Open in app ↗

Search

# 12 Probability Distributions in Data Science
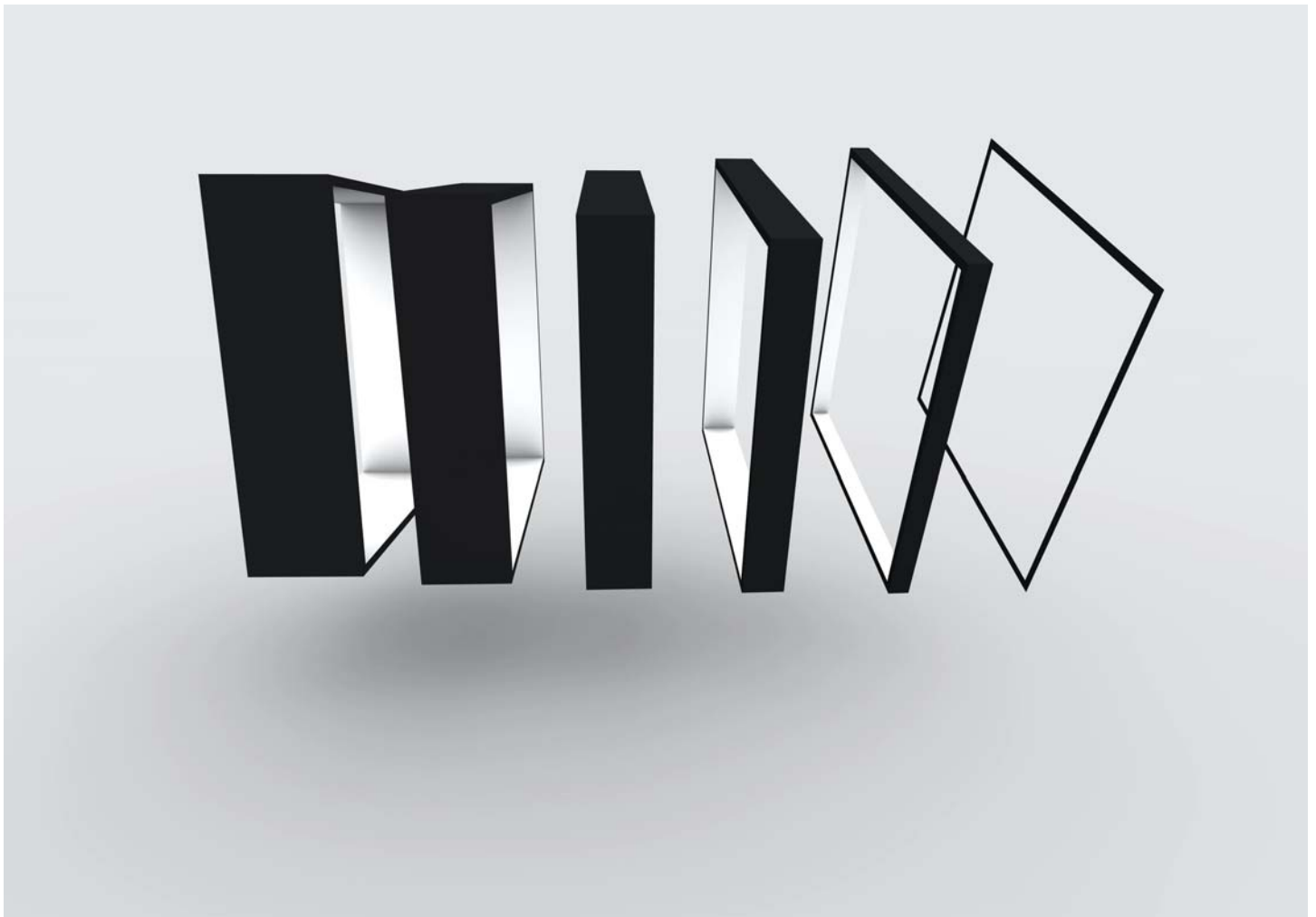
btd · Follow

13 min read · Nov 11, 2023

▷ Listen        ⬆ Share        ••• More



Photo by Mario Verduzco on Unsplash

In the context of data science, distributions refer to the patterns by which data values are spread across a dataset. Understanding the underlying distribution of data is crucial for various statistical analyses and modeling techniques.

