

Mathematical Statistics I

Chapter 2: Probability

Jesse Wheeler

Contents

1	Introduction	1
2	Discrete Random Variables	2
2.1	Common random variables	3
2.1.1	Bernoulli Random Variables	3
2.1.2	Binomial Distribution	4
2.1.3	Geometric and Negative Binomial Distributions	5
2.1.4	Hypergeometric and Poisson Distribution	6
3	Continuous Random Variables	9
3.1	The Uniform Random Variable	10
3.2	The cumulative distribution function	10
3.3	Common random variables	12
3.3.1	The Exponential Distribution	12
3.3.2	The Gamma Density	13
3.3.3	The Normal Distribution	13
3.3.4	The Beta Distribution	15
4	Functions of a random variable	16

1 Introduction

Introduction

- Formally, a *random variable* is a function from a sample space Ω to the real numbers¹.
- That is, for any element $\omega \in \Omega$, a random variable X will map ω to a real number: $X(\omega) \in \mathbb{R}$.
- Most often people think of random variables as random numbers rather than functions; in most instances in this class, this treatment will be sufficient.

Example of a random variable

Consider the experiment of flipping three coins. The sample space is

$$\Omega = \{hhh, hht, hth, thh, htt, tth, ttt\}.$$

- Some possible random variables include (1) the number of heads, (2) the number of tails, (3) the number of heads minus the number of tails.

¹In this class, will assume real-valued spaces, though more generally a random variable can map to any measurable space

- Importantly, a random variable must assign a value to all possible outcomes $\omega \in \Omega$.

Number of Heads

Let X be the random variable representing the number of heads. If the result of the outcome is the event hth , the $X(\{hth\}) = 2$.

A few comments on random variables

- Sometimes in this course I will use the abbreviation RV to mean “random variable”, and you can do so as well.
- It is conventional to use uppercase letters (math text or italics) to denote random variables.
- While a random variable is a function, the outcome of an experiment $\omega \in \Omega$ is random (that’s the point), and we only ever see a single outcome. Thus, the fact that X is a function is often dropped, and we just write X . The realized value of X is random, because the input is random.

2 Discrete Random Variables

Discrete Random Variables

Definition: Discrete random variable

A discrete random variable is a random variable that can take on only a finite or at most a countably infinite number of values.

- Example: The number of heads in three coin flips can only be in the set $\{0, 1, 2, 3\}$. Alternatively, consider flipping a coin indefinitely until you achieve a heads. The possible outcomes are in the set $\{1, 2, 3, \dots\}$, which is countably infinite.

Probabilities

- The probability measure on the sample space determines the probability of the values of X .
- In our example, if a coin is fair, then we can assign a uniform probability measure on the sample set of flipping a coin three times.
- That is, all outcomes are equally likely, each with probability $1/8$.
- The probability that X takes on it’s potential values is easily computed, by counting the number of outcomes that result in the particular value of X :

$$\begin{aligned} P(X = 0) &= \frac{1}{8} \\ P(X = 1) &= \frac{3}{8} \\ P(X = 2) &= \frac{3}{8} \\ P(X = 3) &= \frac{1}{8}. \end{aligned}$$

- More generally, let’s assume that X is a discrete RV, and denote the possible values as x_1, x_2, \dots . There exists a function p such that $p(x_i) = P(X = x_i)$ that satisfies $\sum_i p(x_i) = 1$. This function p is called the *probability mass function* (pmf) of the random variable X .

- We may also be interested in calculating for all values $x \in \mathbb{R}$, the probability $F(x) = P(X \leq x)$; the function F is called the *cumulative distribution function* (cdf). The cdf plays a number of important roles in probability and statistics that we will see later on.

Some notes:

- The cdf is non-decreasing (see Theorem 1.2), and

$$\lim_{x \rightarrow -\infty} F(x) = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} F(x) = 1.$$

- The pmf and cdf are connected: the cdf “jumps” at all values that the pdf $p(x) > 0$.
- Conventionally, the pmf is usually denoted with lower-case letters (e.g., p , f), whereas the cdf is usually denoted with upper-case letters (e.g., F).

See Figures 2.1 and 2.2 of Rice (2007) for a depiction of the pmf and cdf of the 3-coin example.

Independence

- Jumping ahead a little bit, we will define what it means for random variables to be independent (a chapter 3 topic).

Definition: Independent random variables

Let X and Y be discrete random variables defined on the same probability space, taking values x_1, x_2, \dots and y_1, y_2, \dots , respectively. X and Y are said to be independent if, for all i, j ,

$$P(X = x_i, Y = y_j) = P(X = x_i)P(Y = y_j).$$

- This definition follows very similarly to that of independent events. We can also extend this definition to *mutual independence* of many variables if the probabilities of all combinations of variables can be factored.

2.1 Common random variables

2.1.1 Bernoulli Random Variables

Bernoulli Random Variables

- A Bernoulli RV only takes on two values², 0 and 1, with probabilities $1 - p$ and p , respectively. The pmf is therefore

$$\begin{aligned} p(1) &= p \\ p(0) &= 1 - p \\ p(x) &= 0, \quad \text{if } x \neq 0 \text{ and } x \neq 1. \end{aligned}$$

- By using the output of 0 and 1, the pmf is usually written in a more compact form:

$$p(x) = \begin{cases} p^x(1-p)^{1-x}, & \text{if } x = 0 \text{ or } x = 1, \\ 0 & \text{otherwise} \end{cases}$$

²Sometimes you'll see the random variable take values -1 and 1 .

