

# STATISTICS AND PROBABILITY

What is Statistics and Probability explain with examples with respect to Statistical computing

- Statistics and probability are closely related fields that involve the study of data and the likelihood of events occurring. Statistical computing refers to the use of computational methods and tools to analyze and interpret data in the context of statistics and probability.

# Statistics:

- Statistics is the discipline that involves collecting, analyzing, interpreting, presenting, and organizing data. It provides methods to make inferences and predictions about populations based on a sample of data. Statistical techniques help in summarizing and describing data, testing hypotheses, and making informed decisions.
- *Example in Statistical Computing:* Imagine you have a dataset of heights of students in a school. Statistical computing can be used to calculate the mean (average) height, standard deviation (measure of variability), and conduct hypothesis tests to determine if there is a significant difference in heights between male and female students.

# Probability:

- Probability is the branch of mathematics that studies the likelihood of events happening. It assigns a numerical value between 0 and 1 to an event, where 0 indicates impossibility and 1 indicates certainty. Probability theory is used to model uncertainty and randomness.
- *Example in Statistical Computing:* Suppose you are interested in predicting the probability of a student passing an exam based on the number of hours they studied. Statistical computing can be employed to fit a logistic regression model to the data, which provides a probability of passing based on the hours of study.

# STATISTICAL COMPUTING

- **Statistical Computing:** Statistical computing involves the use of computer algorithms and software to perform statistical analysis on data. This includes tasks such as data cleaning, exploratory data analysis, hypothesis testing, regression analysis, and more.
- *Example in Statistical Computing:* You have a dataset of customer purchases and want to analyze the average spending per customer. Using statistical computing, you can calculate the mean and standard deviation of the spending, create visualizations, and perform hypothesis tests to determine if there are significant differences in spending habits across different customer segments.

# SUMMARY

- In summary, statistics and probability are foundational concepts in data analysis, and statistical computing provides the tools and methods to apply these concepts to real-world datasets for making informed decisions and predictions.

# DATA VISUALISATION

- Data visualization is the graphical representation of data to help uncover patterns, trends, and insights that might be challenging to identify in raw, numerical data. It is a crucial aspect of statistical computing because it allows analysts, data scientists, and decision-makers to communicate complex information effectively. Various types of visualizations, ranging from simple charts to intricate interactive dashboards, are used to present data visually.

# Examples of Data Visualization with Statistical Computing:

## 1. Histograms:

- 1. Description:* Histograms are used to display the distribution of a dataset.
- 2. Statistical Context:* In statistical computing, you might create a histogram to visualize the distribution of exam scores in a class and identify patterns such as the central tendency and variability.

## 2. Scatter Plots:

- 1. Description:* Scatter plots show the relationship between two continuous variables.
- 2. Statistical Context:* You could use a scatter plot to visualize the correlation between study hours and exam scores. Statistical computing tools can generate scatter plots and even fit regression lines to explore the relationship quantitatively.

## 1.Box Plots:

1. *Description:* Box plots summarize the distribution of a dataset and highlight key statistics.
2. *Statistical Context:* In statistical computing, a box plot can be created to compare the distribution of salaries across different departments in a company, providing insights into central tendency, spread, and potential outliers.

## 2.Heatmaps:

1. *Description:* Heatmaps represent data in a matrix format using colors to indicate values.
2. *Statistical Context:* Statistical computing tools can generate heatmaps to visualize the correlation matrix between multiple variables, helping to identify patterns and relationships in complex datasets.

## 1. Pie Charts:

1. *Description:* Pie charts represent data as slices of a circular "pie," with each slice corresponding to a proportion of the whole.
2. *Statistical Context:* In statistical computing, a pie chart could be used to show the percentage distribution of different product categories in a sales dataset, providing a quick overview of the composition.

## 2. Line Charts:

1. *Description:* Line charts display data points connected by straight lines, illustrating trends over time or another continuous variable.
2. *Statistical Context:* Statistical computing tools can create line charts to visualize the trend in stock prices over a period, helping analysts identify patterns or anomalies.

- **Interactive Dashboards:**

- *Description:* Dashboards integrate multiple visualizations and allow users to interact with the data dynamically.
- *Statistical Context:* In a business context, statistical computing can be used to create interactive dashboards that display key performance indicators (KPIs) and allow decision-makers to explore and analyze data in real-time

- In summary, data visualization is a powerful tool in statistical computing, enabling effective communication of insights and patterns hidden within data. The choice of visualization depends on the nature of the data and the specific goals of the analysis.