Here are some important probability distributions frequently encountered in data science:

## 1. Normal Distribution (Gaussian Distribution)

The normal distribution, also known as the Gaussian distribution, is a fundamental and widely used probability distribution.

- Symmetrical, bell-shaped curve characterized by its mean (center) and standard deviation (spread).

- The mean ( $\mu$ ) is the center of the distribution, and the standard deviation ( $\sigma$ ) controls the spread or dispersion of the distribution.

- About 68% of the data falls within one standard deviation of the mean, 95% within two standard deviations, and approximately `99.7%` within three standard deviations.

- Empirical Rule (68–95–99.7 Rule) states that for a normal distribution, approximately `68%` of the data falls within one standard deviation of the mean, `95%` within two standard deviations, and `99.7%` within three standard deviations.

- The probability density function (PDF) of the normal distribution forms a bell-shaped curve. This shape is characterized by a single peak at the mean, and it gradually tails off towards positive and negative infinity.

- If a normal distribution has a mean of `0` and a standard deviation of `1`, it is referred to as a standard normal distribution.

- Z-scores can be used to standardize values from any normal distribution to the standard normal distribution.

- Widely used in statistical inference, hypothesis testing, and modeling natural phenomena.

## 2. Bernoulli Distribution

The Bernoulli distribution is a simple and fundamental discrete probability distribution that models a random experiment with two possible outcomes, often referred to as "success" and "failure."

- The experiment associated with a Bernoulli distribution has only two possible outcomes, often denoted as `0` and `1`, where 1 typically represents "success" and 0 represents "failure."

- Examples include a coin flip (heads or tails), the success or failure of a single trial in an experiment, or the outcome of a single question with a yes/no answer.

- The mean (expected value) of a Bernoulli distribution is $E[X] = p$.

- The variance of a Bernoulli distribution is $Var[X] = p(1-p)$.

- The Bernoulli distribution serves as the basis for more complex probability distributions, such as the binomial distribution. When multiple independent Bernoulli trials are conducted, the sum of these trials follows a binomial distribution.

- The Bernoulli distribution is a fundamental building block in probability theory and has applications in various fields, including statistics, machine learning, and decision theory.

### 3. Binomial Distribution

The binomial distribution is a discrete probability distribution that describes the number of successes in a fixed number of independent and identical Bernoulli trials. Each trial is characterized by two possible outcomes, often labeled as "success" and "failure."

- Discrete distribution representing the number of successes in a fixed number of independent Bernoulli trials.

- The distribution models the number of successes in a fixed number of independent trials, denoted by $n$.

- Each trial in a binomial distribution has only two possible outcomes, typically denoted as success (S) and failure (F).

- The trials are assumed to be independent, meaning the outcome of one trial does not affect the outcome of another.

- The probability of success on a single trial is denoted by $p$.

- The mean (expected value) of a binomial distribution is $E[X] = np$.

- The variance of a binomial distribution is $Var[X] = np(1-p)$.

- The binomial distribution is a fundamental concept in probability theory and statistics, and it serves as the basis for the development of more advanced

probability distributions.

- The binomial distribution is widely used to model real-world scenarios where there are two possible outcomes, such as success/failure experiments, the number of heads in a fixed number of coin flips, or defective/non-defective items in a production process.

## 4. Poisson Distribution

The Poisson distribution is a discrete probability distribution that models the number of events occurring in a fixed interval of time or space. It is particularly useful for situations where events are rare and random.

- Discrete distribution representing the number of events occurring in a fixed interval of time or space, denoted by $t$.

- Modeling rare events, such as the number of phone calls at a call center in a given minute.

- Events are assumed to occur independently of each other within the given interval.

- The Poisson distribution is characterized by a single parameter, $\lambda$, which represents the average rate of events per unit interval.

- The mean (expected value) of a Poisson distribution is $E[X] = \lambda$.

- The variance of a Poisson distribution is $Var[X] = \lambda$.

- The Poisson distribution is commonly used in various fields, such as telecommunications (modeling the number of phone calls in a given minute), biology (modeling the number of mutations in a DNA strand), and queuing theory (modeling the number of arrivals at a service point).