Indicator functions

- A common instance of a Bernoulli RV is an *indicator random variable*. Let I_A be the random variable that takes on the value of 1 if the event $A \subset \Omega$ occurs, and 0 otherwise:

$$I_A(\omega) = \begin{cases} 1 & \omega \in A \\ 0 & \text{otherwise} \end{cases}$$

- Here, we see that $P(I_A = 1) = P(A)$.

2.1.2 Binomial Distribution

Binomial Distribution

- Using what we know about independent RVs and Bernoulli RVs, we can derive the pmf for a Binomial distribution.
- Suppose that we have n independent experiments, where n is a fixed (positive) integer. Let each experiment have two outcomes with probabilities p and $1 - p$, respectively, which we call “success” or “failure”. We are interested now in the random variable X , the number of “successes” in n independent trials.
- *Question:* What is the probability that $X = k$, for some $k \in \{0, 1, 2, \dots\}$?

Solution Sketch:

- For $X = k$, we must have *exactly* k successes and $n - k$ failures. By the multiplication law, any one such sequence has probability $p^k(1 - p)^{n-k}$ (for instance suppose $n = 3$. What is the probability of the event SFS ? it's $p \times (1 - p) \times p = p^2(1 - p)^{3-2}$.)
- Because we only care about the *number* of successes, not the order, we now have a counting problem: with n total trials (positions), how many ways can we arrange the k successes? Another way of thinking is: “How many ways can we choose k out of n locations in a sequence to place the successes?”.
- The answer is $\binom{n}{k}$, so in total, the probability that $X = k$ is:

$$p(k) = P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

- The function p above is the pmf for the binomial distribution.

Flipping Coins

Suppose that a coin is flipped 10 times. What is the probability that the coin lands heads exactly 6 times?

Here, $n = 10$, and success = Heads. Assuming the coin is fair, we have

$$P(\text{Num Heads} = 6) = \binom{10}{6} (0.5)^6 (0.5)^4 \approx 210 \times 0.00098 \approx .205$$

- Suppose a five 6-sided (fair) dice are rolled simultaneously. What is the probability that at least two of the dice show the value 6?
- Let X denote the number of 6s in this experiment, which takes values in the set $\{0, 1, \dots, 5\}$. We want the probability that $X \geq 2$.

- Because the different values of X are mutually exclusive events (i.e., $X = 2$ implies $X \neq 3$), we can calculate this as:

$$P(X \geq 2) = \sum_{i \in \{2,3,4,5\}} p(i) \approx 0.1962449,$$

where $p(i)$ is the pmf of the binomial(5, 1/6) distribution.

- Alternatively, we can use the complement set, which is smaller:

$$P(X \geq 2) = 1 - P(X < 2) = 1 - (p(0) + p(1)) \approx 0.1962449$$

- *Note:* A binomial RV can be expressed as the sum of independent Bernoulli RVs. That is, let X_1, X_2, \dots, X_n be independent Bernoulli RVs, each with $P(X_i = 1) = p$. Then, $Y = X_1 + X_2 + \dots + X_n$ is a Binomial RV, with parameters (n, p) .

2.1.3 Geometric and Negative Binomial Distributions

Geometric Distribution

- We can construct a *geometric* RV in a similar way that we did with the binomial distribution.
- Suppose instead of having a fixed number of trials, we continue having a trial until our first success. That means that if $X = k$, we will have $k - 1$ failures, one success, and then stop.
- Thus, the pmf can easily be constructed to be:

$$p(k) = P(X = k) = (1 - p)^{k-1}p, \quad k = 1, 2, 3, \dots$$

Geometric Series

Recall from calculus the geometric series:

$$\sum_{i=0}^{\infty} r^i = \frac{1}{1 - r}, \quad \text{if } 0 < r < 1.$$

This identity occurs in the pmf of the geometric series. Let $0 < p < 1$, then

$$\sum_{k=1}^{\infty} (1 - p)^{k-1}p = p \sum_{j=0}^{\infty} (1 - p)^j = p \frac{1}{1 - (1 - p)} = 1.$$

Negative Binomial Distribution

- The *negative binomial* (NB) distribution can be thought of as a generalization of the geometric distribution; rather than stopping when we have exactly one success, we now will stop when we have r successes.
- For any particular sequence of trials of length k that satisfy this condition, the probability is $p^r(1 - p)^{k-r}$.
- The last trial must be a success (because we stopped), so we need to choose the location of the remaining $r - 1$ successes.
- Thus, if X has a negative binomial distribution, the pmf is:

$$p(k) = P(X = k) = \binom{k-1}{r-1} p^r (1 - p)^{k-r}.$$

- Another way that can be helpful for thinking about the NB-distribution is considering it as the sum of r independent geometric random variables:
- We want to represent the total number of trials until the r th success, which is the sum of the number of trials until (and including) the first success, plus the number of trials from the first success until (and including) the second success, and continued until we get r successes.

