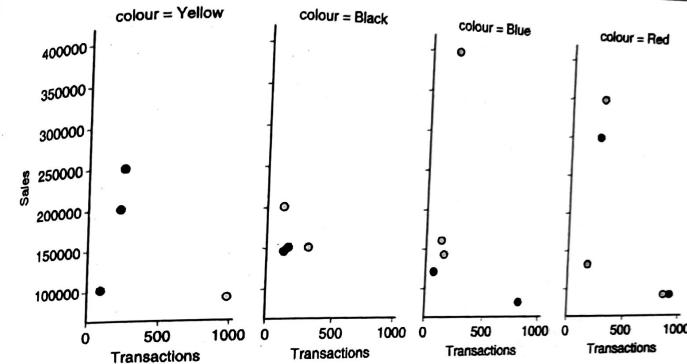
	Colour	Sales	Transactions	Product
0	Yellow	100000	100	Type A
1	Black	150000	150	Type A
2	Blue	80000	820	Type A
3	Red	90000	920	Type A
4	Yellow	200000	230	Type A
5	Black	145000	120	Type A
6	Blue	120000	70	Type A
7	Red	300000	250	Type A
8	Yellow	250000	250	Type A
9	Black	200000	110	Type B
10	Blue	160000	130	Type B
11	Red	90000	860	Type B
12	Yellow	90100	980	Type B
13	Black	150000	300	Type B
14	Blue	142000	150	Type B
15	Red	130000	170	Type B
16	Blue	400000	230	Type B
17	Red	350000	280	Type B

## data\_sales

To summarize the whole data, seaborn library have been used to create a visualization, which includes all the data graphically.

```
graph = sns.FacetGrid(data_sales, col="colour", height=4, hue="product", aspect=.5)
```

```
graph.map(plt.scatter, "transactions", "sales");
```



## Visual Data Summary

## Product

Type A ■■■■■

Type B □□□□□

Fig. 5.2.4

After creating and summarizing the data, as a first step let's move on to the first part of Split-Apply-Combine.

## SPLIT : Create an Object

- In this step we will create the the groups from the dataframe 'data\_sales' by grouping on the basis of the column 'colour'.

```
# SPLIT: Groupby the column 'colour'
data_by = data_sales.groupby('colour')
print(type(data_by))
<class 'pandas.core.groupby.groupby.DataFrameGroupBy'>
```

- Once we apply the groupby() function on the dataframe, it creates the GroupBy object as a result. We can think of this object as a separate dataframe for each group. Each group has been created based on the categories in a grouped column (4 Groups will be created 'Black', 'Blue', 'Red', 'Yellow' from the column 'colour' of the dataframe in our case).

- A GroupBy object stores the data of the individual groups in the form of key value pairs as in dictionary. To know the group names , we can either use attribute 'keys' or use the attribute 'groups' of the GroupBy object.

```
# lets check the names of the groups
```

```
data_gbby.groups
('Black': Int64Index([1, 5, 9, 13], dtype='int64'),
'Blue': Int64Index([2, 6, 10, 14, 16], dtype='int64'),
'Red': Int64Index([3, 7, 11, 15, 17], dtype='int64'),
'Yellow': Int64Index([8, 4, 8, 12], dtype='int64'))
```

For further clarity on Groups and its content, we can run a loop and print the key value pairs.

```
#>>> 'key' is the name of the group and 'value' is the segmented rows from the original DataFrame.
```

```
for key, value in data_gbby:
    print('Groupname: ', key)
    print(value)
    print('-----')
```

```
Groupname: Black
colour  sales  transactions  product
1   Black  150000      150  type A
5   Black  145000      120  type A
9   Black  200000      110  type B
13  Black  150000      300  type B
-----
```

```
Groupname: Blue
colour  sales  transactions  product
2   Blue   120000      80  type A
6   Blue   130000      70  type A
10  Blue   160000      130  type B
14  Blue   142000      150  type B
16  Blue   400000      230  type B
-----
```

```
Groupname: Red
colour  sales  transactions  product
3   Red    90000       90  type A
7   Red   300000      250  type A
11  Red   98000       80  type B
15  Red   130000      170  type B
17  Red   350000      280  type B
-----
```

```
Groupname: Yellow
colour  sales  transactions  product
0   Yellow  180000      80  type A
4   Yellow  200000      230  type A
8   Yellow  250000      250  type A
12  Yellow  90100      980  type B
-----
```

- With the above example, I hope we have developed some clarity on the GroupBy object along with some of its attributes and methods. With this, now let's move forward to the next stage, which is APPLY.

#### APPLY : Apply some function on the Object.

Apply step can performed in three ways : **Aggregation, Transformation, and Filtering**. We all have good amount of experience in using Aggregation with GroupBy objects, but most of us might not have the same experience with the Transformation and Filtering. Here, we will discuss all the three with special focus on Transformation.

#### AGGREGATION

I am assuming that we are already comfortable with applying the aggregation functions with GroupBy object, therefore I will start of with some interesting features of this function.

By choosing multiple columns to create the group, we increase the granularity of the aggregation. For instance, while splitting we created 4 groups based on the column 'colours', which has 4 categories of colours, so we had 4 groups. Now, if include 'product' column, having 2 categories ('type A' and 'type B'), along with the 'colour' column, then we will be having total 8 categories (ex. 'type A-Blue', 'type A-Black' ..) in total (4 x 2). This would be more clear from the below mentioned code.

#### groupby two columns and aggregation

```
data_prod_colour_Index = data_sales.groupby(['product','colour'], as_index=True).sum()
```

```
data_prod_colour_Index
```

# Note: as\_index=True

	sales	transactions
product	colour	
type A	Black	285000
	Blue	200000
	Red	390000
	Yellow	550000
type B	Black	350000
	Blue	702000
	Red	570000
	Yellow	90100

The above code used the aggregation function as sum(), thus we get the sum of sales and the transactions to the level of granularity defined by the combination of the 'product' and 'colour' columns.

It is to be noted that we have used the parameter 'as\_index=True', therefore we can see the 'product' and the 'colour' column as the index. On the contrary, if we take the same parameter as False then in our output we will not get the 'product' and 'colour' columns as the index but as the columns.

#### groupby without index as grouped column

```
data_prod_colour_Noindex = data_sales.groupby(['product','colour'], as_index=False).sum()
```

```
data_prod_colour_Noindex
```

	sales	transactions	
product	colour		
0	type A	Black	285000
1	type A	Blue	200000
2	type A	Red	390000
3	type A	Yellow	550000
4	type B	Black	350000
5	type B	Blue	702000
6	type B	Red	570000
7	type B	Yellow	90100

#### Custom Aggregation grouped by Multiple Columns

In previous example we used only single type of aggregation function for all the columns; however, if we want to aggregate different columns with different aggregation functions then we can use the custom aggregation functionality of the aggregation function. For doing this we can pass on the dictionary to the aggregation function stating the column name as 'key' and function name as 'value'. Interestingly, we can also pass the multiple aggregation functions to a column. Let us see an example code below for more clarity.

```
## Custom Aggregation with GroupBy using Dictionary as a parameter inside aggregation function 'agg()'
data_sales.groupby(['product','colour'], as_index=True).agg({'sales': np.sum, 'transactions':[np.median,'count']})
```

product	colour	sales	transactions	sum	median	count
		sum	median			
type A	Black	265000	135	2		
	Blue	200000	445	2		
	Red	360000	585	2		
type B	Yellow	550000	230	3		
	Black	350000	205	2		
	Blue	782000	150	3		
	Red	570000	280	3		
type C	Yellow	90100	980	1		

## 5.2.4 Pivot Tables

- Pandas :** Pandas is an open-source library that is built on top of the NumPy library. It is a Python package that offers various data structures and operations for manipulating numerical data and time series. It is mainly popular for importing and analyzing data much easier. Pandas is fast and it has high-performance & productivity for users.
- Pivot Tables :** A pivot table is a table of statistics that summarizes the data of a more extensive table (such as from a database, spreadsheet, or business intelligence program). This summary might include sums, averages, or other statistics, which the pivot table groups together in a meaningful way.
- Steps Needed**
  - Import Library (Pandas)
  - Import / Load / Create data.
  - Use Pandas.pivot\_table() method with different variants.

### Pivot tables in Excel

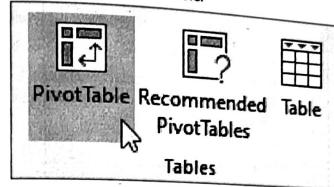
- Pivot tables are one of Excel's most powerful features. A pivot table allows you to extract the significance from a large, detailed data set.
- Our data set consists of 213 records and 6 fields. Order ID, Product, Category, Amount, Date and Country.

A	B	C	D	E	F	G	H
1	Order ID	Product	Category	Amount	Date	Country	
2	1	Carrots	Vegetables	\$4,270	1/6/2016	United States	
3	2	Broccoli	Vegetables	\$8,239	1/7/2016	United Kingdom	
4	3	Banana	Fruit	\$617	1/8/2016	United States	
5	4	Banana	Fruit	\$8,384	1/10/2016	Canada	
6	5	Beans	Vegetables	\$2,626	1/10/2016	Germany	
7	6	Orange	Fruit	\$3,610	1/11/2016	United States	
8	7	Broccoli	Vegetables	\$9,062	1/11/2016	Australia	
9	8	Banana	Fruit	\$6,906	1/16/2016	New Zealand	
10	9	Apple	Fruit	\$2,417	1/16/2016	France	
11	10	Apple	Fruit	\$7,473	1/16/2016	Canada	

### Insert a Pivot Table

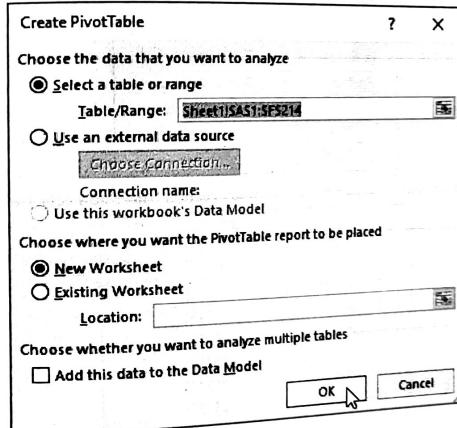
To insert a pivot table, execute the following steps.

- Click any single cell inside the data set.
- On the Insert tab, in the Tables group, click PivotTable.



The following dialog box appears. Excel automatically selects the data for you. The default location for a new pivot table is New Worksheet.

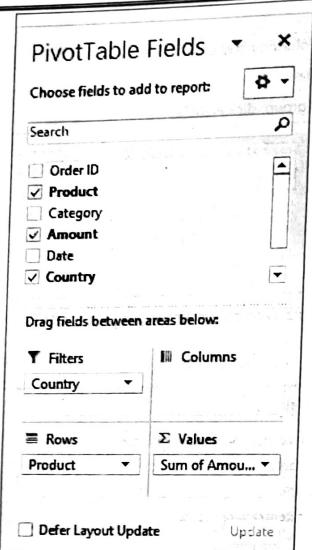
- Click OK.



### Drag fields

The PivotTable Fields pane appears. To get the total amount exported of each product, drag the following fields to the different areas.

- Product field to the Rows area.
- Amount field to the Values area.
- Country field to the Filters area.



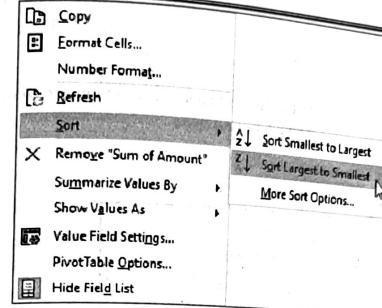
Below you can find the pivot table. Bananas are our main export product. That's how easy pivot tables can be!

A	B	C
1	Country (All)	
2		
3	Row Labels	Sum of Amount
4	Apple	191257
5	Banana	340295
6	Beans	57281
7	Broccoli	142439
8	Carrots	136945
9	Mango	57079
10	Orange	104438
11	Grand Total	1029734
12		

#### Sort

To get Banana at the top of the list, sort the pivot table.

1. Click any cell inside the Sum of Amount column.
2. Right click and click on Sort, Sort Largest to Smallest.



#### Result

A	B	C
1	Country (All)	
2		
3	Row Labels	Sum of Amount
4	Banana	340295
5	Apple	191257
6	Broccoli	142439
7	Carrots	136945
8	Orange	104438
9	Beans	57281
10	Mango	57079
11	Grand Total	1029734
12		

#### Filter

Because we added the Country field to the Filters area, we can filter this pivot table by Country. For example, which products do we export the most to France?

1. Click the filter drop-down and select France.

Result. Apples are our main export product to France.

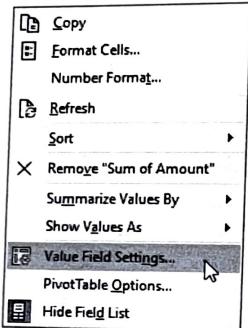
A	B	C
1	Country	France
2		
3	Row Labels	Sum of Amount
4	Apple	80193
5	Banana	36094
6	Carrots	9104
7	Mango	7388
8	Broccoli	5341
9	Orange	2256
10	Beans	680
11	Grand Total	141056
12		

**Note :** you can use the standard filter (triangle next to Row Labels) to only show the amounts of specific products.

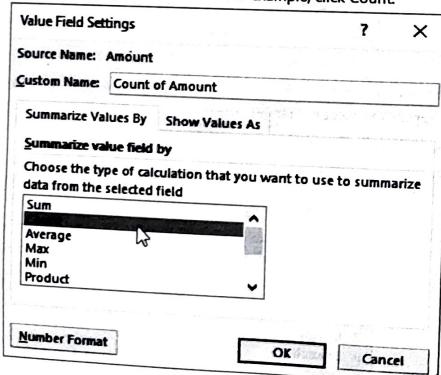
### Change Summary Calculation

By default Excel summarizes your data by either summing or counting the items. To change the type of calculation that you want to use, execute the following steps.