- In certain conditions, the Poisson distribution can be approximated by the binomial distribution when the number of trials ($n$) is large, and the probability of success ($p$) is small, with $\lambda = np$.

- The Poisson distribution is a valuable tool for analyzing and predicting the occurrence of rare events in situations where the events are independent and the average rate of occurrence is known.

## 5. Exponential Distribution

The exponential distribution is a continuous probability distribution that models the time between events in a Poisson process. It is often used to describe the waiting time until the occurrence of the next event in a sequence of independent events that happen at a constant rate.

- Continuous distribution representing the time between events in a Poisson process.

- Modeling the time until an event occurs, such as the time between arrivals at a service point.

- One of the defining features of the exponential distribution is the memoryless property. This means that the probability of an event occurring in the next instant is independent of the past. In other words, the distribution has no "memory" of previous events. `P(T > s+t | T > s) = P(T > t)`.

- The exponential distribution is characterized by a single parameter, $\lambda$, which is the rate parameter. It represents the average number of events per unit of time.

- The mean (expected value) of an exponential distribution is `E[T] = 1/`$\lambda$.

- The variance of an exponential distribution is `Var[T] = 1/`$\lambda^2$.

- If `N(t)` is the number of events in the time interval `[0,t]` following a Poisson process with rate $\lambda$, then `N(t)` follows a Poisson distribution with mean `λt`. The waiting time until the first event follows an exponential distribution with rate $\lambda$.

- The exponential distribution is a valuable tool for modeling and analyzing random processes where events occur independently at a constant rate.

- The exponential distribution is commonly used in reliability engineering, queuing theory, and telecommunications to model the time until the next event occurs. For example, it can represent the time between arrivals at a service point, the time until a radioactive decay event, or the time until a system failure.

## 6. Uniform Distribution

The uniform distribution is a probability distribution where all values in the distribution have an equal probability of occurring. In other words, each outcome within the range of possible values has the same likelihood of being observed.

- All values in the distribution have an equal probability of occurring.

- The probability density function (PDF) is constant over the entire range.

- Used in scenarios where all outcomes are equally likely.

- There are continuous and discrete versions of the uniform distribution. The continuous uniform distribution applies to a continuous range of values, while the discrete uniform distribution is defined for a finite set of discrete values.

- The uniform distribution is commonly used in scenarios where all outcomes are equally likely. Examples include random number generation, modeling uncertainty when the true distribution is unknown, and certain types of simulations.

- The uniform distribution is often used as a basis for generating random numbers with equal likelihood within a specified range. These random numbers can be transformed to simulate other distributions.

- In a continuous uniform distribution, the area under the probability density function is a rectangle, and the total area is equal to `1`. The height of the rectangle is determined by the reciprocal of the width of the interval.

- The uniform distribution is a simple and intuitive model for situations where there is no reason to believe that one outcome is more likely than another within a given range.

## 7. Gamma Distribution

The gamma distribution is a continuous probability distribution that serves as a generalization of the exponential distribution. While the exponential distribution models the time until the first event in a Poisson process (memoryless property), the gamma distribution extends this concept to model the time until the `k-th` event in the same process. The gamma distribution is particularly useful in queueing theory and reliability engineering.

- A generalization of the exponential distribution, often used to model waiting times.

- The gamma distribution is defined by two parameters: `k`, the shape parameter, and `θ`, the scale parameter.

- The shape parameter determines the number of events, and the scale parameter influences the rate of occurrence.

- When $k=1$, the gamma distribution becomes the exponential distribution.

- When $k$ is a positive integer, the gamma distribution is also known as the Erlang distribution.

- The mean (expected value) of a gamma distribution is $E[X] = k\theta$.

- The variance of a gamma distribution is $Var[X] = k\theta^2$.

- In queueing theory, the gamma distribution is often used to model the waiting time until a specified number of events (such as customers arriving at a service point) occur.

- In reliability engineering, the gamma distribution is applied to model the time until $k$ failures in a system.

- The gamma distribution can be expressed as the sum of $k$ independent and identically distributed exponential random variables with rate parameter $1/\theta$.

- If $N(t)$ is the number of events in a Poisson process with rate $\lambda$, then $Tk$, the time until the $k$-th event, follows a gamma distribution with shape parameter $k$ and scale parameter $1/\lambda$.