Negative Binomial Lottery

Suppose that there is a type of lottery where each purchased ticket has equal probability of winning ($p = 1/100$), and there are 3 total prizes to be won. What is the probability that exactly k tickets will be sold until all prizes have been won?

$$P(X = k) = \binom{k-1}{3-1} (0.01)^3 (0.99)^{k-3}.$$

2.1.4 Hypergeometric and Poisson Distribution

The Hypergeometric Distribution

- Suppose that there is a total population of size n , and r have some trait of interest (“success”), and $n - r$ do not (“failure”).
- If we sample m items from the population, then the total number of “successes” in our sample of size m follows a hypergeometric distribution:

$$P(X = k) = \frac{\binom{r}{k} \binom{n-r}{m-k}}{\binom{n}{m}}.$$

- Combinatorially, for $X = k$, we must select k successes out of the total possible r successes in the entire population; there are $\binom{r}{k}$ ways to do this.
- Since we selected m objects in our sample, and we want $m - k$ of them to be failures, we must pick $m - k$ failures from the $n - r$ failures in the population; there are $\binom{n-r}{m-k}$ ways to do this.
- Together, the multiplication principle implies there are $\binom{r}{k} \binom{n-r}{m-k}$ ways that a sample of size m contains k successes from described population.
- Finally, there are a total of $\binom{n}{m}$ ways we can pick our sample:

$$P(X = k) = \frac{\binom{r}{k} \binom{n-r}{m-k}}{\binom{n}{m}}.$$

Balls in a basket

Suppose that there are n balls in a basket, and r balls are black, $n - r$ balls are some other color. If we select $1 \leq m < n$ balls randomly (without replacement), let X denote the number of black balls in our sample of size m . Then, for all $0 \leq k \leq r$,

$$P(X = k) = \frac{\binom{r}{k} \binom{n-r}{m-k}}{\binom{n}{m}}.$$

The Poisson Distribution

- The Poisson distribution is used very frequently in both theory and practice, though the derivation is less intuitive than other distributions, so we will first just provide the pmf:

Definition: Poisson Distribution

The pmf of a random variable X that follows a Poisson distribution with parameter $\lambda > 0$ is

$$p(k) = P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, 2, \dots$$

- Recall from calculus that $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$. Thus, like all pmf's, the pmf of a Poisson distributed RV sums to one:

$$\begin{aligned} \sum_{k=0}^{\infty} p(k) &= \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} \\ &= e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \\ &= e^{-\lambda} e^{\lambda} = 1. \end{aligned}$$

- The value of λ controls the *shape* of the distribution:

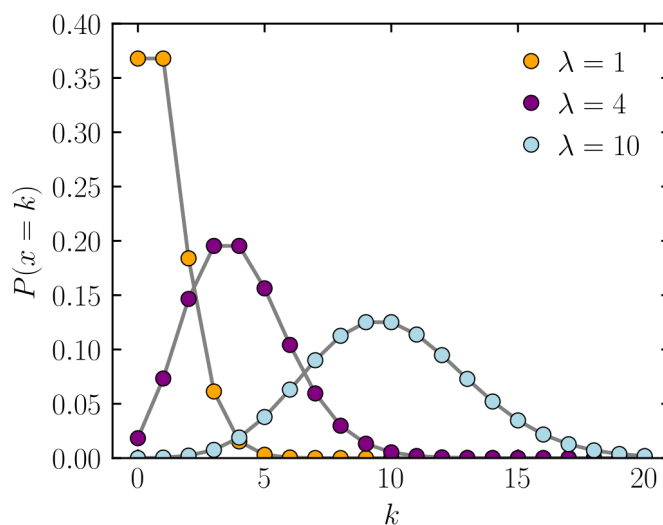


Figure 1: Shape of the Poisson distribution for various values of λ (Wikipedia contributors, 2025b).

- The Poisson distribution can be derived as the limit of a binomial distribution as the number of trials $n \rightarrow \infty$, and $p \rightarrow 0$, such that $np = \lambda$.

Derivation:

Recall the pmf of the binomial distribution can be expressed as

$$p(k) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}.$$

Setting $np = \lambda$, the expression becomes:

$$\begin{aligned} p(k) &= \frac{n!}{k!(n-k)!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{\lambda^k}{k!} \frac{n!}{(n-k)!} \frac{1}{n^k} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k} \end{aligned}$$

Now taking the limit $n \rightarrow \infty$,

$$\begin{aligned} \frac{\lambda}{n} &\rightarrow 0 \\ \frac{n!}{(n-k)!n^k} &\rightarrow 1 \\ \left(1 - \frac{\lambda}{n}\right)^n &\rightarrow e^{-\lambda} \\ \left(1 - \frac{\lambda}{n}\right)^{-k} &\rightarrow 1 \end{aligned}$$

And therefore

$$p(k) \rightarrow \frac{\lambda^k e^{-\lambda}}{k!}.$$

This derivation suggests how a Poisson distribution can arise in practice.

- Let X denote the random variable representing the number of times some event occurs in a fixed time interval.
- Think of dividing the interval into very large number of small sub-intervals of equal length.
- Assume that the sub-intervals are so small that the probability of more than one event in a sub-interval is negligible relative to the probability of one event (which itself is small).
- Finally, assume that the probability of an event in a given sub-interval is identical and independent of that of other sub-intervals.
- Following this, X is nearly binomially distributed, with n being the number of sub-intervals, and $p = \lambda/n$ the probability of the event in each sub-interval.
- Taking the limit, we get something that is nearly Poisson distributed.
- This idea can actually be formalized and made rigorous; you would probably see something like this in a course on stochastic processes.
- The Poisson distribution is often used to model the number of events that occur in a fixed interval.

The Poisson distribution is often good model for the number of events in a fixed time interval if the following conditions are met:

- The occurrence of one event does not affect the occurrence of another.
- The rate at which events occur is fixed.
- Two events cannot occur at the exact same instant.

In this scenario, the random (stochastic) process that generates the data is called a *Poisson process*, which gives rise to the name *rate* for the parameter λ .