1. Click any cell inside the Sum of Amount column.
2. Right click and click on Value Field Settings.



3. Choose the type of calculation you want to use. For example, click Count.



4. Click OK.

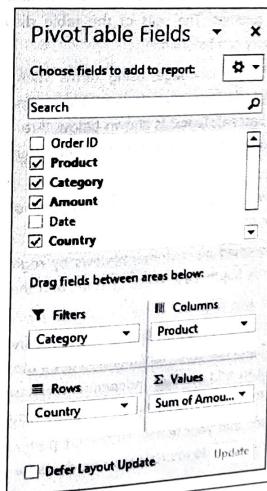
Result: 16 out of the 28 orders to France were 'Apple' orders.

A	B	C
1	Country	France
2		▼
3	Row Labels	Count of Amount
4	Apple	16
5	Banana	7
6	Carrots	1
7	Mango	1
8	Orange	1
9	Beans	1
10	Broccoli	1
11	Grand Total	28
12		

### Two-dimensional Pivot Table

If you drag a field to the Rows area and Columns area, you can create a two-dimensional pivot table. First, insert a pivot table. Next, to get the total amount exported to each country, of each product, drag the following fields to the different areas.

1. Country field to the Rows area.
2. Product field to the Columns area.
3. Amount field to the Values area.
4. Category field to the Filters area.



Below you can find the two-dimensional pivot table.

	A	B	C	D	E	F	G	H	I	J
1	Category	(All)								
2										
3	Sum of Amount	Column								
4	Row Labels	Apple	Banana	Beans	Broccoli	Carrots	Mango	Orange	Grand Total	
5	Australia	20634	52721	14433	17953	8106	9186	8680	131713	
6	Canada	24867	33775		12407		3767	19929	94745	
7	France	80193	36094	680	5341	9104	7388	2256	141056	
8	Germany	9082	39686	29905	37197	21636	8775	8887	155168	
9	New Zealand	10332	40050		4390			12010	66782	
10	United Kingdom	17534	42908	5100	38436	41815	5600	21744	173137	
11	United States	28615	95061	7163	26715	56284	22363	30932	267133	
12	Grand Total	191257	340295	57281	142439	136945	57079	104438	1029734	
13										

## 5.2.5 Cross Tabulation

- Cross tabulation, also known as a contingency table, is another technique used to analyze the relationship between two categorical variables. It presents the frequency distribution of the data for each combination of the variables, providing a clear overview of their associations.
- Cross tabulation is often visualized as a table, with one variable's categories forming the rows and the other variable's categories forming the columns. The cells of the table display the counts or percentages of observations falling into each category combination.
- Cross-tabulation analysis has its unique language, using terms such as "banners", "stubs", "Chi-Square Statistic" and "Expected Values." A typical cross-tabulation table comparing the two hypothetical variables "City of Residence" with "Favorite Baseball Team" is shown below. Are city of residence and being a fan of that team independent? The cells of the table report the frequency counts and percentages for the number of respondents in each cell.
- You typically use cross tabulation when you have categorical variables or data – e.g. information that can be divided into mutually exclusive groups.
- For example, a categorical variable could be customer reviews by region. You divide this information into reviews per geographical area: North, South, East, West, or state, and then analyze the relationships between that data.
- Another example of when to use cross-tabulation is with product surveys – you could ask a group of 50 people "Do you like our products?" and use cross-tabulation to get a more insightful answer. Rather than just recording the 50 responses, you can add another independent variable, such as gender, and use cross-tabulation to understand how the male and female respondents view your product.
- With this information, you might see that your female customers prefer your products more than your male customers. You can then use these insights to improve your products for your male customers.

- As such, using two variables – gender and product likability – along with cross-tabulation, you can get a more comprehensive breakdown of data sets to identify patterns, trends, or other useful information.
- As a result, cross-tabulation is excellent for assessing categorical variables in market research or survey responses, as you can readily compare data sets to discover the relationship between two (or more) seemingly unrelated items.

## Other uses of cross-tabulation

- Employee engagement and satisfaction** : you might want to assess how employees, both male, and female, two or more variables to get a better understanding.
- Exit interviews** : why are employees leaving and what could you do differently? Is there a correlation between employees leaving and progression, for example?
- Departmental issues** : are issues caused by one or more variables and if so, are these issues more prevalent at specific job levels?
- Evaluation surveys at schools** : how do students feel about the course material and is the time spent on the material sufficient enough?

## What are the benefits of cross-tabulation?

As a statistical analysis method that allows categorical evaluation across a data set, cross-tabulation can help to uncover variables or multiple variables that affect a specific result or can aid in improving a specific outcome.

With the examples above, you should now have a good idea of how to cross-tabulation can be used in certain contexts to glean insights. But there are several other benefits to cross-tabulation :

- Error reduction** : analyzing data sets can be confusing, let alone accurately pulling insights from them. Using cross-tabulation, you can make your data sets more manageable at scale (as they simplify them and divide them into representative subgroups).
- More Insights** : cross-tabulation looks at the relationships between one or more categorical variables to uncover more granular insights. These insights might go unnoticed with standard approaches (or require more work to reveal).
- Actionable information** : as cross-tabulation simplifies data sets and allows you to quickly compare the relationships between them, you can uncover insights faster and apply new strategies as necessary. These are the main benefits of cross-tabulation, but as a statistical analysis method, it can be applied to a wide range of research areas and disciplines to help you get more from your data.
- This technique is particularly valuable when exploring categorical data, as it helps identify patterns and dependencies between different variables.

## 5.2.6 Time Series

- Time series data refers to a sequence of data points ordered by time. Each data point represents an observation taken at a specific time interval. Time series data is common in fields such as finance, economics, climate science, and many others.
- Analyzing time series data requires specialized techniques due to its temporal nature. Here are some key aspects of time series analysis:
  - Time Series Basics :** Understanding the basic characteristics of time series data is essential for analysis. This includes identifying trends, seasonality (repeating patterns), and other temporal patterns that may exist in the data.
  - Date and Time Data Types and Tools :** Working with time series data requires using appropriate date and time data types and tools. Most programming languages and data analysis libraries provide specific data structures and functions for handling dates and times efficiently.
  - Resampling and Frequency Conversion :** Resampling involves changing the frequency of the time series data by aggregating or interpolating data to different time intervals. This operation is useful when the original data has a different granularity than the desired analysis frequency. Frequency conversion is the process of changing the data frequency to a specified frequency, such as converting daily data to monthly data.
  - Time Zone Handling :** Time zone handling is crucial when dealing with data collected from different regions with different time zones. Converting and standardizing time zones is necessary to maintain consistency and accuracy in time-based analyses.
  - Time Series Modeling :** Time series modeling aims to identify and model patterns and trends in the data to make predictions or forecasts. Techniques like autoregressive integrated moving average (ARIMA) and seasonal decomposition of time series (STL) are commonly used for time series modeling.

## 5.3 Data Analysis

Data analysis is the practice of working with data to glean useful information, which can then be used to make informed decisions.

### Data analysis process

As the data available to companies continues to grow both in amount and complexity, so too does the need for an effective and efficient process by which to harness the value of that data. The data analysis process typically moves through several iterative phases. Let's take a closer look at each.

- **Identify** the business question you'd like to answer. What problem is the company trying to solve? What do you need to measure, and how will you measure it?
- **Collect** the raw data sets you'll need to help you answer the identified question. Data collection might come from internal sources, like a company's client relationship management (CRM) software, or from secondary sources, like government records or social media application programming interfaces (APIs).

Clean the data to prepare it for analysis. This often involves purging duplicate and anomalous data, reconciling inconsistencies, standardizing data structure and format, and dealing with white spaces and other syntax errors.

Analyze the data. By manipulating the data using various data analysis techniques and tools, you can begin to find trends, correlations, outliers, and variations that tell a story. During this stage, you might use data mining to discover patterns within databases or data visualization software to help transform data into an easy-to-understand graphical format.

Interpret the results of your analysis to see how well the data answered your original question. What recommendations can you make based on the data? What are the limitations to your conclusions?

Date and time data types and tools are essential components in data analysis and programming. They provide the capability to work with temporal data efficiently, handle date-related operations, and manipulate time-based information accurately. These data types and tools are crucial for a wide range of applications, including financial analysis, event scheduling, time series analysis, and more. Let's explore date and time data types and some popular tools used for working with them:

### 5.3.1 Date and Time Data Types

- Date :** A date data type represents a specific day, typically in the format of "year-month-day" (YYYY-MM-DD). It is used to store information about calendar dates, such as birthdays, transaction dates, or event dates.
- Time :** A time data type represents a specific time of day, often in the format of "hour:minutes:second" (HH:MM:SS). It is used to store time-related information, like event timings or the duration of an activity.
- DateTime :** A DateTime data type combines both date and time information. It represents a specific point in time, including both the date and the time of day.
- Timestamp :** A timestamp data type represents the number of seconds (or milliseconds) elapsed since a predefined reference point in time, often referred to as the "epoch." Timestamps are particularly useful for precise time-based calculations and are commonly used in databases and programming languages.

### 5.3.2 Date and Time Tools

- Python - datetime module :** Python, a popular programming language, provides the "datetime" module for working with date and time data. It offers classes and functions to create, manipulate, and format date and time objects. The "datetime" module supports various functionalities like parsing dates and times from strings, performing arithmetic operations with time intervals, and converting between different formats.

```
import datetime

# Creating a date object
today = datetime.date.today()

# Creating a time object
current_time = datetime.time(12, 30, 0)

# Creating a datetime object
current_datetime = datetime.datetime.now()

# Formatting dates and times
formatted_date = today.strftime("%Y-%m-%d")
print(formatted_date) # Output: 2023-07-30
```

- b) pandas – date time functionality :** The pandas library in Python provides extensive support for handling date and time data. It introduces specialized data structures like "Timestamp" and "DatetimeIndex" for time series operations. With pandas, you can easily parse date strings, perform resampling and frequency conversion, calculate time differences, and manipulate time series data efficiently.

```
import pandas as pd

# Creating a date range
date_range = pd.date_range(start='2023-07-01', end='2023-07-31', freq='D')

# Creating a DataFrame with date index
data = {'sales': [100, 150, 200, 180, 250]}
df = pd.DataFrame(data, index=date_range)

# Resampling data to weekly frequency
weekly_sales = df.resample('W').sum()
print(weekly_sales)
```

- c) SQL - Date and Time Functions :** In SQL databases, there are specific functions to handle date and time data, such as DATE, TIME, DATETIME, TIMESTAMP, and various date-related functions like DATEADD, DATEDIFF, and DATEPART. These functions allow for filtering and aggregating data based on date and time conditions, making it easier to perform temporal queries.

```
SELECT *
FROM sales_data
WHERE DATE(order_date) = '2023-07-15';

SELECT COUNT(*)
FROM sales_data
WHERE MONTH(order_date) = 7 AND YEAR(order_date) = 2023;
```

### 5.3.3 Time Series Basic

Time series basics refer to the fundamental concepts and characteristics of time series data. A time series is a sequence of data points collected over successive time intervals, where each data point corresponds to a specific time. Time series data is commonly encountered in various fields, such as finance, economics, climate science, and more. Understanding the basics of time series data is crucial for effective analysis, modeling, and forecasting. Let's explore the key aspects of time series basics :

- Temporal Ordering :** The most critical characteristic of time series data is its temporal ordering. Data points in a time series are arranged chronologically, with each observation representing a specific point in time. The time intervals between data points are typically regular (e.g., hourly, daily, monthly) or irregular, depending on the application.

- Time Domain and Value Domain :** In time series data, the "Time Domain" refers to the time axis, which represents the time intervals at which data points were collected. The "Value Domain" represents the measured or observed values associated with each time point. The value domain could be numerical (e.g., temperature, sales, stock prices) or categorical (e.g., event types, weather conditions).
- Trend :** The trend in a time series refers to the long-term pattern or movement in the data over time. Trends can be upward (increasing), downward (decreasing), or relatively constant. Understanding the trend is essential for identifying the overall direction or behavior of the data.
- Seasonality :** Seasonality in a time series refers to recurring patterns or fluctuations that occur at regular intervals, often related to calendar seasons, months, days of the week, or other periodic events. Seasonality can significantly influence the data, and detecting it is crucial for proper modeling and forecasting.
- Cyclical Patterns :** Cyclical patterns are longer-term fluctuations in a time series that do not have a fixed frequency. Unlike seasonality, which occurs at regular intervals, cyclical patterns are more irregular and can be attributed to economic cycles or other non-calendar-based factors.
- Stationarity :** Stationarity is an important concept in time series analysis. A time series is considered "strictly stationary" if its statistical properties, such as mean, variance, and autocorrelation, remain constant over time. Stationarity simplifies modeling and facilitates accurate forecasts. If a time series is non-stationary, transformations or differencing techniques may be applied to achieve stationarity.
- Autocorrelation :** Autocorrelation measures the correlation between a time series and its lagged versions. In other words, it assesses the relationship between a data point and its previous observations. Autocorrelation is crucial for identifying patterns and dependencies in the time series, especially in cases where historical data points impact future values.
- White Noise :** White noise is a special type of time series with no discernible patterns, trends, or seasonality. It exhibits random behavior, and each data point is independent of other data points. White noise serves as a benchmark to evaluate the presence of meaningful patterns in a time series.

#### 5.3.4 Data Ranges

Date ranges refer to a sequence of dates within a specified period, forming a continuous interval of time. Date ranges are commonly used in data analysis, time series modeling, and other applications where it is necessary to work with a specific time frame or period. Understanding date ranges is crucial for conducting various temporal analyses and aggregations. Let's explore date ranges in more detail :

- Defining a Date Range :** A date range is typically defined by two dates: a start date and an end date. The start date represents the beginning of the range, while the end date marks its conclusion. The range includes all dates between the start and end dates, inclusive of both endpoints.
- Regular and Irregular Date Ranges :** Date ranges can be regular or irregular, depending on the time intervals between consecutive dates. A regular date range has constant intervals between each date, such as daily, weekly, or monthly intervals. For example, a daily date range from January 1 to January 7 would include seven consecutive days.