# What is common probability distributions.Explain with some examples

- Probability distributions describe the likelihood of different outcomes in a random experiment or process. Common probability distributions are used in statistical modeling to represent various types of data. Here are some examples of common probability distributions:

# Examples

## **1. Normal Distribution (Gaussian Distribution):**

- 1. Description:* The normal distribution is symmetric and bell-shaped. It is characterized by its mean (average) and standard deviation.
- 2. Example:* Heights of a population, IQ scores, errors in measurements.

## **2. Binomial Distribution:**

- 1. Description:* The binomial distribution models the number of successes in a fixed number of independent and identical Bernoulli trials (experiments with two possible outcomes).
- 2. Example:* Coin flips (success: heads, failure: tails), the number of successful drug trials in a series of experiments.

## **1. Poisson Distribution:**

- 1. Description:* The Poisson distribution models the number of events occurring in a fixed interval of time or space.
- 2. Example:* Number of phone calls received at a call center in an hour, occurrences of earthquakes in a region.

## **2. Exponential Distribution:**

- 1. Description:* The exponential distribution models the time between events in a Poisson process, where events occur continuously and independently at a constant average rate.
- 2. Example:* Time between arrivals of customers at a service point, time between radioactive decay events.

## **1. Uniform Distribution:**

- 1. Description:* The uniform distribution represents outcomes that are equally likely over a range.
- 2. Example:* Rolling a fair six-sided die, selecting a random number between 0 and 1.

## **2. Gamma Distribution:**

- 1. Description:* The gamma distribution generalizes the exponential distribution and is often used to model the sum of gamma-distributed random variables.
- 2. Example:* Waiting time until the  $k$ -th arrival in a Poisson process, time until a fixed number of failures in a series of independent and identically distributed exponential random variables.

## **1. Beta Distribution:**

- 1. Description:* The beta distribution is a family of continuous probability distributions defined on the interval  $[0, 1]$ .
- 2. Example:* Modeling the distribution of proportions, success/failure experiments.

## **2. Logistic Distribution:**

- 1. Description:* The logistic distribution is an S-shaped distribution that is symmetric and similar to the normal distribution.
- 2. Example:* Modeling growth processes, logistic regression in statistics.

- **Hypergeometric Distribution:**

- *Description:* The hypergeometric distribution models the number of successes in a fixed-size sample drawn without replacement from a finite population.

- *Example:* Drawing a certain number of red balls from an urn containing both red and blue balls.

# What is common probability mass functions. Explain with examples

- A probability mass function (PMF) is a function that describes the probability of discrete random variables taking on specific values. Here are some common probability mass functions along with examples:

# Examples

## 1. Bernoulli Distribution:

1. *PMF:*  $P(X=k) = p^k \cdot (1-p)^{1-k}$ , where  $k$  is the outcome (0 or 1).
2. *Example:* A single coin flip (success: heads with probability  $p$ , failure: tails with probability  $1-p$ ).

## 2. Binomial Distribution:

1. *PMF:*  $P(X=k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$ , where  $n$  is the number of trials and  $k$  is the number of successes.
2. *Example:* The number of heads in  $n$  coin flips.

## 3. Poisson Distribution:

1. *PMF:*  $P(X=k) = \frac{e^{-\lambda} \cdot \lambda^k}{k!}$ , where  $\lambda$  is the average rate of events.
2. *Example:* The number of emails received in an hour with an average rate of  $\lambda$  emails per hour.

## 1. Geometric Distribution:

1. *PMF:*  $P(X=k) = (1-p)^{k-1} \cdot p$ , where  $k$  is the number of trials until the first success.

2. *Example:* The number of coin flips needed to get the first head.

## 2. Hypergeometric Distribution:

1. *PMF:*  $P(X=k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$ , where  $N$  is the population size,  $K$  is the number of successes in the population, and  $n$  is the sample size.

2. *Example:* Drawing  $n$  cards without replacement from a deck of  $N$  cards, where  $K$  cards are of a specific type.

## 1. Uniform Distribution:

1. *PMF:*  $P(X=k) = \frac{1}{b-a+1}$ , where  $k$  can take on any integer value in the range  $[a, b]$ .
2. *Example:* Rolling a fair six-sided die (each outcome has an equal probability of  $1/6$ ).