- The gamma distribution provides flexibility in modeling waiting times and time until a specified number of events in various fields, making it a valuable tool in probability theory and statistics.

### 8. Beta Distribution

The beta distribution is a continuous probability distribution defined on the interval $[0,1][0,1]$. It is a versatile distribution commonly used in Bayesian statistics to model random variables representing proportions or probabilities. The beta distribution is parameterized by two positive shape parameters, denoted by $\alpha$ and $\beta$.

- The shape parameters $\alpha$ and $\beta$ determine the shape of the distribution. Higher values of $\alpha$ bias the distribution towards higher values, while higher values of $\beta$ bias it towards lower values.

- When both $\alpha$ and $\beta$ are set to 1, the beta distribution becomes a uniform distribution over `[0,1][0,1]` .

- The mean (expected value) of a beta distribution is `E[X]= α / (α+β)` .

- The variance of a beta distribution is `Var[X]= αβ / ((α+β)² (α+β+1))` .

- The beta distribution is often employed to model proportions and probabilities. It is suitable for situations where the random variable represents the success probability in a binary outcome.

- In Bayesian statistics, the parameters $\alpha$ and $\beta$ of the beta distribution can be updated as new data is observed, allowing for Bayesian estimation and updating of beliefs.

- The beta distribution is widely used in Bayesian statistics as a conjugate prior for the binomial distribution. This means that if you have a beta-distributed prior and you observe data from a binomial distribution, the posterior distribution is also beta-distributed.

- The beta distribution is a flexible and powerful tool in Bayesian statistics, providing a natural way to model uncertainty about probabilities and proportions. Its ability to represent a wide range of shapes makes it applicable in various statistical and probabilistic scenarios.

### 9. Logistic Distribution

The logistic distribution is a continuous probability distribution that resembles the normal distribution but has heavier tails. It is commonly used in statistics and machine learning, especially in logistic regression, where it serves as the cumulative distribution function (CDF) for modeling binary outcomes.

- Continuous distribution resembling the normal distribution but with heavier tails.

- The logistic distribution is symmetric around its mean ( $\mu$ ) and has heavier tails compared to the normal distribution.

- Logistic regression, modeling binary outcomes.

- The mean (expected value) of a logistic distribution is `E[X] = μ`.

- The variance of a logistic distribution is `Var[X] = (s² π²) / 3`.

- The logistic distribution has an S-shape curve, making it suitable for modeling probabilities and binary outcomes.

- In logistic regression, the logistic distribution is used as the CDF to model the probability of a binary outcome. The logistic regression model transforms a linear combination of predictor variables using the logistic function.

- The logistic distribution is particularly useful for modeling binary outcomes, such as success/failure or yes/no.

- The tails of the logistic distribution are heavier than those of the normal distribution. This can be advantageous in modeling situations where extreme values are more likely to occur.

- The logistic distribution is related to the logistic function, which is the link function used in logistic regression. The logistic function maps real-valued numbers to the range `(0, 1)`, making it suitable for modeling probabilities.

- The logistic distribution is a key component in logistic regression, a widely used statistical method for binary classification problems. Its shape and properties make it well-suited for modeling situations where outcomes are binary and have a sigmoidal relationship with predictor variables.

## 10. Cauchy Distribution

The Cauchy distribution is a continuous probability distribution known for its heavy tails and the property that both its mean and variance are undefined. It is characterized by a symmetric bell-shaped curve, but its tails extend indefinitely.

- Heavy-tailed distribution with undefined mean and variance.

- The Cauchy distribution has undefined mean and variance. The integral of the distribution diverges, leading to the lack of finite moments.

- One of the defining characteristics of the Cauchy distribution is its heavy tails. This means that extreme values are more likely to occur compared to distributions with finite variance.

- The Cauchy distribution is symmetric around its median ( $x0$ ). This symmetry holds regardless of the values of the location and scale parameters.

- The Cauchy distribution is associated with the Cauchy principal value, which is used in physics and engineering to handle integrals that converge conditionally.

- Unlike many other distributions, the Cauchy distribution lacks moment stability due to its undefined mean and variance. This can make it challenging to work with in certain statistical analyses.

- The Cauchy distribution does not satisfy the conditions of the Central Limit Theorem, as its variance is undefined. Therefore, the sample mean of a Cauchy-distributed sample may not converge to a normal distribution as sample size increases.