In contrast, an irregular date range may have varying intervals between dates. For example, an irregular date range may include dates representing significant events or occurrences without adhering to a fixed time interval.

- 3. Generating Date Ranges :** Generating date ranges can be done programmatically in various programming languages and libraries. Many programming languages, such as Python and R, provide built-in functions or libraries for working with date ranges. For instance, in Python, the pandas library's "date\_range" function allows users to create regular date ranges with specified start and end dates and a desired frequency.

```
import pandas as pd

# Create a daily date range for a week starting from 2023-07-01
date_range = pd.date_range(start='2023-07-01', periods=7, freq='D')
print(date_range)
```

- 4. Using Date Ranges in Data Analysis :** Date ranges are valuable in data analysis to filter, aggregate, and analyze data within specific time periods. They allow analysts to focus on a particular range of dates and perform calculations or visualizations accordingly.

For example, a sales dataset may contain transactions from various dates. By using a date range from January 1 to January 31, analysts can compute total sales for that specific month, enabling comparisons with other months or identifying monthly trends.

- 5. Time Series Analysis :** In time series analysis, date ranges play a crucial role in defining the time frame for which historical data is available. Analyzing data within a specific date range helps in understanding patterns and trends over that period. Additionally, using date ranges is essential for conducting time-based forecasts and predictions.

- 6. Handling Missing Data :** Date ranges can also be used to identify missing data within a time series. By generating a complete date range with regular intervals, analysts can identify gaps in the data where observations are missing, allowing for proper handling or imputation of missing values.

### 5.3.5 Frequencies and Shifting

- 1. Frequencies :** In time series analysis, frequency refers to the time interval between consecutive data points. It determines how often data is collected or observed. Frequencies can be regular or irregular, depending on the time intervals between data points.

- a) **Regular Frequencies :** Regular frequencies have a constant time interval between consecutive data points. Common examples include daily, weekly, monthly, and annual frequencies. Regular frequencies are often used in financial data, economic indicators, and other regularly sampled time series.

- b) **Irregular Frequencies :** Irregular frequencies have varying time intervals between consecutive data points. This occurs when data is collected at irregular time intervals, such as in event-based data or when data is missing for certain time periods.

- 2. Shifting :** Shifting, also known as lagging or time shifting, involves moving data points forward or backward in time. It is a common operation in time series analysis and is used for various purposes, such as comparing time series at different time lags or aligning time series for temporal analysis.

- a) **Forward Shifting :** Forward shifting involves moving data points to later time periods. For example, shifting a time series forward by one day would mean that each data point is replaced by the data point that occurs one day later in the time series.

Forward shifting is useful for conducting time-based comparisons, such as comparing current data with data from the previous day or week.

- b) **Backward Shifting :** Backward shifting involves moving data points to earlier time periods. For example, shifting a time series backward by one day would mean that each data point is replaced by the data point that occurred one day earlier in the time series.

Backward shifting is often used for lagged analysis, such as studying the impact of past events on current or future data.

### 3. Applications of Frequencies and Shifting

- a) **Resampling and Frequency Conversion :** Frequencies are essential in resampling time series data. Resampling involves changing the frequency of the time series by aggregating or interpolating data to different time intervals. For example, converting daily data to monthly data involves resampling at a different frequency.

- b) **Aligning Time Series :** Shifting is commonly used to align multiple time series for comparison and analysis. By aligning time series at the same time points, analysts can perform meaningful calculations, such as calculating differences or correlations between related time series.

- c) **Time Series Modeling :** In time series modeling, shifting can be used to create lagged variables, which are often useful as input features for predictive models. Lagged variables represent past observations and can help capture temporal dependencies in the data.

### 5.3.6 Time Zone Handling

Time zone handling is a critical aspect of working with time and date data, especially when dealing with data collected from different regions or time zones. Time zone handling ensures that temporal data is accurately represented and converted across different time zones, enabling consistent analysis and proper interpretation of time-based information.

#### 1. Time Zones and UTC

- Time zones are geographical regions that have standardized offsets from Coordinated Universal Time (UTC), also known as Greenwich Mean Time (GMT). Each time zone has a specific offset, represented in hours and minutes, indicating the difference between local time and UTC.
- For example, Eastern Standard Time (EST) in the United States is UTC-5, meaning it is five hours behind UTC. Conversely, Central European Time (CET) is UTC+1, meaning it is one hour ahead of UTC.

#### 2. Time Zone Awareness

- When working with time and date data, it's crucial to be time zone aware. Time zone awareness means explicitly recording the time zone information associated with each timestamp, ensuring that the data accurately reflects the local time at the point of observation.
- Time zone awareness is essential for applications where the analysis depends on the temporal relationships between events across different time zones or when data is collected from multiple regions.

### 3. UTC as a Reference

- To handle time zones effectively, it is common to convert all temporal data to and from UTC. UTC serves as a reference time standard, providing a common reference point for converting and comparing time across different time zones.
- By converting data to UTC, temporal data from various time zones can be standardized, enabling easy comparisons and consistent analysis.

### 4. Time Zone Conversions

- Time zone conversions involve converting time and date data between different time zones. This is necessary when displaying or analyzing data collected in one time zone for users or applications in another time zone.
- For example, a company with offices in multiple countries may need to convert timestamps recorded in each office's local time to a unified time zone for central reporting and analysis.

### 5. Daylight Saving Time (DST)

- Daylight Saving Time (DST) is a practice in some regions where clocks are adjusted forward or backward by one hour during certain periods of the year to make better use of natural daylight. DST introduces complexities in time zone handling, as some time zones have DST transitions while others do not.
- To handle DST transitions accurately, it's essential to use time zone libraries or services that take into account the historical and future DST changes for each time zone.

### 6. Time Zone Libraries

- Most programming languages provide built-in or external time zone libraries to handle time zone conversions and awareness. For example, Python has the "pytz" library, while JavaScript has the "moment-timezone" library.
- These libraries offer functions and utilities for converting between time zones, dealing with DST transitions, and ensuring time zone-awareness in applications.

#### 5.3.7 Periods and Periods Arithmetic

In time series data analysis, "periods" refer to fixed-length time intervals, such as days, months, quarters, or years, that are used to organize and aggregate temporal data. Periods are distinct from time stamps, which represent specific points in time. Periods are particularly useful for resampling, frequency conversion, and performing time-based calculations. Periods arithmetic involves operations that manipulate and interact with these fixed-length time intervals. Let's explore periods and periods arithmetic in more detail:

##### 1. Periods

- In the context of time series, periods represent a time interval of fixed length. For example, a period of "one month" represents a time interval that spans from the first day of a month to the last day of that month. Similarly, a period of "one year" represents a time interval that covers an entire calendar year.
- Periods are useful for organizing time series data into regular intervals, which simplifies various time-based analyses, such as calculating monthly averages, identifying seasonality, or comparing data between specific time periods.

### 2. Periods Arithmetic

Periods arithmetic involves performing mathematical operations on periods to obtain new periods or derive meaningful insights from temporal data. Some common operations in periods arithmetic include:

- Addition and Subtraction :** Periods can be added or subtracted to move forward or backward in time. For example, adding one month to a given period would shift it to the next month, while subtracting one month would move it to the previous month.
- Difference :** By taking the difference between two periods, we can calculate the duration or time span between them. For example, the difference between two quarterly periods can determine the number of quarters between those periods.
- Multiplication and Division :** Periods can also be multiplied or divided by a scalar value. For instance, multiplying a monthly period by 3 would result in a quarterly period, while dividing a yearly period by 4 would yield a quarterly period.

#### 3. Periods Arithmetic Example

```
import pandas as pd

# Create periods representing the first day of each month from January
months_period = pd.period_range(start='2023-01-01', end='2023-12-01')

# Add 3 months to the first period to get the 4th month
fourth_month_period = months_period[0] + 3
print(fourth_month_period) # Output: Period('2023-04', 'M')

# Calculate the difference between the 6th and 2nd periods
period_difference = months_period[5] - months_period[1]
print(period_difference) # Output: <5 * MonthEnds>

# Multiply the 3rd period by 4 to get the 12th period (one year)
twelfth_month_period = months_period[2] * 4
print(twelfth_month_period) # Output: Period('2023', 'A-DEC')
```

In the context of mathematics and numerical operations, "periods" typically refer to different concepts depending on the specific field of mathematics or context in which they are used. Here are some common interpretations of "periods" and "periods arithmetic":

##### 1. Periods in Trigonometry and Geometry

- In trigonometry, a "period" refers to the interval over which a trigonometric function, such as sine or cosine, repeats itself. For example, the sine function has a period of  $2\pi$  radians, which means it repeats its values every  $2\pi$  radians.
- In geometry, "periodic tessellation" refers to a repeating pattern of tiles or shapes that cover a plane without any gaps or overlaps. These tessellations often have periodic characteristics.

##### 2. Periods in Time and Frequency

In the context of time and frequency analysis, a "period" usually refers to the time it takes for a repeating event or wave to complete one full cycle. For example, the period of a simple harmonic motion can be the time it takes for a pendulum to swing back and forth once.

### 3. Periods in Financial Mathematics:

- In finance, "periods" often refer to discrete time intervals, such as months, quarters, or years. Financial calculations often involve compound interest, where the number of periods plays a crucial role in determining the final amount.
- Now, let's briefly discuss "periods arithmetic," which might refer to mathematical operations or calculations involving periods:

### 4. Periodic Functions Arithmetic

When working with periodic functions like sine and cosine, you can perform arithmetic operations involving these functions. For example, you can add, subtract, multiply, or divide periodic functions while considering their periods.

### 5. Time-Series Arithmetic:

In time-series analysis, arithmetic operations can be performed on data points collected at regular time intervals (periods). These operations might include calculating averages, growth rates, or changes in values over specific periods of time.

### 6. Financial Periods Arithmetic:

- In finance, calculations often involve periods, such as compounding interest over a specific number of periods, calculating the net present value of cash flows occurring at different periods, or determining the future value of investments over multiple periods.
- To perform arithmetic operations with periods, it's important to understand the specific context and units involved, as different fields and situations may have their own conventions and formulas for working with periods.

## 5.3.8 Resampling and Frequency Conversion

Resampling and frequency conversion are important techniques used in time series data analysis to change the time intervals at which data is recorded or observed. These operations are useful for adjusting the granularity of time series data, aggregating data at different time intervals, and preparing data for specific time-based analyses. Let's explore resampling and frequency conversion in more detail :

**1. Resampling :** Resampling involves changing the frequency of a time series by aggregating or interpolating data to different time intervals. Resampling is often necessary when the original data is collected at a different frequency than the desired analysis frequency. There are two primary types of resampling :

**a) Downsampling :** Downsampling involves reducing the frequency of the data, converting it to a lower time resolution. This is typically done to summarize data over longer time periods. For example, converting daily data to weekly or monthly data involves downsampling.

In downsampling, data is aggregated or combined within each new time interval. Common aggregation functions include sum, mean, median, or other statistical measures.

**b) Upsampling :** Upsampling involves increasing the frequency of the data, converting it to a higher time resolution. This is often done to add more data points within a given time period. For example, converting monthly data to daily data involves upsampling.

In upsampling, new data points are interpolated or filled in between existing data points. Interpolation methods may include linear interpolation, polynomial interpolation, or other techniques depending on the nature of the data.

### 2. Frequency Conversion :

- Frequency conversion refers to changing the time intervals of a time series to a new frequency, which can be either higher or lower than the original frequency. Frequency conversion is a combination of both downsampling and upsampling operations.
- For example, converting daily data to hourly data involves frequency conversion by upsampling, where new data points are inserted between the existing daily data points.
- Frequency conversion is useful when analyzing data at a different time granularity, or when aligning multiple time series with different frequencies for comparison.

### 3. Resampling and Frequency Conversion Example :

Let's illustrate resampling and frequency conversion using Python with the pandas library :

```
python
```

```
import pandas as pd
```

```
# Create a daily time series with random data
```

```
dates = pd.date_range(start='2023-01-01', periods=10, freq='D')
daily_data = pd.Series(range(10), index=dates)
```

```
# Downsample to weekly data by taking the sum of each week
```

```
weekly_data = daily_data.resample('W').sum()
print(weekly_data)
```

```
# Upsample to hourly data using linear interpolation
```

```
hourly_dates = pd.date_range(start='2023-01-01', end='2023-01-10', freq='H')
hourly_data = daily_data.resample('H').interpolate(method='linear')
print(hourly_data)
"
```

In this example, we first create a daily time series with random data. We then downsample the data to weekly frequency by taking the sum of data points within each week. Next, we upsample the data to hourly frequency using linear interpolation to fill in new data points between the existing daily data points.