## 2. Multinomial Distribution:

1. *PMF:*  
$$P(X_1=k_1, X_2=k_2, \dots, X_r=k_r) = \frac{n!}{k_1! \cdot k_2! \cdot \dots \cdot k_r!} \cdot p_1^{k_1} \cdot p_2^{k_2} \cdot \dots \cdot p_r^{k_r}$$
, where  $n$  is the number of trials,  $r$  is the number of categories,  $k_i$  is the number of occurrences in category  $i$ , and  $p_i$  is the probability of category  $i$ .
2. *Example:* Rolling a fair six-sided die  $n$  times and counting the occurrences of each number.

# What is Bernoulli Distribution

- The Bernoulli distribution is a discrete probability distribution that models a random experiment with two possible outcomes, often referred to as "success" and "failure." It is named after the Swiss mathematician Jacob Bernoulli. The outcomes are typically denoted as 0 and 1, where 1 represents success and 0 represents failure.

# Key characteristics of the Bernoulli distribution:

1. **Parameter:** The distribution is characterized by a single parameter  $p$ , which represents the probability of success.

2. **Probability Mass Function (PMF):**

3.  $p$  & \text{if } k = 1 \ \ 1 - p & \text{if } k = 0 \ \end{cases}

4. where  $k$  is the outcome (0 or 1)

• **Mean and Variance:** Mean( $X$ )= $p$  Mean( $\mu$ )= $p$   
Variance( $X^2$ )= $p \cdot (1 - p)$  Variance( $\sigma^2$ )= $p \cdot (1 - p)$

# Example:

- Consider a single coin flip, where success (1) is getting heads and failure (0) is getting tails. If the probability of getting heads is  $p=0.6$ , then the Bernoulli distribution for this experiment can be expressed as:
  - $0.6 & \text{if } k = 1 \ \ 0.4 & \text{if } k = 0 \ \ \text{end{cases}}$
- This means that the probability of success (getting heads) is 0.6, and the probability of failure (getting tails) is 0.4.

# Summary

- The Bernoulli distribution serves as the foundation for more complex distributions, such as the binomial distribution, which models the number of successes in a fixed number of independent Bernoulli trials. Additionally, it is a fundamental distribution in probability theory and is often used in various statistical applications, especially when dealing with binary outcomes.

# What is Binomial Distribution .Explain with examples

- The binomial distribution is a discrete probability distribution that models the number of successes in a fixed number of independent and identical Bernoulli trials, where each trial has only two possible outcomes: success or failure. The outcomes are often denoted as 1 for success and 0 for failure. The trials are independent, meaning the outcome of one trial does not affect the outcome of another, and they are identical, meaning each trial has the same probability of success.

# Key characteristics of the Binomial Distribution:

## 1. Parameters:

1.  $n$ : The number of trials.
2.  $p$ : The probability of success on a single trial.

## 2. Probability Mass Function (PMF):

$P(X=k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$   
where  $k$  is the number of successes.

**3. Mean and Variance:** Mean( $\mu$ ) =  $np$

Variance( $\sigma^2$ ) =  $np(1-p)$

# Example:

- Suppose you flip a fair coin (where the probability of heads,  $p$ , is 0.5) five times. The number of heads you get in these five flips follows a binomial distribution.
- **Parameters:**  $n=5$  (number of coin flips),  $p=0.5$  (probability of heads on a single flip).
- **Probability Mass Function:**  
 $P(X=k) = \binom{5}{k} \cdot 0.5^k \cdot 0.5^{5-k}$  for  $k=0, 1, 2, 3, 4, 5$ .
- **Example Calculations:**
  - $P(X=3) = \binom{5}{3} \cdot 0.5^3 \cdot 0.5^2 = 10 \cdot 0.125 \cdot 0.25 = 0.3125$
  - $P(X \leq 2) = P(X=0) + P(X=1) + P(X=2)$

- This means that the probability of getting exactly 3 heads in 5 coin flips is 0.3125, and the probability of getting 2 or fewer heads is the sum of the probabilities for  $k=0,1,2$ .
- The binomial distribution is widely used in statistics and probability theory to model scenarios with binary outcomes and a fixed number of trials, such as the number of successful trials in a series of independent experiments.

# What is Poisson distributions. Explain with examples

- The Poisson distribution is a probability distribution that describes the number of events that will occur within a fixed interval of time or space, given the average number of occurrences in that interval. It is named after the French mathematician Siméon Denis Poisson.
- The Poisson distribution is characterized by a single parameter, often denoted as  $\lambda$  (lambda), which represents the average rate of occurrence of the event

# The probability mass function (PMF) of the Poisson distribution is given by:

- $P(X=k) = \frac{\lambda^k}{k!} e^{-\lambda}$
- where:
- $P(X=k)$  is the probability of observing  $k$  events,
- $e$  is the base of the natural logarithm (approximately 2.71828),
- $\lambda$  is the average rate of events,
- $k$  is a non-negative integer representing the number of events.