Example: Telephone calls

Suppose that an office receives telephone calls as a Poisson process with $\lambda = 0.5$ calls per minute. The number of calls in a 5-min. interval follows a Poisson distribution with parameter $5\lambda = 2.5$. Thus, the probability of no calls in a 5-min. interval is $p(0) = e^{-2.5} \approx .082$; the probability one call is $p(1) = 2.5e^{-2.5} \approx .205$

3 Continuous Random Variables

Introduction

- Because discrete RVs take only a finite number of possibilities, they are relatively simple to define.
- In many situations, however, we are interested in random variables that can take on a continuum of values rather than a finite (or countably infinite) number.

Example: Lifetime of electronic

We might be interested in the lifetime of an electronic component; the total lifetime may be random, but may take on any positive real number.

Density function

- For continuous random variables, we no longer have a pmf (which maps all values of the random variable to their corresponding probabilities).
- Instead, the role of the pmf is taken by a *probability density function* (pdf), which we will denote $f(x)$.

Basic properties of a pdf

If $f(x)$ is a pdf, then $f(x) \geq 0$ for all x , f is piece-wise continuous, and $\int_{-\infty}^{\infty} f(x)dx = 1$.

Probabilities

- If X is a random variable with a density function f , then for any $a \leq b$, the probability that X falls in the interval (a, b) (with the treatment that if $a = b$, the interval collapses to the set $\{a\}$) is given by:

$$P(a < X < b) = \int_a^b f(x)dx.$$

- An immediate consequence of this definition is that $P(X = a) = 0$ for any $a \in \mathbb{R}$.

Example: Continuous uniform random variable

- By *uniform* probability, we mean that all outcomes in the given set are equally as likely.
- For example, if X is a RV with uniform distribution on the interval $[0, 1]$, then any real number in this interval is equally likely, and the probability that X is in a sub-interval of length h should be equal to h .
- You can verify that the following density satisfies this condition:

$$f(x) = \begin{cases} 1, & 0 \leq x \leq 1 \\ 0 & x < 0 \text{ or } x > 1. \end{cases}$$

- For instance, we can pick some $c \in [0, 1]$, and $h \in (0, 1 - c)$. Then the probability that $X \in (c, c + h)$ is given by:

$$P(c < X < c + h) = \int_c^{c+h} 1_{0 \leq x \leq 1} dx = \int_c^{c+h} 1 dx = (c + h) - c = h.$$

- The previous density can be generalized to any interval $[a, b]$, such that $a < b$.

3.1 The Uniform Random Variable

Continuous uniform density

If X is a RV uniformly distributed on an interval $[a, b]$, where $a < b$, then the corresponding density function is:

$$f(x) = \begin{cases} 1/(b-a) & a \leq x \leq b \\ 0 & x < a \text{ or } x > b. \end{cases}$$

- One important thing to note is that if X is a continuous RV, then

$$P(a < X < b) = P(a \leq X < b) = P(a < X \leq b).$$

- This is because the probability X being any particular value is zero; if this were not the case, then the probability of the entire set would infinite (and probabilities must sum to one).

3.2 The cumulative distribution function

Cumulative distribution function

- The cumulative distribution function of a continuous random variable X is defined in the same way as for a discrete random variable:

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(x)dx.$$

- Thus, we can connect the cdf and the pdf of a continuous random variable using the fundamental theorem of calculus.
- Specifically, if $f(x)$ is continuous at x , then $F'(x)$, and

$$P(a \leq X \leq b) = \int_a^b f(x)dx = F(b) - F(a).$$

- This derivation gives some hints at properties of the cdf (both continuous and discrete RVs). Let F be a distribution function. Then the following properties hold:
 - F is right-continuous.
 - F is monotonically increasing (non-decreasing).
 - $F : \mathbb{R} \rightarrow [0, 1]$ and satisfies $\lim_{x \rightarrow -\infty} F(x) = 0$, and $\lim_{x \rightarrow \infty} F(x) = 1$.
- *Note:* every probability distribution supported on the real numbers is uniquely identified by its distribution function F (more to come).

cdf of continuous uniform density

From the definition, we can calculate the cdf of the continuous uniform density rather easily. Suppose that X is uniformly distributed on $[0, 1]$. Then the cdf F is:

$$\begin{aligned} F(X \leq x) &= \int_{-\infty}^x f(x)dx \\ &= \int_{-\infty}^x 1[0 \leq x \leq 1]dx \\ &= \begin{cases} 0 & x < 0 \\ x & 0 \leq x \leq 1 \\ 1 & x > 1. \end{cases} \end{aligned}$$

Percentiles

- You're probably familiar with the term *median*. For any given sample, the median defines the "mid-point", meaning that half of the values are larger, half are smaller.
- This same concept applies to distribution functions.
- That is, the median of a distribution F is defined to be that value $x_{.5}$ such that $P(X < x_{.5}) = 0.5$.
- Formally, the sample median is the same as the definition above, using the *empirical* distribution function (Wikipedia contributors, 2025a).

It is important to note that, as defined, the median value may not be unique!

Definition: Percentile

Let F be the cdf of a continuous random variable. The p th quantile of the distribution F is defined to be any value x_p such that $F(x_p) = P(X \leq x_p) = p$. If F is strictly increasing, then x_p is unique and we say that $F^{-1}(p) = x_p$.

If F is not strictly increasing, then x_p may not be unique; in this case, all such values are considered percentiles. If an inverse function is needed in this case, we will define $F^{-1}(p) = \inf\{x \in \mathbb{R} : F(x) \geq p\}$.

- The last bit of the definition is just some important book keeping to ensure the inverse function exists in odd examples, though I don't think it comes up in this course.