## 5.3.9 Moving Window Functions

Moving window functions, also known as rolling or sliding window functions, are essential tools in time series data analysis. They involve applying a specific function or operation to a fixed-size window of data that moves along the time axis. These functions are used to calculate rolling statistics, smooth data, identify trends, and perform other time-based computations. Moving window functions are particularly valuable when dealing with noisy or volatile time series data. Let's explore moving window functions in more detail :

- Basic Concept :** The basic concept of moving window functions involves defining a fixed-size window that spans a specific number of consecutive data points in the time series. The window moves one step at a time along the time axis, continuously updating the window's data points and recalculating the function's output.
- Applications :** Moving window functions have various applications in time series data analysis, including :
  - Rolling Statistics :** Calculating rolling statistics such as rolling mean, rolling median, rolling standard deviation, or other aggregated measures within the moving window. These statistics help smooth data and identify underlying patterns and trends.
  - Moving Averages :** Computing moving averages by taking the mean of the data within the window. Moving averages are commonly used to reveal underlying trends or eliminate noise in the data.
  - Exponential Moving Averages (EMA) :** Similar to moving averages, EMA assigns exponentially decreasing weights to data points within the window, giving more weight to recent observations. EMA is widely used in financial analysis and trend detection.
  - Rolling Sum :** Calculating the sum of data points within the window to analyze cumulative values or trends.
- Window Size :** The size of the moving window is an essential parameter in moving window functions. The window size determines the number of data points included in the computation at any given time. A larger window size results in a smoother output but may lead to slower responsiveness to changes in the data. On the other hand, a smaller window size provides a more responsive output but may be more sensitive to noise.
- Handling Boundary Effects :** When applying moving window functions, boundary effects must be considered. At the beginning and end of the time series, there may not be enough data points to form a complete window. Different strategies can be used to handle this, such as padding missing values or using weighted windows that give more weight to available data points.

In addition to the above two boundary effects of electrostatic and hydrodynamic nature respectively, the presence of the boundary also has a geometric confining effect upon the ion distribution around the particle if the double layer is thick enough to reach the boundary. This affects the electric driving force once an external electric field is applied upon the system. This perhaps is the most important issue in the study of boundary effect in electrokinetic motion, as the double layer polarization/deformation determines the ultimate particle motion in general. And obviously the thicker the double layer is, the more significant this factor is, as the deformation of the double layer can be much more profound.

### 5.3.9(A) Moving Window Functions Example

- Time series data is a series of data points recorded with a time component (temporal) present. Majority of the time these data points are recorded at a fixed time interval.
- Many real-world datasets like stock market data, weather data, geography datasets, earthquake datasets, etc are time series datasets.
- While working with time series datasets, we need to perform various operations on them to analyze datasets from different perspectives. The two most common operations are resampling and moving window functions.

- Time series Resampling is the process of changing frequency at which data points(observations) are recorded. Resampling is generally performed to analyze how time series data behaves under different frequencies.
- Moving window functions are aggregate functions applied to time series datasets by moving window of fixed /variable size through them. Moving window functions can be used to smooth time series to handle noise.

Let's illustrate moving window functions using Python with the pandas library :

```
import pandas as pd

# Create a time series with random data
dates = pd.date_range(start='2023-01-01', periods=10, freq='D')
data = pd.Series([10, 15, 20, 25, 30, 35, 40, 45, 50, 55], index=dates)

# Calculate the rolling mean with a window size of 3
rolling_mean = data.rolling(window=3).mean()
print(rolling_mean)
```

- In this example, we create a time series with random data and calculate the rolling mean using a window size of 3. The rolling mean smooths the data by taking the average of every three consecutive data points.

### Review Questions

- What is mean by group by mechanics?
- Explain data aggregation.
- Discuss various uses of data aggregation.
- Explain three steps of the split-apply-combine strategy.
- Explain various date and time tools.
- What is time series of data?
- Discuss about data ranges.



# 6

# Data Analysis of Visualization and Modeling

## Syllabus

Reconstruction, Visualization and Analysis of Medical Images

**Introduction :** PET Images, Ultrasound Images, Magnetic Resonance Images, Conclusion and Discussion,

**Case Study :** ER/Studio, Erwin data modeler, DbSchema Pro, Archi, SQL Database Modeler, LucidChart, Pgmodeler.

## 6.1 Reconstruction

### Reconstruction and Its Types in Data Modelling and Visualization

Data modeling and visualization are critical components of the data analysis process. They involve creating representations of data to aid in understanding patterns, relationships, and trends. In certain scenarios, data may need to be reconstructed or transformed to better suit the analysis or visualization process.

Reconstruction refers to the process of transforming and reshaping data to derive meaningful insights and present information in a more understandable manner. In this note, we will explore the concept of reconstruction and its various types in the context of data modeling and visualization, along with relevant examples.

#### 1. Data Reconstruction :

Data reconstruction involves the manipulation and modification of raw data to enhance its usability for analysis and visualization purposes. This process may entail aggregating, filtering, transforming, or cleaning data to uncover valuable insights that might not be immediately apparent from the original dataset.

Reconstruction in data modeling and visualization refers to the process of generating or creating representations of data, often in a more structured or informative manner. It involves transforming raw data into visualizations or models that are easier to interpret, understand, and analyze. Here are some examples of reconstruction in data modeling and visualization :

(i) **3D Reconstruction from Imaging Data :** In fields like medical imaging or computer vision, 3D reconstruction involves taking a series of 2D images (such as MRI or CT scans) and reconstructing a 3D representation of the object or scene. This can be used for creating detailed anatomical models for medical diagnoses, archaeological reconstructions, or even virtual reality environments.

(ii) **Financial Data Visualization :** Financial data often consists of large, complex datasets. Reconstruction in this context could involve transforming raw financial data into interactive dashboards, charts, and graphs that provide insights into trends, patterns, and anomalies. For instance, a company might reconstruct their financial data into visualizations to analyze revenue trends over time.

- (iii) **Molecular Visualization in Chemistry :** Chemists might use data from various analytical techniques to reconstruct molecular structures. This could involve generating 3D molecular models based on spectroscopic data, crystallography, or computational simulations. These visualizations aid in understanding molecular behavior and interactions.
- (iv) **GIS and Terrain Reconstruction :** Geographic Information Systems (GIS) are used to model and visualize geographical data. Terrain reconstruction involves transforming elevation data from sources like satellite imagery into 3D terrain models, which can be used for urban planning, environmental assessment, or simulation purposes.
- (v) **Biological Systems Modeling :** In systems biology, researchers might reconstruct intricate models of biological pathways or cellular processes. These models can include molecular interactions, protein networks, and gene expression patterns, helping scientists understand how various components of a biological system work together.
- (vi) **Historical Data Reconstruction :** Historians might reconstruct historical events or scenarios based on fragmented data sources. For instance, a historian might piece together historical records, artifacts, and accounts to reconstruct a detailed timeline of a past event.
- (vii) **Meteorological Data Visualization :** Weather data, such as temperature, humidity, and wind speed, can be transformed into dynamic visualizations like weather maps, heatmaps, and animations. These reconstructions help meteorologists and the public understand and predict weather patterns.
- (viii) **Social Network Analysis :** Social networks generate vast amounts of data on connections and interactions between individuals. Researchers can reconstruct visualizations of these networks, revealing insights into information flow, community structure, and influence dynamics.
- (ix) **CAD Model Rendering :** Computer-Aided Design (CAD) models are often reconstructed into realistic renderings for architectural visualization, product design, and engineering simulations. These visualizations help designers and engineers understand how a design will look and function in the real world.
- (x) **Astrophysical Simulations :** Astronomers use complex simulations to reconstruct the behavior of celestial bodies and cosmic phenomena. These simulations generate visualizations that aid in understanding galaxy formation, black hole dynamics, and more.

In all these examples, reconstruction involves converting raw data into meaningful visual or structural representations that enhance understanding and facilitate analysis. The process often requires a combination of domain knowledge, data manipulation, and visualization techniques.

#### 2. Types of Reconstruction

- a. **Temporal Reconstruction :** Temporal reconstruction involves organizing data over time. For instance, it could be transforming daily sales data into monthly or yearly summaries to identify seasonal trends more effectively.
- b. **Spatial Reconstruction :** Spatial reconstruction involves representing data in various geographic locations. Geospatial data, such as GPS coordinates, can be transformed into visual maps to identify regional patterns or to display data on interactive dashboards.

- c. **Hierarchical Reconstruction** : Hierarchical reconstruction involves arranging data into hierarchical structures. For instance, data representing organizational departments can be organized into a tree-like structure to facilitate better understanding and analysis.
- d. **Aggregated Reconstruction** : Aggregated reconstruction involves summarizing data by aggregating information from multiple records. It is useful for obtaining overall trends and insights without losing key details.
- e. **Data Cleansing and Imputation** : Data reconstruction may also involve data cleansing and imputation techniques to handle missing or erroneous values. For example, missing values in a dataset can be replaced with appropriate estimations or averages to maintain data integrity.
- f. **Dimensionality Reduction** : Dimensionality reduction is a type of reconstruction that involves transforming high-dimensional data into a lower-dimensional representation. Principal Component Analysis (PCA) is a common technique used to achieve this goal while preserving the most critical information.
- g. **Data Normalization and Standardization** : Data normalization and standardization are methods used to scale and transform data into a common range, making it easier to compare and visualize variables with different units and scales.
- h. **Time-Series Forecasting** : Reconstructing time-series data can include forecasting future values based on historical patterns, enabling businesses to make informed decisions and anticipate future trends.
- i. **Feature Engineering** : Feature engineering involves creating new features or variables based on existing data to improve the performance of predictive models or enhance data visualization.
- j. **Joining and Merging Data** : Data reconstruction often involves combining multiple datasets through joins or merges to create a comprehensive dataset that encompasses all relevant information.
- k. **Unpivoting Data** : Unpivoting is a technique used to transform data from a wide format (multiple columns) to a long format (a single column), making it easier to perform operations like filtering and aggregation.
- l. **Data Interpolation** : Data interpolation involves estimating unknown data points between known data points based on existing values, ensuring a smoother representation of data.

### 3. Examples of Reconstruction in Data Modeling and Visualization

- a. Consider a retail company that collects daily sales data. By using temporal reconstruction, this data can be transformed into weekly, monthly, or yearly summaries to identify seasonal sales trends, peak selling periods, and overall performance.
- b. In geospatial analysis, spatial reconstruction can be applied to GPS data collected from mobile devices. This data can be visualized on a map to determine hotspots of user activity, helping businesses optimize their marketing strategies.
- c. In an organizational setting, hierarchical reconstruction can be employed to create an organizational chart based on employee roles and departments. This helps visualize the reporting structure and understand communication pathways.
- d. Aggregated reconstruction can be used in customer data analysis to identify customer segments with common traits or preferences, which is valuable for targeted marketing campaigns.

- e. In data preprocessing, data cleansing and imputation techniques can be employed to handle missing or erroneous data points, ensuring the integrity and accuracy of the final analysis.
- f. Dimensionality reduction techniques, such as PCA, can be used in visualizations to convert high-dimensional data, like customer behavior data, into a lower-dimensional representation that preserves essential patterns and relationships.
- g. By applying time-series forecasting to historical stock market data, investors can make informed decisions and predict potential market trends.
- h. Feature engineering can be applied to social media data, where new features can be generated, like sentiment scores, to understand customer feedback and improve sentiment analysis.
- i. Data reconstruction through joining and merging can be seen in healthcare, where patient data from different sources, such as electronic health records and lab reports, are combined for comprehensive analysis.
- j. Unpivoting data can be useful when dealing with survey responses where each question's multiple answers are transformed into individual records for better analysis.
- k. Data interpolation can be applied to weather data to estimate missing temperature readings between known data points, providing a more continuous representation of weather conditions.

In conclusion, data reconstruction plays a pivotal role in data modeling and visualization. It involves transforming raw data into more usable and informative formats, enabling better analysis and visualization. By employing various types of reconstruction techniques, analysts can extract meaningful insights from complex datasets and present them in a comprehensible manner to support decision-making processes and drive business growth.

## 6.2 Visualization and Analysis of Medical Images

### Definition of Visualization in Data Modelling

Visualization is a powerful tool in data modeling that allows analysts and stakeholders to gain valuable insights from complex datasets by presenting information in a visual and intuitive manner. Data modeling involves creating a representation of data that captures the relationships, patterns, and structures within the dataset. Visualization complements data modeling by providing a visual representation of the modeled data, making it easier to understand, analyze, and communicate findings. In this note, we will explore the importance of visualization in data modeling and its key benefits.

- 1. **Enhanced Data Understanding** : Visualization aids in understanding the underlying patterns and trends within the data. By representing data points graphically, analysts can quickly grasp the distribution, central tendencies, and outliers, enabling them to make data-driven decisions with greater confidence.
- 2. **Identification of Patterns and Anomalies** : Visualizations help reveal hidden patterns and anomalies that may not be immediately apparent from the raw data. Scatter plots, line charts, and histograms can uncover relationships and trends that might be missed when examining numerical data alone.
- 3. **Simplified Communication** : Complex data models can be challenging to explain to non-technical stakeholders. Visualization simplifies the communication of complex ideas and findings, making it easier for decision-makers to understand and act upon the insights derived from the data.

- 4. Exploratory Data Analysis :** Visualization is an essential tool in exploratory data analysis (EDA). It allows analysts to interactively explore the data, drill down into specific aspects, and gain a deeper understanding of the dataset's characteristics and relationships.
- 5. Decision Support :** Data modeling aims to aid decision-making processes, and visualization plays a key role in this regard. By presenting data in a visually compelling manner, decision-makers can quickly grasp the implications of different options and choose the most appropriate course of action.
- 6. Storytelling with Data :** Visualization enables analysts to tell compelling stories with data. By creating interactive dashboards or infographics, analysts can present data insights in a narrative format, making it easier for audiences to engage with and understand the information.
- 7. Validation of Models :** Visualizing data models allows analysts to validate the accuracy and effectiveness of the models. Visualization can reveal whether the model aligns with the actual data and can help identify areas for improvement or refinement.
- 8. Insight Into Data Relationships :** Data modeling involves creating relationships between different variables. Visualization makes these relationships tangible and visible, helping analysts and stakeholders understand how changes in one variable affect others.
- 9. Data-driven Decision Making :** Effective visualization empowers data-driven decision-making processes. When data is presented in a visually appealing and understandable way, decision-makers are more likely to base their choices on data-driven insights, leading to better outcomes for organizations.

In conclusion, visualization is a crucial component of data modeling that enhances data understanding, identifies patterns and anomalies, simplifies communication, supports exploratory data analysis, aids decision-making, and allows for storytelling with data. By using visualization in data modeling, analysts can leverage the power of visuals to derive valuable insights, communicate complex ideas, and make informed decisions based on the analyzed data.