- Here are a couple of examples to illustrate the concept:
- **Example 1:**
- Suppose you are observing the number of emails you receive per hour, and on average, you receive 5 emails per hour. The distribution of the number of emails per hour can be modeled as a Poisson distribution with  $\lambda = 5$ . You can then calculate the probability of receiving a specific number of emails in an hour using the Poisson distribution.

- **Example 2:**

- Imagine you are monitoring a radioactive substance, and on average, you observe 2 decay events per minute. The number of decay events in a minute can be modeled as a Poisson distribution with  $\mu = 2\lambda = 2$ . You can use the Poisson distribution to calculate the probability of observing a particular number of decay events in a given minute.
- In both examples, the Poisson distribution helps model situations where events occur independently, and the average rate of occurrence is known.

# what is common probability density functions in statistical computing

- In statistical computing, several common probability density functions (PDFs) are frequently used to model the distribution of random variables. Here are some of the most common ones:

## 1. Normal (Gaussian) Distribution:

1. Parameters: Mean ( $\mu$ ), Standard Deviation ( $\sigma$ )
2. PDF:  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$
3. The normal distribution is widely used due to the Central Limit Theorem, making it applicable in many real-world scenarios.

## 2. Binomial Distribution:

1. Parameters: Number of trials ( $n$ ), Probability of success ( $p$ )
2. PMF:  $P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$
3. Models the number of successes in a fixed number of independent Bernoulli trials.

## 1. Poisson Distribution:

1. Parameter: Average rate ( $\lambda$ )
2. PMF:  $P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!}$
3. Models the number of events occurring in a fixed interval of time or space.

## 2. Exponential Distribution:

1. Parameter: Rate ( $\lambda$ )
2. PDF:  $f(x) = \lambda e^{-\lambda x}$
3. Describes the time until an event of interest occurs in a Poisson process.

## 3. Uniform Distribution:

1. Parameters: Minimum ( $a$ ), Maximum ( $b$ )
2. PDF:  $f(x) = \frac{1}{b-a}$
3. All values within the range are equally likely.

## 1. Gamma Distribution:

1. Parameters: Shape ( $k$ ), Rate ( $\theta$ )
2. PDF:  $f(x; k, \theta) = \frac{\theta^k}{\Gamma(k)} x^{k-1} e^{-\theta x}$
3. Generalizes the exponential distribution and is used for modeling waiting times.

## 2. Beta Distribution:

1. Parameters: Shape parameters ( $\alpha, \beta$ )
2. PDF:  $f(x; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$
3. Often used as a prior distribution in Bayesian statistics.

## 3. Log-Normal Distribution:

1. Parameters: Location ( $\mu$ ), Scale ( $\sigma$ )
2. PDF:  $f(x; \mu, \sigma) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{(\ln(x) - \mu)^2}{2\sigma^2}}$
3. Models the distribution of a random variable whose logarithm is normally distributed.

# PURPOSE

- These distributions are foundational in statistical modeling and are extensively used in various fields for data analysis, simulation, and hypothesis testing.

# Certainly! Let's consider some real-world examples where these probability density functions (PDFs) might be applied:

## 1. Normal Distribution:

1. **Example:** The heights of a population. Heights tend to follow a normal distribution, with most people clustered around the average height.

## 2. Binomial Distribution:

1. **Example:** Flipping a biased coin. If you have a coin with a known probability of landing heads (success), the number of heads in a fixed number of flips follows a binomial distribution.

## 3. Poisson Distribution:

1. **Example:** The number of customer arrivals at a service desk in a given time period, assuming arrivals are random and independent.

## 4. Exponential Distribution:

1. **Example:** The time between arrivals of events in a Poisson process, such as the time between phone calls at a call center.