Some important percentiles have their own names, including:

- Median: $p = 1/2$.
- Quartiles (lower and upper): $p = 1/4$, and $p = 3/4$, resp.
- Min: $p = 0$.
- Max: $p = 1$.

Note that the inverse cdf is sometimes called the *quantile function*.

Calculating the inverse cdf

Suppose that

$$F(x) = \begin{cases} 0 & x < 0 \\ x^2 & 0 < x < 1 \\ 1 & x > 1 \end{cases}$$

for $0 \leq x \leq 1$. Find the inverse distribution function F^{-1} .

Solution:

First, let's check that this is a valid distribution function. First, $\lim_{x \rightarrow -\infty} F(x) = 0$, and $\lim_{x \rightarrow \infty} F(x) = 1$.

Now we note that it is trivially monotonically increasing from $(-\infty, 0)$ and $[1, \infty)$. Now on $[0, 1]$, x^2 is also increasing, and hence we have $F(x)$ is a monotonically increasing function.

Finally, we need to check that it is right-continuous. In this case it is trivial, because the only points of potential discontinuity are $x = 0$ and $x = 1$, but the limit clearly exists and equals the function value at both of these points.

Now to find the inverse function, we will focus on the more interesting part of the function, and solve $y = F(x) = x^2$ for x , obtaining $x = F^{-1}(y) = \sqrt{y}$. This provides the inverse function for all points in $(0, 1)$, but what about the endpoints? These are not unique, so we take:

$$F^{-1}(0) = \inf\{x \in \mathbb{R} : F(x) \geq 0\} = \inf_x [0, \infty) = 0,$$

and

$$F^{-1}(1) = \inf\{x \in \mathbb{R} : F(x) \geq 1\} = \inf_x [1, \infty) = 1.$$

Fortunately, these points already match what we found in the mid-point with the square-root function: $\sqrt{0} = 0$ and $\sqrt{1} = 1$ (you could have guessed this would happen since the function is always continuous), so we get

$$F^{-1}(p) = \sqrt{p}.$$

Thus, the inverse function defined on the interval $[0, 1]$ is

3.3 Common random variables

3.3.1 The Exponential Distribution

Exponential Distribution

- The exponential density function is:

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

- Like the Poisson distribution, the exponential density function depends on a single parameter λ .
- When this is the case, we refer to it as the *family* of exponential densities that is *indexed* by the parameter λ .
- The cdf is easily found via the fundamental theorem of calculus:

$$F(x) = \int_{-\infty}^x f(u) du = \begin{cases} 1 - e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

- From this, we can easily find quantiles of the distribution, such as the median: by solving $F(x_{(.5)}) = 1/2$

$$1 - e^{-\lambda x_{(.5)}} = \frac{1}{2} \implies x_{(.5)} = \frac{\log 2}{\lambda}.$$

- The exponential distribution is often used to model lifetimes or waiting times (time-to-event).
- In this context, it's conventional to replace the variable x with t .
- The exponential distribution has a unique property known as the *memoryless* property.
- That is, if something follows an exponential distribution and has already lasted a time of s , then the probability that it will last another t units of time does not depend on s :

Memoryless property: Let T be an exponentially distributed RV, and $s, t > 0$. Calculate $P(T > t+s | T > s)$.

$$\begin{aligned} P(T > t+s | T > s) &= \frac{P(T > t+s \text{ and } T > s)}{P(T > s)} \\ &= \frac{P(T > t+s)}{P(T > s)} \\ &= \frac{1 - F(t+s)}{1 - F(s)} \\ &= \frac{e^{-\lambda(t+s)}}{e^{-\lambda s}} \\ &= e^{-\lambda t} = P(T > t). \end{aligned}$$

- It can be shown that any continuous RV with the *memoryless* property must be exponentially distributed.
- Similarly, it can be shown that any discrete RV with the *memoryless* property must be geometrically distributed (maybe a HW question?)
- The exponential distribution is also related to the *Poisson process* that we have discussed.
- Consider a poisson process with rate λ over an interval $\mathcal{T} \subset \mathbb{R}$.
- While the number of events in any interval $T_0 \subset \mathcal{T}$ of length t follows a Poisson distribution, the time-to-next-event T follows an exponential distribution:
- Suppose that an event occurs at time $t_0 \in T_0$, and let T denote the time until next event. Then:

$$\begin{aligned} P(T > t) &= P(\text{no events in } (t_0, t_0 + t)) \\ &= P(X = 0), \quad \text{where } X \sim \text{Pois}(\lambda t). \\ &= e^{-\lambda t}, \end{aligned}$$

- Therefore T is exponentially distributed with parameter λ .

3.3.2 The Gamma Density

The Gamma Density

- The *gamma* density function depends on two parameters $\alpha > 0$ and $\lambda > 0$.

$$g(t) = \frac{\lambda^\alpha}{\Gamma(\alpha)} t^{\alpha-1} e^{-\lambda t}, \quad t \geq 0.$$

- For $t < 0$, we define $g(t) = 0$.
- The *gamma function*, is defined as:

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du, \quad x > 0.$$

- Note that if $\alpha = 1$, then the gamma density coincides with the exponential density.
- The parameter $\alpha > 0$ in this formulation is called the *shape* parameter.
- The parameter $\lambda > 0$ is called the *scale* parameter.
- As the names suggest, α changes the *shape* of the density function, whereas λ changes the scale of the density (i.e., can be used to change from inches to feet).

3.3.3 The Normal Distribution

The Normal Density

- The normal distribution plays a central role in both probability and statistics.
- It is also called *the Gaussian* distribution, after Carl Friedrich Gauss, who used it as a model for modeling measurement errors.
- One reason it is so important is the *Central Limit Theorem* (CLT, Chapter 6), which suggests that the normal distribution is useful in a large number of settings.