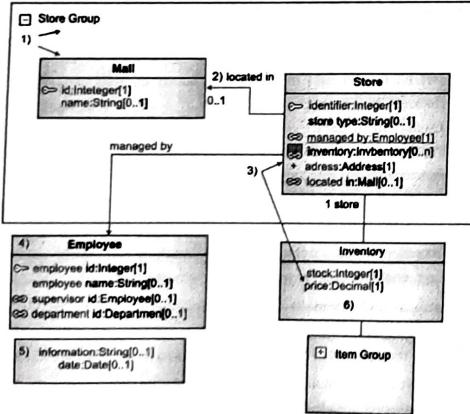


Fig. 6.2.1 : Data-driven decision-making processes

- The analysis of medical images plays a pivotal role in modern healthcare, facilitating accurate diagnosis, treatment planning, and monitoring of various medical conditions. Data modeling and visualization are instrumental in extracting meaningful insights from these images and presenting them in an easily interpretable manner.
  - Data modeling in medical image analysis involves creating sophisticated algorithms and models to process and interpret the complex image data. Image segmentation is a crucial aspect of data modeling, where identifying and quantifying abnormalities, allowing for precise diagnosis and treatment recommendations. Additionally, data modeling techniques like image registration align multiple images from different modalities or time points, facilitating a comprehensive assessment of changes over time or fusion of information for a more accurate diagnosis.
  - Visualization plays a vital role in presenting the results of medical image analysis to healthcare professionals and patients. Techniques like volume rendering, surface rendering, and multi-planar reformatting enable the visualization of 3D medical image data, allowing for a detailed exploration of anatomical structures and pathological changes. Visualization enhances understanding and communication, making it easier for radiologists, surgeons, and other medical practitioners to collaborate and make informed decisions.
  - The combination of data modeling and visualization allows for quantitative measurements of anatomical structures and pathological changes. These measurements are critical in assessing disease progression, planning interventions, and monitoring treatment outcomes. Image-based modeling, derived from the analysis of medical images, enables the creation of patient-specific anatomical models, providing a foundation for simulation and personalized treatment planning.
  - Moreover, the application of data modeling and visualization in computer-aided diagnosis (CAD) systems has revolutionized medical imaging. CAD systems utilize advanced algorithms to assist radiologists and clinicians in detecting and diagnosing medical conditions accurately. These systems act as a valuable second opinion, reducing human error and improving the overall diagnostic accuracy.
  - Medical imaging encompasses various techniques such as X-ray, Computed Tomography (CT), Magnetic Resonance Imaging (MRI), Ultrasound, Positron Emission Tomography (PET), and Single-Photon Emission Computed Tomography (SPECT). Each technique provides unique information about the internal structures of the body, and their combination allows for a comprehensive assessment of medical conditions.
- 1. Image Segmentation :** One of the essential tasks in medical image analysis is image segmentation. Segmentation involves dividing the image into meaningful regions or structures, which is crucial for identifying and quantifying abnormalities accurately. Techniques like thresholding, region-growing, and deep learning-based segmentation algorithms are commonly used in this process.
  - 2. Disease Detection and Diagnosis :** Medical image analysis helps detect and diagnose various diseases and conditions. For instance, mammograms are used for breast cancer screening, while CT and MRI scans aid in the diagnosis of neurological disorders and musculoskeletal conditions.

- 3. Image Registration :** Image registration is the process of aligning multiple medical images from different modalities or time points. It allows for better visualization of changes over time or fusion of information from different imaging modalities, enabling comprehensive diagnosis and treatment planning.
- 4. Quantitative Measurements :** Medical image analysis allows for precise quantitative measurements of anatomical structures and pathological changes. For example, cardiac MRI can provide accurate measurements of heart chamber volumes and ejection fraction, essential in assessing cardiac function.
- 5. Image-Based Modeling :** Medical image data can be used to create patient-specific anatomical models for simulation and treatment planning. These models find applications in surgical planning, implant design, and personalized medicine.
- 6. Computer-Aided Diagnosis (CAD) :** CAD systems utilize medical image analysis algorithms to assist radiologists and clinicians in making diagnoses. These systems can highlight potential abnormalities, assisting in early detection and reducing human error.
- 7. Visualization Techniques :** Visualization is critical in medical image analysis to present complex information in an easily interpretable manner. Techniques such as volume rendering, surface rendering, and multi-planar reformatting aid in visualizing 3D medical image data.
- 8. Image Enhancement :** Image enhancement techniques are employed to improve the quality and visibility of medical images. These techniques aim to reduce noise, enhance contrast, and highlight specific features, thus aiding in more accurate diagnosis.
- 9. Feature Extraction :** Feature extraction techniques are used to identify relevant patterns and characteristics in medical images. These extracted features can be fed into machine learning algorithms for disease classification and prediction.
- 10. Data Integration :** Medical image analysis often involves integrating imaging data with other clinical data, such as patient demographics, laboratory results, and genomic information. This integration allows for a comprehensive understanding of patient health and disease progression.
- 11. Data Mining and Research :** Large repositories of medical images are valuable resources for data mining and research purposes. By analyzing these datasets, researchers can identify population trends, develop new diagnostic tools, and gain insights into disease mechanisms.
- In conclusion, the analysis of medical images is an integral part of modern healthcare. By employing various techniques, such as image segmentation, disease detection, image registration, and quantitative measurements, medical image analysis provides valuable information for diagnosis, treatment planning, and research.
  - Visualization techniques play a crucial role in presenting the results of image analysis in a manner that is easily understandable to healthcare professionals and patients. Additionally, the integration of medical image data with other clinical information and its application in data modeling and machine learning further enhances our ability to deliver personalized and precise healthcare solutions.
  - Random variables are a fundamental concept in statistics and probability theory. Discrete random variables take on countable, specific values, while continuous random variables assume uncountably infinite values. Understanding the properties of both types is crucial in many statistical applications.

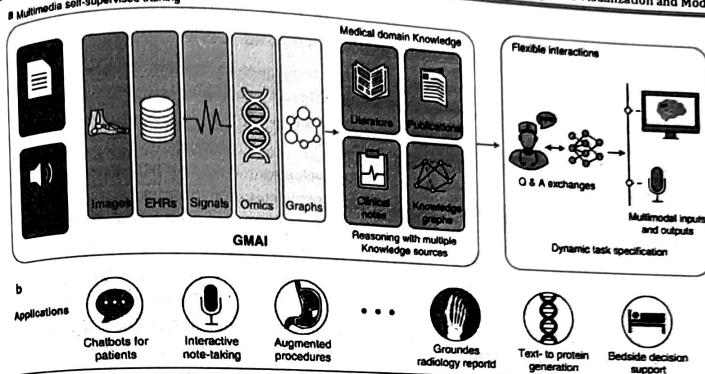


Fig. 6.2.2

### 6.3 Introduction : PET Images

- Positron Emission Tomography (PET) is a powerful medical imaging technique used to visualize metabolic processes and functional changes in tissues at a molecular level. PET images provide valuable information about the activity of specific biomolecules, such as glucose metabolism or neurotransmitter receptor binding, and are extensively used in various applications in data modeling and visualization in the field of medicine and medical research.
- PET images play a crucial role in cancer diagnosis and treatment planning. PET scans can detect areas with increased glucose uptake, which is characteristic of many cancer cells. By overlaying PET images with other imaging modalities, such as CT or MRI, clinicians can precisely localize tumor regions and assess the extent of cancer spread. This information is vital for treatment decisions and evaluating treatment response over time.
- Data modeling in PET imaging involves advanced image reconstruction algorithms to generate high-resolution images with enhanced signal-to-noise ratios. These reconstructed images provide better visualization and more accurate quantification of tracer uptake, enabling more precise analysis and interpretation of the metabolic activity in various tissues.
- PET images are essential in neuroscience research, particularly in the study of brain function and neurological disorders. By using radiotracers that bind to specific neurotransmitter receptors, researchers can investigate the distribution and activity of these receptors in the brain. PET data modeling helps quantify receptor binding potentials and create brain activation maps, facilitating the study of various neurological conditions, such as Alzheimer's disease and schizophrenia.
- In cardiovascular research, PET imaging is used to assess myocardial perfusion and metabolism. Radiotracers like Rubidium-82 and Nitrogen-13 ammonia are employed to evaluate blood flow and metabolic activity in the heart. These data are valuable for diagnosing coronary artery disease, evaluating myocardial viability, and guiding treatment strategies.

- PET images are also utilized in assessing various other diseases, including infectious and inflammatory conditions. Infection-specific radiotracers, like Fluorine-18 FDG, can highlight areas of inflammation or infection, aiding in the diagnosis and monitoring of infectious diseases.
- Data modeling in PET research involves kinetic modeling techniques to quantify the pharmacokinetics of radiotracers and extract relevant physiological parameters. These models help researchers understand tracer distribution and clearance mechanisms, leading to a deeper understanding of disease processes and treatment responses.
- PET images find applications in pharmacology and drug development. Preclinical PET studies are conducted to evaluate the distribution and kinetics of potential drug candidates in animal models. The data obtained from these studies aid in selecting the most promising compounds for further development and clinical trials.
- In personalized medicine, PET imaging is used to stratify patients based on their molecular characteristics and predict their response to specific treatments. By using PET scans to assess the expression of certain biomarkers, clinicians can tailor treatment plans to individual patients, optimizing therapeutic outcomes.
- Data visualization plays a critical role in PET image analysis and interpretation. Techniques like voxel-based analysis and region-of-interest (ROI) analysis allow for the visualization of specific regions with altered metabolic activity, aiding in the identification of pathological changes.
- PET images can be co-registered with other imaging modalities, such as MRI or CT, to provide multimodal visualizations that combine anatomical and functional information. These fused images are invaluable for understanding the relationships between structure and function in various tissues and organs.
- Quantitative visualization of PET data is essential for accurate diagnosis and research analysis. Software tools allow for the extraction of standardized uptake values (SUVs) and other quantitative measures, facilitating comparisons across patients and research studies.
- In summary, PET images are a vital resource in data modeling and visualization in various medical and research applications. The combination of data modeling and visualization techniques enables the accurate quantification and interpretation of metabolic activity and functional changes in tissues. This information is crucial for diagnosing diseases, planning treatments, understanding disease mechanisms, and developing personalized therapeutic approaches, contributing to advancements in medical science and patient care.

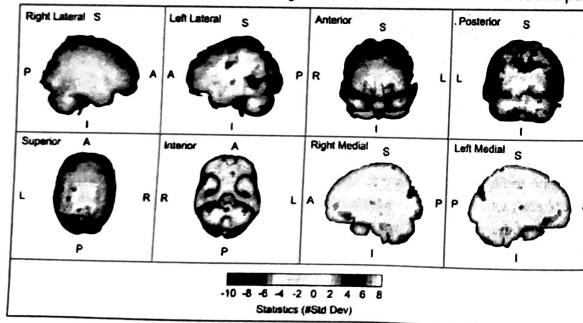


Fig. 6.3.1 : In PET image analysis

**6.3.1 UltraSound Images**

- Ultrasound imaging, also known as sonography, is a non-invasive medical imaging technique that uses high-frequency sound waves to produce real-time images of internal body structures. Ultrasound images are widely used in various medical applications, and their analysis through data modeling and visualization plays a crucial role in diagnosis, treatment planning, and research.
- One of the primary applications of ultrasound images is in obstetrics, where they are used to monitor fetal development during pregnancy. Data modeling in obstetric ultrasound involves creating algorithms to estimate gestational age, measure fetal growth, and assess fetal well-being. Visualization of these ultrasound images allows healthcare providers to identify any potential abnormalities or complications in the developing fetus.
- In radiology, ultrasound images aid in the diagnosis of various conditions affecting internal organs, such as the liver, kidneys, and gallbladder. Data modeling techniques are used to improve image quality and enhance the diagnostic accuracy of these images. Visualization of ultrasound scans enables radiologists to detect tumors, cysts, and other abnormalities, helping in the early detection and management of diseases.
- Ultrasound images also play a vital role in the assessment of cardiovascular health. In echocardiography, ultrasound images of the heart help evaluate cardiac structure and function. Data modeling techniques are utilized to calculate parameters such as ejection fraction and cardiac output, providing valuable information for diagnosing heart diseases and guiding treatment decisions. Visualization of these images allows cardiologists to visualize the heart's movements and assess any cardiac abnormalities.
- In urology, ultrasound images are used to examine the urinary system, including the kidneys, bladder, and prostate. Data modeling techniques in urologic ultrasound help measure kidney size, assess kidney function, and diagnose conditions like kidney stones and urinary tract infections. Visualization of these images aids urologists in guiding biopsies, placing drainage catheters, and planning surgical interventions.
- In musculoskeletal imaging, ultrasound is used to visualize tendons, ligaments, muscles, and joints. Data modeling techniques enable accurate measurement of structures and facilitate diagnosis of conditions like tendinitis, bursitis, and joint effusions. Visualization of musculoskeletal ultrasound images assists orthopedic surgeons and sports medicine specialists in planning treatments and monitoring patients' progress.
- Ultrasound images have applications in emergency medicine, where they are used to assess trauma patients and guide medical interventions. Data modeling techniques help identify internal bleeding, evaluate organ injuries, and visualize foreign objects. Visualization of these images assists emergency physicians in making rapid and informed decisions in critical situations.
- In interventional radiology, ultrasound is used as a real-time imaging guide during minimally invasive procedures. Data modeling techniques allow for precise needle placements and catheter insertions, reducing the risk of complications. Visualization of ultrasound images in interventional procedures enhances safety and accuracy during the interventions.
- Ultrasound images are valuable in evaluating breast health. Data modeling techniques are employed in breast ultrasound to differentiate between benign and malignant masses and to aid in breast cancer staging. Visualization of these images helps radiologists in detecting breast abnormalities and guiding biopsy procedures.
- In gastroenterology, ultrasound imaging is used to assess the liver, gallbladder, pancreas, and other gastrointestinal organs. Data modeling techniques enable accurate measurement of liver steatosis, fibrosis, and tumor characterization. Visualization of these images aids gastroenterologists in diagnosing liver and biliary diseases and guiding therapeutic procedures.