## **1. Uniform Distribution:**

**1. Example:** Randomly selecting a number between 1 and 6 with a fair six-sided die.

## **2. Gamma Distribution:**

**1. Example:** The time until a light bulb burns out, assuming the bulb's lifetime follows a gamma distribution.

## **3. Beta Distribution:**

**1. Example:** Modeling the distribution of the probability of success in a Bernoulli trial based on prior knowledge.

## **4. Log-Normal Distribution:**

**1. Example:** Stock prices. The log-normal distribution is often used to model the distribution of stock prices.

# UNIFORM DISTRIBUTIONS

- A uniform distribution, also known as a rectangular distribution, is a probability distribution where all values within a specified range are equally likely to occur. In other words, every outcome in the range has an equal probability of occurring. The probability density function (PDF) of a continuous uniform distribution is constant within the interval and zero outside the interval.
- For a continuous uniform distribution on the interval  $[a, b]$ , the probability density function  $f(x)$  is given by:

- For a continuous uniform distribution on the interval  $[a, b]$ , the probability density function  $f(x)$  is given by:
- $f(x) = \frac{1}{b-a}$
- where:
- $a$  is the minimum value in the interval,
- $b$  is the maximum value in the interval,
- $b > a$ .

- The probability density function ensures that the total area under the curve is equal to 1, reflecting the fact that the entire probability space is covered within the specified interval.
- For a discrete uniform distribution, where each value in a finite set is equally likely, the probability mass function (PMF) is given by:
- $\sum_{x \in S} P(X=x) = 1$   $P(X=x) = \frac{1}{n}$   
where:
- $n$  is the number of distinct values in the set.

# Example:

- Suppose you roll a fair six-sided die. The outcome of each roll follows a discrete uniform distribution, where each number (1, 2, 3, 4, 5, 6) has an equal probability of  $\frac{1}{6}$ . If you represent the random variable  $X$  as the result of the die roll, the PMF is:
- $P(X=x) = \frac{1}{6}$
- This means that each number on the die has an equal chance of being rolled, and the probability distribution is uniform.

# What is normal distribution, Explain with examples

- A normal distribution, also known as a Gaussian distribution, is a continuous probability distribution that is symmetric and bell-shaped. It is characterized by two parameters: the mean ( $\mu$ ), which represents the center of the distribution, and the standard deviation ( $\sigma$ ), which measures the spread or dispersion of the distribution.

# The probability density function (PDF) of the normal distribution is given by:

- $f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$
- The normal distribution is often denoted as  $N(\mu, \sigma^2)$ , where  $\sigma^2$  is the variance.
- Key properties of the normal distribution include:
  1. It is symmetric around the mean ( $\mu$ ).
  2. About 68% of the data falls within one standard deviation of the mean, 95% within two standard deviations, and 99.7% within three standard deviations (empirical rule or 68-95-99.7 rule).
  3. The mean, median, and mode are all equal and located at the center of the distribution.

# Example:

- Suppose we have a population of adult males and we are interested in their heights. If the heights of these individuals follow a normal distribution with a mean ( $\mu$ ) of 175 cm and a standard deviation ( $\sigma$ ) of 8 cm, we can use the normal distribution to make various probability statements about the heights.

## 1. Probability of a Random Male Being Taller than 180 cm:

$P(X > 180) = \int_{180}^{\infty} f(x; 175, 8) dx$  This calculates the probability that a randomly selected male is taller than 180 cm.

## 2. Range of Heights Containing 95% of the Population: According to the empirical rule, about 95% of the heights fall within

$\pm 2\sigma$ . So, the range would be  $[175 - 2(8), 175 + 2(8)] = [159, 191]$  cm.

# What is student's t-distribution. Explain with examples

- Student's t-distribution, often simply referred to as the t-distribution, is a probability distribution that arises in the context of statistical inference when the sample size is small, and the population standard deviation is unknown. It is used in situations where the underlying population follows a normal distribution, but the sample size is not large enough for the normal distribution to be applied.

The t-distribution is characterized by a parameter known as the degrees of freedom ( $df$ ). The probability density function (PDF) of the t-distribution is given by:

- $f(t, df) = \frac{\Gamma(df)}{\sqrt{\pi} \Gamma(df - 1/2)} (1 + \frac{t^2}{df})^{-df}$
- Here,  $\Gamma$  denotes the gamma function.
- As the degrees of freedom increase, the t-distribution approaches the standard normal distribution.

# Example:

- Let's consider an example where the t-distribution is commonly used: confidence intervals for the mean.
- Suppose you have a small sample (e.g.,  $n=15$ ) from a population with an unknown standard deviation. You want to estimate the population mean and provide a confidence interval. In this case, you would use the t-distribution because the sample size is small.

- In summary, the t-distribution is particularly useful when dealing with small sample sizes and provides wider intervals compared to the normal distribution due to the increased uncertainty associated with estimating the population standard deviation from a small sample.