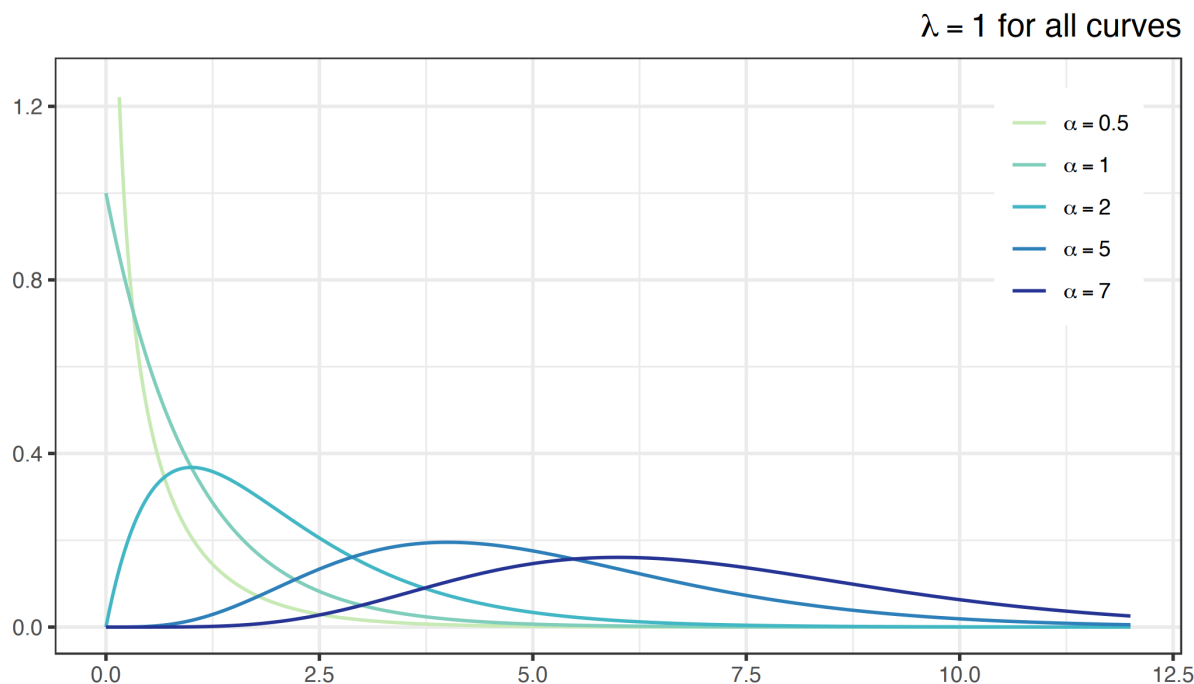


Figure 2: A few Gamma pdf functions for various levels of α .

- Roughly speaking, the CLT states that large(ish) sums (or averages) of independent random variables will be approximately normally distributed.

The density function for the normal distribution depends on two parameters:

- μ : the mean of the distribution.
- σ : the standard deviation of the distribution.

The Normal Density

Let X be a random variable that is normally distributed with mean μ and standard deviation σ . The corresponding probability density function is

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}, \quad -\infty < x < \infty$$

A few notes about this density / distribution:

- You'll often see the shorthand $X \sim N(\mu, \sigma^2)$ to mean “ X follows a normal distribution with parameters μ and σ ”.
- The density is symmetric about μ , meaning $f(\mu - x) = f(\mu + x)$; $x = \mu$ is also the maximum value of the density (mode). The “spread” of the distribution is determined by σ .
- When the mean $\mu = 0$ and standard deviation $\sigma = 1$, we call this the *standard normal distribution*.
- There is no closed form expression for the cdf of the normal distribution. The cdf of the standard normal is usually denoted $\Phi(\cdot)$, and the pdf $\phi(\cdot)$.

3.3.4 The Beta Distribution

The Beta Desnity

- The *beta* distribution is useful for modeling random variables that are restricted to the interval $[0, 1]$:
- The density function depends on two parameters, $\alpha, \beta > 0$.

$$f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad 0 \leq x \leq 1.$$

- Note that if $\alpha = \beta = 1$, the distribution becomes uniformly distributed.
- Both α and β are shape parameters, and the distribution is quite flexible.