- Ultrasound is also used in assessing thyroid nodules and other neck structures. Data modeling techniques help distinguish benign from malignant nodules, reducing unnecessary biopsies. Visualization of thyroid ultrasound images aids endocrinologists in diagnosing thyroid disorders and determining appropriate treatments.
- In data modeling for ultrasound imaging, algorithms are developed for image enhancement, noise reduction, and feature extraction. These modeling techniques improve the quality and diagnostic value of ultrasound images, making them more reliable for clinical decision-making.
- Visualization of ultrasound images is essential for interpreting the acquired data. Real-time visualization allows healthcare professionals to observe dynamic movements and changes within the body, facilitating immediate feedback and diagnosis.
- Three-dimensional (3D) and four-dimensional (4D) ultrasound visualization techniques enable a more comprehensive understanding of complex structures and anatomical relationships. These techniques are particularly useful in fetal imaging, cardiac evaluation, and guided interventions.
- In research settings, data modeling and visualization of ultrasound images are used to develop new imaging techniques, study disease mechanisms, and evaluate the effectiveness of novel therapeutic interventions. These applications contribute to advancements in medical science and the development of improved diagnostic and treatment methods.
- In conclusion, ultrasound images have a wide range of applications in medicine and medical research. Data modeling and visualization of ultrasound images play a crucial role in accurate diagnosis, treatment planning, and research. Through the use of sophisticated algorithms and visualization techniques, ultrasound imaging continues to evolve, providing valuable insights for healthcare professionals and contributing to improved patient care.

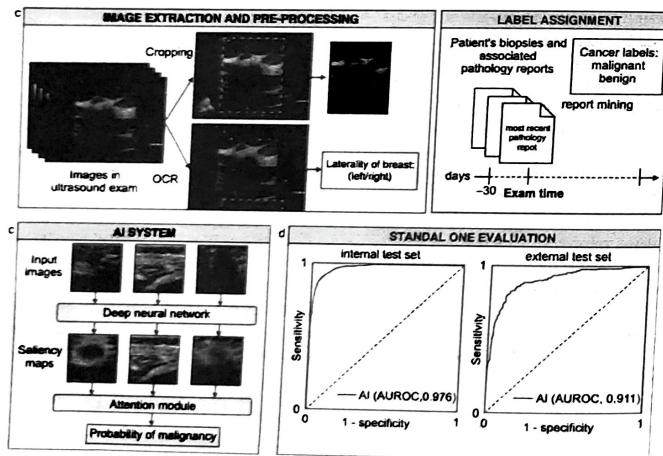


Fig. 6.3.2 : Image Extraction and Pre-Processing

- Magnetic Resonance Imaging (MRI) is a non-invasive medical imaging technique that uses magnetic fields and radio waves to generate detailed images of the internal structures of the body. MRI provides excellent soft tissue contrast, making it a valuable tool in various medical applications. Data modeling and visualization of MRI images play a critical role in diagnosis, treatment planning, and research.
- One of the primary applications of MRI is in neuroimaging, where it is used to study the structure and function of the brain and spinal cord. Data modeling techniques in neuroimaging involve advanced image processing algorithms to segment and analyze brain structures, identify abnormalities, and quantify brain volumes and connectivity. Visualization of MRI data in neuroimaging allows researchers and clinicians to explore brain networks and identify regions of interest in neurological disorders like Alzheimer's disease, multiple sclerosis, and stroke.
- In oncology, MRI is used to diagnose and stage various cancers, such as brain tumors, breast cancer, and prostate cancer. Data modeling techniques in oncologic MRI involve computer-aided analysis for tumor segmentation, assessment of tumor characteristics, and prediction of treatment responses. Visualization of MRI data in oncology enables precise tumor localization and evaluation of tumor progression, guiding treatment decisions and monitoring therapy effectiveness.
- MRI is also widely used in musculoskeletal imaging to assess joint and soft tissue conditions, such as ligament tears, cartilage injuries, and herniated discs. Data modeling techniques in musculoskeletal MRI involve the quantification of cartilage thickness, joint space, and bone morphology. Visualization of MRI data in musculoskeletal imaging allows orthopedic surgeons and sports medicine specialists to plan surgical interventions and monitor treatment outcomes.
- In cardiovascular imaging, MRI provides detailed information about the heart's structure and function. Data modeling techniques in cardiac MRI include cardiac motion analysis, evaluation of blood flow, and assessment of myocardial viability. Visualization of MRI data in cardiac imaging helps cardiologists diagnose heart diseases, evaluate the impact of coronary artery blockages, and plan cardiac surgeries.
- MRI plays a vital role in abdominal imaging, enabling the evaluation of organs such as the liver, kidneys, and pancreas. Data modeling techniques in abdominal MRI involve automated segmentation of organs and the quantification of tissue characteristics. Visualization of MRI data in abdominal imaging assists radiologists and gastroenterologists in diagnosing liver diseases, kidney disorders, and pancreatic tumors.
- Functional MRI (fMRI) is a specialized MRI technique used to study brain activity during specific tasks or at rest. Data modeling techniques in fMRI involve statistical analysis to identify brain regions activated during different cognitive processes. Visualization of fMRI data allows researchers to create activation maps, revealing how specific brain areas are involved in tasks like language processing, memory retrieval, and motor functions.
- Diffusion-Weighted Imaging (DWI) is an MRI technique that assesses the movement of water molecules in tissues. Data modeling in DWI involves diffusion tensor imaging (DTI) to generate fiber tractography maps. DWI visualization aids in studying brain connectivity and understanding the effects of brain injuries and neurodegenerative diseases.

- MRI is used in fetal imaging to assess fetal development and detect congenital abnormalities. Data modeling techniques in fetal MRI involve motion correction and artifact reduction to obtain clear images of the moving fetus. Visualization of fetal MRI data helps obstetricians and pediatric specialists diagnose fetal anomalies and plan appropriate care for expectant mothers and their babies.
- In research settings, data modeling and visualization of MRI images are used to develop novel imaging techniques, study disease mechanisms, and evaluate the effectiveness of experimental treatments. These applications contribute to advancements in medical science and the development of improved diagnostic and therapeutic methods.
- Machine learning and artificial intelligence algorithms are increasingly being applied to MRI data to assist in automated segmentation, disease detection, and prediction of treatment outcomes. These data modeling techniques have the potential to improve diagnostic accuracy and efficiency, enabling more personalized and effective healthcare.

### Conclusion

- In conclusion, Magnetic Resonance Imaging (MRI) is a versatile medical imaging technique with numerous applications in various fields of medicine. Data modeling and visualization of MRI images are essential for accurate diagnosis, treatment planning, and research. By employing sophisticated algorithms and visualization techniques, MRI continues to evolve as a powerful tool in modern healthcare, providing valuable insights for medical professionals and contributing to improved patient care.

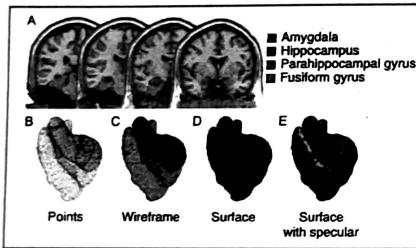


Fig. 6.3.3 : Functional MRI

### Discussion

#### Title : Discussion on Different Types of Images and Techniques Used in Data Modelling and Visualization

In the realm of data science, data modelling and visualization play a critical role in understanding complex datasets and extracting meaningful insights. Images, as a form of visual representation, have proven to be powerful tools for conveying information and patterns in various fields, from medicine to finance, and from astronomy to engineering. This discussion aims to explore different types of images and the techniques employed in data modelling and visualization to enable effective analysis and decision-making.

- Types of Images In Data Modelling :** Data in various domains can be represented using different types of images. For instance, medical imaging relies on techniques like X-rays, MRI, CT scans, and ultrasounds to visualize internal structures and diagnose diseases. In astronomy, telescopic images reveal celestial objects,

while satellite imagery provides valuable geospatial data. In computer vision, images are used to train machine learning models for tasks such as object detection and image classification. Furthermore, graphs and charts, often presented as visual images, depict trends and patterns in numerical data, making them accessible to a broader audience.

- Techniques in Image Preprocessing :** Before data modelling and visualization, image preprocessing is crucial to enhance the quality and extract relevant information. Techniques like resizing, normalization, and noise reduction help standardize image inputs for more accurate analyses. Image segmentation is also employed to partition images into meaningful regions, allowing researchers to focus on specific areas of interest and simplifying the data analysis process.
- Image Feature Extraction :** Extracting relevant features from images is vital for effective data modelling. Techniques such as edge detection, texture analysis, and feature descriptors like SIFT (Scale-Invariant Feature Transform) and HOG (Histogram of Oriented Gradients) aid in identifying distinctive patterns and structures within images. Feature extraction reduces the dimensionality of data, facilitating efficient modeling and visualization.
- Image Classification :** Image classification involves categorizing images into predefined classes or labels. Machine learning algorithms, such as convolutional neural networks (CNNs), have revolutionized image classification tasks. CNNs learn hierarchical features from images and can accurately classify objects, making them invaluable in various applications like facial recognition, autonomous vehicles, and medical diagnoses.
- Object Detection :** Object detection is the process of identifying and locating specific objects within images. Techniques like the region-based R-CNN (Region-based Convolutional Neural Network) and single-stage detectors like YOLO (You Only Look Once) enable real-time object detection, enabling applications such as surveillance, robotics, and augmented reality.
- Image Segmentation :** Image segmentation aims to divide an image into coherent and meaningful regions. Algorithms like U-Net, FCN (Fully Convolutional Network), and Mask R-CNN have shown significant advancements in semantic and instance segmentation tasks. This technique is widely used in medical image analysis, autonomous driving, and computer-aided design.
- Image Generative Models :** Generative models, such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), can synthesize new images based on the patterns and structures learned from the training data. These models have fascinating applications in art generation, data augmentation, and anomaly detection.
- Data Visualization Techniques :** Apart from image-based visualization, other techniques are used to represent complex data visually. Graphs, scatter plots, heatmaps, and treemaps are commonly employed to display relationships, distributions, and hierarchical structures within data. Interactive visualizations using tools like D3.js enhance user engagement and facilitate data exploration.
- Augmented Reality and Virtual Reality :** Augmented Reality (AR) and Virtual Reality (VR) are revolutionizing data visualization by overlaying or immersing users in virtual environments. These technologies have been applied in education, gaming, and architecture to create compelling visual experiences and simulate real-world scenarios.

- 10. Visual Analytics :** Visual analytics combines data visualization with advanced data analysis techniques. It enables users to interactively explore data, discover hidden patterns, and gain insights from complex datasets. This interdisciplinary approach has found applications in healthcare, finance, and business intelligence.
- 11. Challenges in Data Modelling and Visualization :** Despite the advancements in data modelling and visualization, several challenges persist. These include dealing with large-scale image datasets, ensuring interpretability of deep learning models, handling imbalanced classes in image classification, and addressing ethical concerns related to privacy and bias in visual data.
- 12. Ethical Considerations :** As data modelling and visualization techniques become more sophisticated, ethical considerations become paramount. It is essential to handle sensitive image data responsibly, protect individual privacy, and avoid perpetuating biases that could lead to unfair decisions and discrimination.
- 13. Future Directions :** The future of data modelling and visualization lies in the integration of various technologies, such as AI, AR, and VR, to create more immersive and intuitive visual experiences. Improving the explainability of AI models and addressing ethical challenges will be crucial for the widespread adoption of these technologies.
- 14. Importance of Collaboration :** Collaboration between domain experts, data scientists, and visualization specialists is essential for successful data modelling and visualization projects. Combining domain knowledge with technical expertise leads to more accurate insights and actionable recommendations.
- 15. Conclusion :** In conclusion, data modelling and visualization with different types of images have become indispensable in understanding complex datasets and making informed decisions. The combination of advanced algorithms, interactive visualizations, and emerging technologies promises a future where data-driven insights are accessible and actionable across various industries. However, it is essential to remain mindful of ethical considerations and challenges as we harness the power of visual data for the betterment of society.

## 6.4 Case Study

### 6.4.1 ER/Studio

- ER/Studio is a comprehensive and powerful data modeling tool developed by Quest Software, now part of IDERA. It serves as a robust platform for data professionals to design, manage, and document data assets through the process of data modeling. Data modeling is a crucial phase in the database development lifecycle, where conceptual, logical, and physical data models are created to represent the structure and relationships of data in an organization. ER/Studio provides a user-friendly interface, making it accessible to both technical and non-technical stakeholders involved in the data modeling process.
- One of the primary uses of ER/Studio is in conceptual data modeling. It allows data architects to create high-level business requirements and entities, fostering better communication between business users and data professionals. The tool facilitates logical data modeling, capturing the relationships and attributes of data entities, irrespective of the underlying database management system. Logical data models bridge the gap between business requirements and technical implementation, ensuring alignment between stakeholders.