- The Cauchy distribution is used in physics and engineering, particularly in signal processing and resonance phenomena. It can arise in problems involving natural frequencies, resonances, and phase transitions.

- While the Cauchy distribution has interesting mathematical properties and applications in specific domains, its lack of mean and variance can make it less suitable for certain statistical analyses where these moments are crucial. Researchers often prefer more well-behaved distributions when working with real-world data.

## 11. Chi-squared Distribution

The chi-squared distribution is a continuous probability distribution that arises in the context of hypothesis testing. It is a special case of the gamma distribution and is commonly used in statistical inference, particularly for conducting tests related to goodness-of-fit and independence in contingency tables.

- Continuous distribution that arises in the context of hypothesis testing.

- The chi-squared distribution is characterized by a parameter $k$, which represents the degrees of freedom. The degrees of freedom determine the shape of the distribution.

- The chi-squared distribution is central to hypothesis testing, especially in situations involving categorical data. It is often used in tests of independence in contingency tables and goodness-of-fit tests.

- In a goodness-of-fit test, the chi-squared distribution is used to assess whether an observed frequency distribution differs significantly from a theoretical (expected) distribution.

- In tests of independence, the chi-squared distribution is employed to examine whether there is a significant association between two categorical variables.

- Contingency tables are commonly used in the context of chi-squared tests. These tables summarize the joint distribution of two or more categorical variables.

- For large degrees of freedom, the chi-squared distribution approximates a normal distribution, allowing for the application of the central limit theorem in large-sample situations.

- The chi-squared distribution is also related to multivariate statistics, where it is used in Wilk's Lambda and Hotelling's $T2$ tests.

- The chi-squared distribution is a fundamental tool in statistics for comparing observed and expected frequencies in categorical data and assessing the independence of categorical variables. It plays a crucial role in hypothesis testing and is widely applied in various fields, including biology, social sciences, and epidemiology.

## 12. Student's t-Distribution

Student's t-distribution, often simply referred to as the t-distribution, is a continuous probability distribution used in statistical inference, especially when dealing with small sample sizes and situations where the population variance is unknown. It is a fundamental distribution for conducting t-tests and constructing confidence intervals.

- Continuous distribution used for statistical inference when the sample size is small and population variance is unknown.

- The t-distribution is characterized by a parameter $n$, which represents the degrees of freedom. The degrees of freedom determine the shape of the distribution.

- The t-distribution is central to t-tests, which are statistical tests used to compare means of two groups or to test the hypothesis about a single population mean. Common t-tests include the one-sample t-test, two-sample t-test, and paired-sample t-test.

- The t-distribution is used to construct confidence intervals for the population mean. The width of the interval depends on the sample size, and the uncertainty

associated with estimating the population variance from a small sample is accounted for by using the t-distribution.

- The t-distribution is bell-shaped and symmetric, similar to the normal distribution. However, it has heavier tails, especially for smaller sample sizes. As the sample size increases, the t-distribution approaches the normal distribution.

- As the degrees of freedom decrease (indicating a smaller sample size), the t-distribution has fatter tails, resulting in wider confidence intervals compared to intervals constructed using the normal distribution.

- Critical values for hypothesis testing and constructing confidence intervals using the t-distribution are often obtained from t-tables, which provide values for various degrees of freedom and significance levels.

- The t-distribution is a crucial tool in statistics, especially when working with small sample sizes. It provides a more accurate model for the distribution of sample means when the population variance is unknown. As the sample size increases, the t-distribution approaches the normal distribution, and the distinction becomes less important for large samples.

These distributions play a crucial role in various statistical analyses, hypothesis testing, and machine learning models. Depending on the nature of the data and the problem at hand, different distributions may be more appropriate for modeling and analysis.

Probability Distributions        Data Science        Statistical Analysis        Hypothesis Testing

Machine Learning

Follow

# Written by btd

764 Followers

Learning & making lists

---

## More from btd



 btd

## 13 Statistical Analysis Methods for Data Analysts & Data Scientists

Statistical analysis techniques encompass a wide range of methods used to analyze data, make inferences, and draw conclusions about...

✦ · 14 min read · Nov 9, 2023

👏 217        💬                                          🔖⁺        •••

---