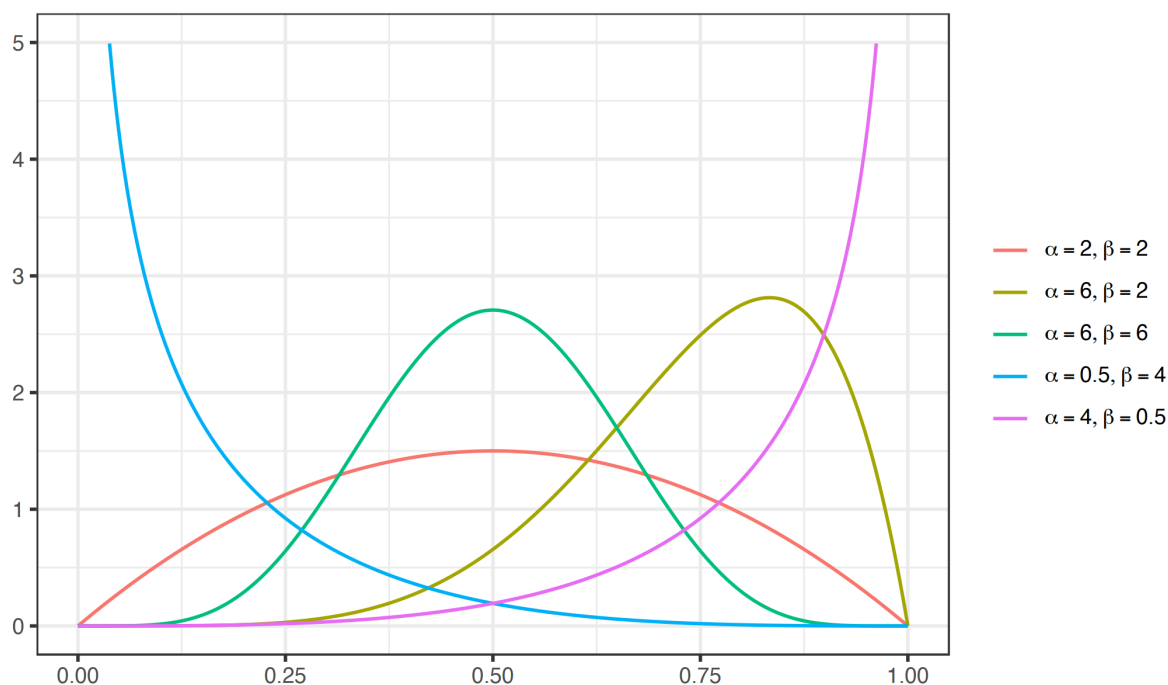


Figure 3: A few Beta pdf functions for various levels of α and β .

Comments on density functions

- So far, we have been using f and F to denote the pdf and cdf of a random variable, respectively. If there is more than one random variable, say X and Y , we may want to distinguish between the functions, and we do so like: $f_X(x)$, or $F_Y(y)$.
- There exist alternative *parameterizations* of the common pmf / pdf functions we have discussed. For example, we introduced the negative binomial distribution with parameters p and r . This is known as the size-probability parameterization, but there also exists an alternative with parameters μ and k , known as the mean-dispersion parameterization (commonly used in Ecology).

4 Functions of a random variable

Variable transformations

- Often we will be interested in function of a random variable.
- Let X be a random variable, and g an arbitrary function.
- Our goal is to find the distribution of the random variable $Y = g(X)$.

Kinetic energy

Let X denote a random variable representing the velocity of a particle of mass m ; we might be interested in the distribution of $Y = \frac{1}{2}mX^2$, the particle's kinetic energy.

- We will eventually provide a rule for a general transformation, $g : \mathbb{R} \rightarrow \mathbb{R}$, but we will first build some intuition using simple transformations.

Example: Linear Transformations (Part I)

Let X be a random variable with pdf and cdf f_X and F_X , respectively. Find the density of $Y = aX + b$, where $a > 0$, and $b \in \mathbb{R}$.

We will use what I like to call “the cdf method”. It’s a simple idea that uses the connection between the cdf and pdf of a random variable, and the definition of the cdf.

$$\begin{aligned}F_Y(y) &= P(Y \leq y) \\&= P(aX + b \leq y) \\&= P\left(X \leq \frac{y - b}{a}\right) \\&= F_X\left(\frac{y - b}{a}\right).\end{aligned}$$

On the left, we have the cdf of Y , and we have equated this to the cdf of X . Now we can differentiate the equation with respect to y :

$$\begin{aligned}f_Y(y) &= \frac{d}{dy}F_X\left(\frac{y - b}{a}\right) \\&= \frac{1}{a}f_X\left(\frac{y - b}{a}\right).\end{aligned}$$

- This same idea is how we will build a general formula for finding the pdf of a transformed random variable.
- Now let’s change it slightly, and see some potential pitfalls.

Example: Linear Transformations (Part II)

Let X be a random variable with pdf and cdf f_X and F_X , respectively. Find the density of $Y = aX + b$, where $a < 0$, and $b \in \mathbb{R}$.

$$\begin{aligned}F_Y(y) &= P(Y \leq y) \\&= P(aX + b \leq y) \\&= P\left(X \geq \frac{y - b}{a}\right) \\&= 1 - F_X\left(\frac{y - b}{a}\right).\end{aligned}$$

On the left, we have the cdf of Y , and we have equated this to the cdf of X . Now we can differentiate the equation with respect to y :

$$\begin{aligned} f_Y(y) &= 0 - \frac{d}{dy} F_X\left(\frac{y-b}{a}\right) \\ &= -\frac{1}{a} f_X\left(\frac{y-b}{a}\right). \end{aligned}$$

- Note that all pdf functions have to be positive. In this example, we have $a < 0$, and therefore the negative sign in the final result helps us ensure that $f_Y(y)$ is positive.
- We can also see that $\left|\frac{1}{a}\right| = -\frac{1}{a}$, so we can write:

$$f_Y(y) = f_X\left(\frac{y-b}{a}\right) \left|\frac{1}{a}\right|.$$

- Now, more generally, assume that g is a strictly monotonically increasing function. Then: $P(g(X) \leq y) = P(X \leq g^{-1}(y))$.
- If g is a strictly monotonically decreasing function, then $P(g(X) \leq y) = P(X \geq g^{-1}(y))$.
- These two examples can be combined to give us a new proposition

Proposition 2.1: monotonic transformations

Let X be a continuous random variable with density f_X , and let $Y = g(X)$, where g is differentiable, and strictly monotonic on an interval I . Then Y has the density function

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|.$$

for y such that $y = g(x)$, and $f_Y(y) = 0$ if $y \neq g(x)$ for any $x \in I$.

proof

Suppose that g is increasing. Then

$$\begin{aligned} F_Y(y) &= P(Y \leq y) \\ &= P(g(X) \leq y) \\ &= P(X \leq g^{-1}(y)) \\ &= F_X(g^{-1}(y)). \end{aligned}$$

Taking the derivative,

$$f_Y(y) = f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|,$$

where the last equality is a result of the fact that if g is increasing, the derivative is always positive, and hence equal to its absolute value.