- ER/Studio enables the conversion of logical data models into physical data models that are specific to a particular database platform. This ensures compatibility and efficient data storage and retrieval. By supporting collaborative data modeling efforts, ER/Studio allows multiple users to work on the same data model simultaneously, promoting teamwork and accelerating the data modeling process.
- Another key feature is version control and change management, enabling data professionals to track model changes, compare versions, and revert to previous states when necessary. This maintains a comprehensive history of data models and ensures data integrity throughout the development lifecycle.
- ER/Studio is also instrumental in performing impact analysis. Data professionals can assess the potential effects of data changes on the overall data architecture, aiding in understanding the implications of data modifications and making informed decisions.
- The tool plays a crucial role in supporting data governance initiatives. With a centralized repository for data models, ER/Studio ensures consistent data modeling standards and compliance with regulatory requirements, such as GDPR and CCPA.
- Data lineage tracking is facilitated by ER/Studio, showing the origin and transformation of data throughout its lifecycle. Additionally, the tool aids in generating comprehensive documentation for data models, enhancing understanding and communication with stakeholders.
- Another valuable application of ER/Studio is reverse engineering, which allows data professionals to transform existing databases into visual data models. This is particularly useful for understanding legacy systems and maintaining data consistency during migration projects.
- ER/Studio's visualization capabilities enable data professionals to present data models and metadata in a visually appealing and easily understandable format. This enhances communication with stakeholders and supports effective decision-making processes.
- The tool supports integration with various database management systems, allowing seamless synchronization of data models with the actual databases. This ensures that the implemented database matches the intended data model, minimizing discrepancies.
- ER/Studio is also employed in data warehouse design and the definition of Extract, Transform, Load (ETL) processes. Data modeling for data warehousing is essential for creating efficient data structures and optimizing query performance.
- Forward engineering is another vital feature of ER/Studio, allowing data professionals to generate database scripts and DDL (Data Definition Language) statements from data models. This simplifies the implementation process and reduces errors.
- In conclusion, ER/Studio is a versatile and indispensable data modeling tool with various applications in data modeling and visualization. It supports the entire data modeling lifecycle, from conceptual to physical data modeling, promoting collaboration, compliance, and data governance. With robust visualization and modeling, ER/Studio facilitates communication among stakeholders and supports data documentation capabilities. ER/Studio integrates with different DBMS platforms, supports impact analysis, and drives decision-making. Its integration with ETL design makes it a valuable asset in the field of data modeling and visualization.

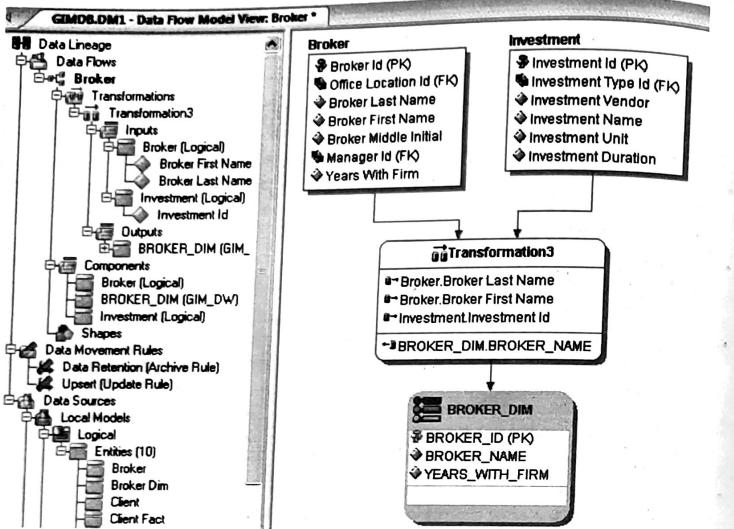


Fig. 6.4.1

#### 6.4.2 Erwin Data Modeler

- Erwin Data Modeler is a powerful and widely-used data modeling tool developed by Erwin, Inc. It is designed to assist organizations in effectively designing, managing, and visualizing their data models. Data modeling is a crucial step in the database development lifecycle, where data professionals create conceptual, logical, and physical data models to represent the structure and relationships of data within an organization.
- Erwin Data Modeler offers a user-friendly interface, making it accessible to both technical and non-technical users. Its drag-and-drop functionality simplifies the process of creating and modifying data models, allowing for efficient data modeling.
- One of the key features of Erwin Data Modeler is its support for conceptual, logical, and physical data models. This enables data professional to capture high-level business requirements, define data structures, and generate database-specific models, bridging the gap between business needs and technical implementation.
- The tool facilitates reverse engineering, allowing users to extract data models from existing databases. This feature is particularly valuable for understanding and documenting legacy systems, ensuring data consistency during migrations, and supporting data integration initiatives.

- Forward engineering is another essential capability of Erwin Data Modeler. It enables users to generate SQL scripts and Data Definition Language (DDL) statements from data models, streamlining the process of implementing databases and ensuring alignment with the intended data model.
- Erwin Data Modeler also supports collaborative data modeling efforts. Multiple users can work on the same data model simultaneously, promoting teamwork and accelerating the data modeling process.
- Version control and change management features ensure that data professionals can track model changes, compare versions, and maintain data integrity, making it easier to manage and document the evolution of data models.
- Impact analysis is another critical functionality provided by Erwin Data Modeler. Users can assess the potential effects of data changes on the overall data architecture, helping them understand data dependencies and make informed decisions.
- Data visualization capabilities in Erwin Data Modeler enable users to present data models in a visually appealing and easily understandable format. This aids in effective communication with stakeholders and facilitates data-driven decision-making.
- The tool also supports customization, allowing users to tailor data models according to their organization's specific needs. It enables the generation of comprehensive documentation and reports, aiding in data governance and compliance with regulations.
- Integration with various Database Management Systems (DBMS) ensures that data models are compatible with the selected database platforms, making implementation seamless and minimizing discrepancies.
- Data lineage tracking is facilitated by Erwin Data Modeler, providing insights into the origin and transformation of data throughout its lifecycle. Additionally, the tool assists in performing impact analysis, identifying the consequences of data changes.
- In data warehouse design, Erwin Data Modeler is instrumental in creating efficient data structures for optimal query performance, supporting organizations in gaining insights from data and making informed decisions.
- For business intelligence solutions, the tool can be used to design data models, enabling organizations to extract valuable insights from their data and make data-driven decisions.
- Erwin Data Modeler also supports data integration initiatives and the design of Extract, Transform, Load (ETL) processes, ensuring seamless data movement between systems.
- Data governance initiatives are enhanced through Erwin Data Modeler's ability to maintain consistency and compliance with data regulations, making it an essential tool in supporting data governance efforts.
- In conclusion, Erwin Data Modeler is a versatile and feature-rich data modeling tool that plays a significant role in data modeling and visualization. Its user-friendly interface, support for conceptual, logical, and physical data models, collaboration features, and data visualization capabilities make it a valuable asset for data professionals seeking efficient data management and decision-making. Its integration with various DBMS platforms, support for impact analysis, assistance in data warehouse and ETL design, and contributions to data governance initiatives further underline its importance in the field of data modeling and visualization.

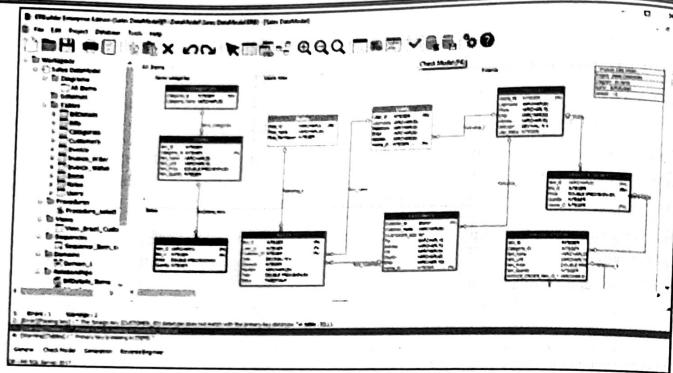


Fig. 6.4.2

#### 6.4.3 DbSchemaPro

DbSchema Pro is a comprehensive and feature-rich data modeling and visualization tool that caters to the needs of database professionals and data architects. With a user-friendly interface and powerful capabilities, DbSchema Pro streamlines the process of designing, managing and visualizing complex databases, making it an invaluable asset in the field of data modeling and visualization.

##### 1. Intuitive User Interface

DbSchema Pro offers an intuitive and visually appealing user interface, enabling users to navigate and interact with the tool effortlessly. Its drag-and-drop functionality and well-organized features make data modeling tasks more accessible to both technical and non-technical users.

##### 2. Multi-Platform Support

DbSchema Pro supports multiple database management systems (DBMS), including MySQL, PostgreSQL, Oracle, Microsoft SQL Server, and more. This cross-platform compatibility allows users to work with different databases seamlessly.

##### 3. Visual Data Modeling

One of the standout features of DbSchema Pro is its ability to create visual data models. It offers various diagramming tools to design conceptual, logical, and physical data models, representing complex database structures visually.

##### 4. Interactive Diagrams

DbSchema Pro's interactive diagrams enable users to explore and navigate through the data models with ease. The tool allows users to drill down into specific entities, attributes, and relationships, aiding in better understanding of the data structure.

5. **Reverse Engineering**  
DbSchema Pro facilitates reverse engineering, allowing users to import existing databases and automatically generate visual data models. This feature is valuable for understanding and documenting legacy databases.
6. **Synchronization and Migration**

DbSchema Pro supports database synchronization, enabling users to compare and merge database structures to keep them in sync. It also assists in database migration, ensuring smooth transitions between different database versions.

##### 7. Data Query and Visualization

The tool offers a built-in SQL editor that allows users to write and execute SQL queries directly within the application. Query results can be visualized using charts, graphs, and pivot tables for enhanced data analysis.

##### 8. Schema Documentation

DbSchema Pro generates comprehensive and visually appealing documentation for data models, facilitating communication with stakeholders and ensuring data governance.

##### 9. Collaboration and Teamwork

The tool enables collaborative data modeling efforts by supporting concurrent editing and version control. Multiple team members can work on the same data model simultaneously, enhancing teamwork and productivity.

##### 10. Virtual Foreign Keys

DbSchema Pro supports virtual foreign keys, allowing users to define relationships between tables even if there are no explicit foreign keys in the database. This feature aids in data exploration and analysis.

##### 11. Data Import and Export

DbSchema Pro offers a range of data import and export options, enabling users to migrate data between various databases and formats with ease.

##### 12. Data Visualization with Charts

DbSchema Pro provides built-in charting capabilities that enable users to visualize data from queries or database tables in the form of pie charts, bar charts, line charts, and more.

##### 13. ER Diagram Generation

DbSchema Pro automatically generates Entity-Relationship (ER) diagrams based on imported database structures, providing users with a clear and organized representation of data relationships.

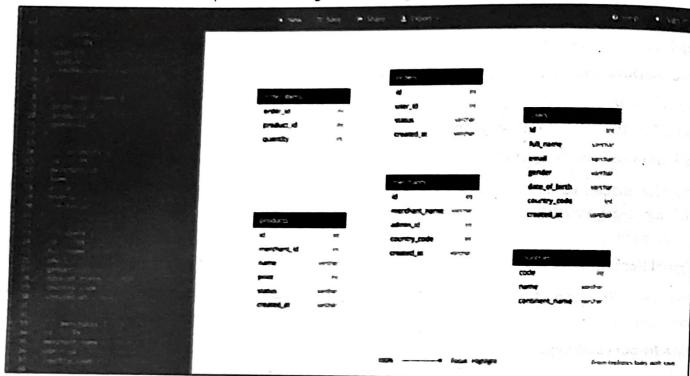
##### 14. Security and Data Privacy

The tool supports data masking and anonymization, helping users protect sensitive data during the modeling and visualization process.

##### 15. Virtual Database Access

DbSchema Pro allows users to connect to virtual databases, such as AWS RDS, Microsoft Azure, and Google Cloud SQL, providing a secure and efficient way to interact with cloud-based databases.

- In conclusion, DbSchema Pro is a versatile and powerful data modeling and visualization tool that empowers database professionals and data architects to design, manage, and visualize complex databases with ease. Its intuitive user interface, multi-platform support, interactive diagrams, and comprehensive documentation capabilities make it a valuable asset in the data modeling process.
- With features like reverse engineering, synchronization, data query, and collaboration support, DbSchema Pro streamlines the data modeling workflow and enhances teamwork and productivity. The tool's built-in charting, virtual foreign keys, and data import/export functionalities further contribute to effective data visualization and analysis. Additionally, DbSchema Pro's focus on security and data privacy ensures that sensitive information remains protected during the modeling and visualization journey.



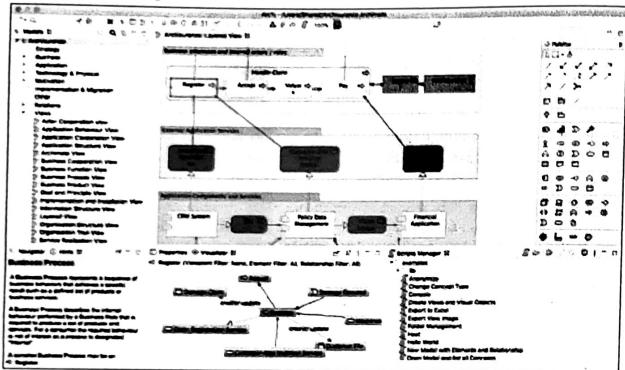
#### 6.4.4 Archi

- Archi is an open-source data modeling and visualization tool that is specifically designed for enterprise architecture and business process modeling. As a powerful and flexible tool, Archi serves as a valuable asset for data professionals and enterprise architects seeking to model and visualize complex systems, processes, and data structures.
- One of the key features of Archi is its user-friendly interface, which allows both technical and non-technical users to navigate and interact with the tool effortlessly. With its drag-and-drop functionality and straightforward modeling components, Archi simplifies the data modeling process and ensures seamless user experience.
- Archi excels in conceptual data modeling, where it enables users to represent high-level business concepts and their relationships. By creating visual models of business processes and systems, stakeholders can better understand the interdependencies and flow of data within an organization, facilitating effective communication and decision-making.

The tool's support for UML (Unified Modeling Language) and BPMN (Business Process Model and Notation) standards makes it a comprehensive solution for data modeling and visualization. Users can create UML class diagrams, use case diagrams, activity diagrams, and BPMN diagrams to model different aspects of the enterprise architecture.

- Reverse engineering is another essential feature offered by Archi, allowing users to import existing systems and databases to automatically generate visual models. This capability is crucial for understanding and documenting legacy systems, facilitating their integration with newer technologies.
- Archi's data visualization capabilities enable users to create interactive diagrams and flowcharts that provide a clear representation of data structures, processes, and relationships. These visualizations aid in comprehending complex data architectures and communicating them to stakeholders more effectively.
- The tool supports version control, ensuring that data models are maintained and tracked throughout the development lifecycle. Users can compare different versions of data models, identify changes, and revert to previous states as needed, promoting data integrity and consistency.
- Archi allows for collaborative data modeling efforts, enabling multiple users to work on the same project simultaneously. This feature fosters teamwork and facilitates brainstorming and input from various stakeholders, leading to better data models and informed decisions.
- Data modeling with Archi extends beyond individual projects, as it supports the creation of reusable data components and templates. These reusable assets can be applied across multiple projects, saving time and effort in the modeling process and promoting consistency across the organization.
- The tool's documentation capabilities are also noteworthy, as it generates detailed and visually appealing reports for data models and architectural diagrams. These reports can be shared with stakeholders, project teams, and management to provide valuable insights into the enterprise architecture.
- Archi offers extensive customization options, allowing users to tailor data models and diagrams according to the organization's specific needs and branding guidelines. Customization ensures that data models align with the organization's unique requirements and presentation standards.
- For effective data visualization, Archi allows users to create visualizations in different formats, including PNG, SVG, and PDF. These visualizations can be used in presentations, documentation, and reports to provide a comprehensive understanding of the data and its relationships.
- The tool's integration with enterprise architecture frameworks, such as TOGAF (The Open Group Architecture Framework), provides users with a structured approach to data modeling and enterprise architecture design. This integration ensures adherence to industry best practices and standards.
- Archi's flexibility extends to its support for add-ons and plugins, enabling users to extend the tool's functionality based on their specific needs. These add-ons and plugins enhance the modeling and visualization capabilities of Archi, catering to various use cases and industries.
- As an open-source tool, Archi benefits from an active community of users and contributors who continuously improve and enhance its features. This active community support ensures that Archi remains relevant and up-to-date in the ever-evolving field of data modeling and visualization.
- In conclusion, Archi is a powerful and versatile data modeling and visualization tool that caters to the needs of enterprise architects, business analysts, and data professionals. With its user-friendly interface, support for UML and BPMN standards, reverse engineering capabilities, and collaborative features, Archi simplifies data modeling and promotes effective communication of complex data architectures.

- Its integration with enterprise architecture frameworks, documentation capabilities, and support for reusable components make it an ideal choice for modeling enterprise-level data structures and processes. As an open-source tool with a thriving community, Archi continues to evolve and adapt to the changing demands of data modeling and visualization, making it a valuable asset for organizations seeking efficient and well-structured enterprise architecture design.



#### 6.4.5 SQL Database Modeler

- SQL Database Modeler is a specialized tool designed for data modeling and visualization in the context of SQL databases. It provides database professionals, data architects, and developers with powerful features and capabilities to design, manage, and visualize complex database structures efficiently.
- One of the key features of SQL Database Modeler is its ability to create visual data models that represent the structure and relationships of SQL databases. Users can design conceptual, logical, and physical data models, ensuring a comprehensive understanding of the database architecture.
- The tool supports various database management systems (DBMS), including MySQL, PostgreSQL, Microsoft SQL Server, Oracle, and others, making it compatible with a wide range of SQL database platforms.
- SQL Database Modeler allows for reverse engineering, enabling users to import existing databases and generate visual data models automatically. This feature is beneficial for understanding and documenting legacy databases and ensuring smooth integration with newer technologies.
- Interactive and navigable diagrams are integral to SQL Database Modeler, providing users with a clear visualization of the database schema, tables, columns, and relationships. This visual representation aids in comprehending complex data structures and communicating them effectively to stakeholders.
- The tool's support for SQL scripting and query execution allows users to write and execute SQL queries directly within the application. This feature facilitates data analysis, data manipulation, and validation directly within the modeling environment.

- SQL Database Modeler supports database synchronization, enabling users to compare and merge database structures to maintain consistency across various environments.
- Data migration is made more accessible with SQL Database Modeler, as it assists in converting data models between different database platforms, streamlining the process of database migration.
- Comprehensive documentation generation is a key feature, as SQL Database Modeler automatically creates detailed reports and visual documentation for data models, making it easier to share insights with stakeholders and maintain data governance.
- The tool's visualization capabilities extend to charting, allowing users to create various types of charts and graphs to visualize query results or database statistics.
- SQL Database Modeler allows for customization, enabling users to tailor data models and diagrams according to their organization's unique requirements and branding guidelines.
- Data visualization with SQL Database Modeler encompasses exporting visualizations in different formats such as PNG, SVG, and PDF, enabling seamless integration into presentations, documentation, and reports.
- The tool's integration with version control systems enables users to track changes in data models, compare different versions, and revert to previous states when necessary, ensuring data integrity throughout the modeling process.
- Collaborative data modeling is facilitated by SQL Database Modeler, as it supports concurrent editing and team collaboration, promoting effective teamwork and project management.
- As an essential tool for database professionals and data architects, SQL Database Modeler provides a seamless and efficient solution for data modeling and visualization in SQL databases. With its visual data modeling capabilities, reverse engineering, interactive diagrams, SQL scripting support, and database synchronization, SQL Database Modeler streamlines the data modeling process and enhances productivity.
- Its compatibility with various DBMS platforms, data migration support, comprehensive documentation generation, and charting capabilities make it a valuable asset for organizations seeking to manage and visualize complex SQL databases effectively.

#### 6.4.6 LucidChart

- Lucidchart is a web-based diagramming and visualization tool that offers a wide range of applications in data modeling and visualization. As a popular choice among professionals, Lucidchart provides intuitive features and collaborative capabilities that enhance data modeling and visualization processes for database professionals, data architects, and business analysts.
- One of the key features of Lucidchart is its user-friendly interface, allowing both technical and non-technical users to create complex data models and visualizations easily. Its drag-and-drop functionality and extensive library of shapes and templates make it accessible to users with varying levels of expertise.
- Lucidchart is ideal for conceptual, logical, and physical data modeling. Users can design high-level business requirements and convert them into detailed data structures, enabling effective communication between business stakeholders and technical teams.

- With its extensive set of predefined templates and diagramming options, Lucidchart supports various modeling techniques, including Entity-Relationship (ER) diagrams, UML diagrams, and Data Flow Diagrams (DFD). This versatility allows users to represent data relationships, processes, and flows with precision.
- Collaborative features are at the core of Lucidchart, enabling multiple team members to collaborate in real-time on the same data model. This facilitates teamwork, streamlines decision-making, and ensures data models are up-to-date and accurate.
- Lucidchart's data visualization capabilities extend beyond data modeling, allowing users to create charts, graphs, and dashboards from live data sources, supporting data-driven decision-making and analysis.
- The tool's integration with cloud-based storage services like Google Drive, Microsoft OneDrive, and Dropbox enables seamless storage and sharing of data models and visualizations.
- Lucidchart offers reverse engineering functionality, allowing users to import existing databases and generate visual data models automatically. This feature is valuable for understanding and documenting legacy databases and systems.
- Lucidchart supports SQL code generation, allowing users to generate SQL scripts and DDL (Data Definition Language) statements from data models, facilitating database implementation and synchronization.
- Users can import and export data models in various file formats, ensuring compatibility with other modeling and database management tools.
- The platform's extensive security features, including role-based access controls, single sign-on (SSO), and data encryption, ensure the safety and confidentiality of sensitive data.
- Lucidchart provides data-driven automation capabilities, allowing users to define rules that automatically adjust data model elements based on changes or updates in the database.
- The tool's revision history and version control features allow users to track changes in data models, compare versions, and revert to previous states when needed, ensuring data integrity and accuracy.
- Lucidchart's presentation mode enables users to showcase data models and visualizations in a professional and visually engaging manner during meetings and presentations.
- The platform supports real-time collaboration with stakeholders, allowing non-technical users to provide input, feedback, and insights during the data modeling process.
- Lucidchart's seamless integration with other business productivity tools, such as Google Workspace, Microsoft Office, and Atlassian Jira, enhances the overall workflow and efficiency of data modeling and visualization projects.
- In conclusion, Lucidchart is a versatile and powerful web-based diagramming and visualization tool that offers extensive applications in data modeling and visualization. With its user-friendly interface, collaborative features, support for various modeling techniques, and integration with cloud-based storage and productivity tools, Lucidchart streamlines the data modeling process and facilitates effective communication among stakeholders.
- Its data visualization capabilities, reverse engineering support, SQL code generation, and data-driven automation features further enhance its value for database professionals and data architects. Lucidchart is a leading choice for data modeling and visualization, enabling organizations to design, manage, and communicate complex data structures and relationships with ease and precision.

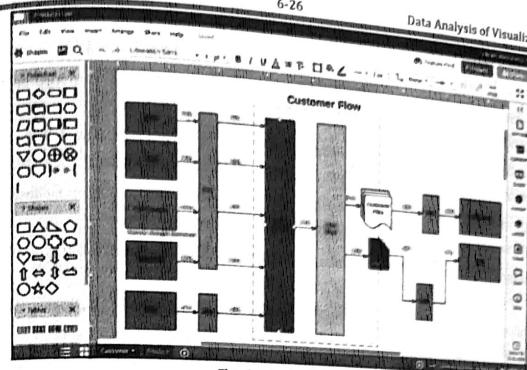
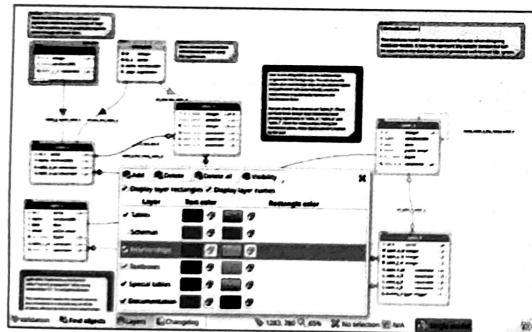


Fig. 6.4.3 : Lucid Chart

#### 6.4.7 PgModeler

- PgModeler is an advanced open-source data modeling and visualization tool specifically designed for PostgreSQL databases. As a feature-rich platform, PgModeler empowers database professionals, data architects, and developers to efficiently design, manage, and visualize complex data structures within PostgreSQL databases.
- One of the key features of PgModeler is its support for visual data modeling, enabling users to create conceptual, logical, and physical data models. The tool provides a comprehensive and intuitive interface for users to represent high-level business requirements and convert them into detailed data structures.
- PgModeler's reverse engineering capabilities allow users to import existing databases and automatically generate visual data models. This feature is invaluable for understanding and documenting legacy databases, ensuring smooth integration with modern technologies.
- With interactive and navigable diagrams, PgModeler provides a clear visualization of the PostgreSQL database schema, tables, columns, and relationships. This visual representation aids in comprehending complex data structures and effectively communicating them to stakeholders.
- The tool's SQL scripting support allows users to write and execute SQL queries directly within the application. This feature facilitates data analysis, data manipulation, and validation directly within the data modeling environment.
- PgModeler supports version control and integration with Git, allowing users to track changes in data models, compare different versions, and revert to previous states when needed. This ensures data integrity throughout the modeling process.
- Comprehensive documentation generation is another key strength of PgModeler. The tool automatically generates detailed reports and visual documentation for data models, making it easier to share insights with stakeholders and maintain data governance.
- PgModeler offers data migration support, allowing users to convert data models between different PostgreSQL database versions. This streamlines the process of migrating databases to newer versions.

- The platform supports the generation of SQL scripts and DDL (Data Definition Language) statements from data models. This feature simplifies the database implementation and synchronization process.
  - PgModeler's database synchronization capabilities enable users to compare and merge database structures to maintain consistency across various environments.
  - Data visualization in PgModeler extends beyond data modeling, as the tool provides charting capabilities, allowing users to create various types of charts and graphs to visualize query results or database statistics.
  - PgModeler provides efficient data modeling and visualization with its user-friendly interface, interactive diagrams, reverse engineering support, SQL scripting, and charting capabilities.
  - The tool's integration with PostgreSQL databases ensures compatibility and seamless communication with the database engine.
  - PgModeler allows for collaborative data modeling efforts, as it supports concurrent editing and team collaboration, promoting effective teamwork and project management.
  - The platform's flexibility and customization options enable users to tailor data models and diagrams according to their organization's specific requirements and branding guidelines.
  - As an open-source tool, PgModeler benefits from an active community of users and contributors, ensuring continuous improvement and enhancement of its features.
  - In conclusion, PgModeler is a robust and versatile open-source data modeling and visualization tool dedicated to PostgreSQL databases. With its support for visual data modeling, reverse engineering, interactive diagrams, SQL scripting, and charting capabilities, PgModeler streamlines the data modeling process and enhances productivity for database professionals and data architects.
  - Its integration with version control, data migration support, and comprehensive documentation generation further underline its value in managing and visualizing complex PostgreSQL databases. The platform's collaborative features, customization options, and seamless integration with PostgreSQL databases make PgModeler an invaluable asset for organizations seeking to optimize data modeling and visualization efforts in PostgreSQL environments.



**Fig. 6.4.4**

### **Review Questions**

- 0.1 Give examples of Reconstruction in data modeling and visualization.
  - 0.2 Explain various types of examples of reconstruction in data modeling and visualization.
  - 0.3 Explain medical images and application in data modelling and visualization.
  - 0.4 What are PET Images?
  - 0.5 What are Ultrasound Images?
  - 0.6 What are MRI Images?
  - 0.7 What is Erwin Data Modeler?
  - 0.8 What is ER/Studio? What are the Uses or Application of ER/Studio?
  - 0.9 Explain about DbSchemaPro data modeler and visualization tool.