Now suppose that g is decreasing. Then

$$\begin{aligned} F_Y(y) &= P(Y \leq y) \\ &= P(g(X) \leq y) \\ &= P(X \geq g^{-1}(y)) \\ &= 1 - F_X(g^{-1}(y)). \end{aligned}$$

Taking the derivative,

$$f_Y(y) = -f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|,$$

where the last equality is a result of the fact that if g is decreasing, the derivative is always negative, and hence equal to its absolute value when multiplied by negative one.

- Proposition 2.1 is useful to have / know, but it's often easier to just work from scratch (we'll see a few examples of this).
- Proposition 2.1 can also be extended for use with non-monotonic functions g . The idea is that you break $g : \mathcal{X} \rightarrow \mathcal{Y}$ into a series of functions $g_i : \mathcal{X}_i \rightarrow \mathcal{Y}$ where $\mathcal{X}_i \subset \mathcal{X}$ is a set such that g_i is monotonic. Then, the final density f_Y can be found by summing these functions:

$$f_Y(y) = \sum_i f_X(g_i^{-1}(y)) \left| \frac{d}{dy} g_i^{-1}(y) \right|.$$

See Casella and Berger (2024, Theorem 2.1.8) for more details.

- Here are some common RV transformations that really come in handy.

Example: Normal Distribution I

Let $X \sim N(\mu, \sigma^2)$. Then if $Y = aX + b$, $Y \sim N(a\mu + b, a^2\sigma^2)$.

proof: direct consequence of Linear Transformation examples (parts I and II)

Proposition 2.2: Uniform CDF

Let X be a random variable with cdf F . Then $Z = F(X)$ has a uniform distribution on $[0, 1]$.

proof: $P(Z \leq z) = P(F(X) \leq z) = P(X \leq F^{-1}(z)) = F(F^{-1}(z)) = z$

Proposition 2.3: Inverse Uniform CDF

Let U be uniform on $[0, 1]$, and let $X = F^{-1}(U)$. Then the cdf of X is F .

proof: $P(X \leq x) = P(F^{-1}(U) \leq x) = P(U \leq F(x)) = F(x)$

- This last proposition is very useful. We can use it to *generate random numbers* (pseudorandom).
- Many computer packages have ways of generating numbers uniformly on $U[0, 1]$; the proposition implies that to generate from any arbitrary distribution with cdf F , all we need to do is apply F^{-1} to uniform $[0, 1]$ random numbers.

Chi-square distribution

- If $Z \sim N(0, 1)$, find the density of $X = Z^2$.

$$\begin{aligned} F_X(x) &= P(X \leq x) \\ &= P(-\sqrt{x} \leq Z \leq \sqrt{x}) \\ &= \Phi(\sqrt{x}) - \Phi(-\sqrt{x}). \end{aligned}$$

- We now find the density by differentiating the cdf with respect to x .


$$\begin{aligned} f_X(x) &= \frac{1}{2}x^{-1/2}\phi(\sqrt{x}) + \frac{1}{2}x^{-1/2}\phi(-\sqrt{x}) \\ &= x^{-1/2}\phi(\sqrt{x}) \quad \text{since } \phi \text{ is symmetric} \\ &= \frac{x^{-1/2}}{\sqrt{2\pi}} e^{-x/2}, \quad x \geq 0. \end{aligned}$$

- If you recall that $\Gamma(1/2) = \sqrt{\pi}$, you may recognize that this is a particular instance of the gamma density, with $\alpha = \lambda = 1/2$.
- We call this density the *chi-square density* with 1 degree of freedom.

Final remarks

- We have introduced some concepts of random variables, but a full rigorous discussion about random variables requires background in measure theory. If you are interested in learning more, a good textbook for a statistics student is Resnick (2019); this is considered a graduate level text, and some background in analysis will be beneficial.
- We have discussed only discrete and continuous random variables. In practice, we often run into random variables that have both a discrete and continuous component. For instance, consider a zero-inflated continuous random variable X , where $X = 0$ with probability p (discrete component), but $X \sim N(0, 1)$ with probability $1 - p$ (continuous component).

Acknowledgments

- Compiled on August 15, 2025 using R version 4.5.1.
- Licensed under the [Creative Commons Attribution-NonCommercial license](#).  Please share and remix non-commercially, mentioning its origin.
- We acknowledge [students and instructors for previous versions of this course / slides](#).

References

- Casella G, Berger R (2024). *Statistical inference*. Chapman and Hall/CRC. 25
- Resnick S (2019). *A probability path*. Springer. 27
- Rice JA (2007). *Mathematical statistics and data analysis*, volume 371. 3 edition. Thomson/Brooks/Cole Belmont, CA. 5
- Wikipedia contributors (2025a). “Empirical distribution function — Wikipedia, The Free Encyclopedia.” https://en.wikipedia.org/w/index.php?title=Empirical_distribution_function&oldid=1300940691. [Online; accessed 14-August-2025]. 19
- Wikipedia contributors (2025b). “Poisson distribution.” Accessed 14 August 2025, URL https://en.wikipedia.org/wiki/Poisson_distribution. 1