

# Introduction to Simulation for Biologists: Probability Distributions

Jun Ishigohoka

2024-06-11

**Your chances of getting killed by a  
cow are low, but never zero**



# Contents

Summary	1
Binomial distribution	1
Geometric distribution	3
Poisson distribution	5
Exponential distribution	7
Other distributions used in the course	7
Hypergeometric distribution . . . . .	7
Normal distribution . . . . .	8

## Summary

- Binomial distribution: discrete number of successes
  - $n$ : number of trials (`size` in `rbinom()`)
  - $p$ : prob. of success
  - `rbinom(1, size = n, prob = p)` for discrete random variable  $k$  successes  $k \sim B(n, p)$
- Geometric distribution: discrete number of failures before the first success
  - $p$ : prob. of success
  - `rgeom(1, prob = p)` for discrete random variable  $k$  trials  $k \sim G(p)$
- Poisson distribution: discrete number of events
  - $\lambda$ : rate of events per unit of continuous time/space.
  - `rpois(1, lambda = lambda)` for discrete random variable  $k$  events  $k \sim Pois(\lambda)$
- Exponential distribution: continuous waiting time until the first event
  - $\lambda$ : rate of events per unit of continuous time/space.
  - `rexp(1, rate = lambda)` for continuous random variable  $x$  units  $x \sim Exp(\lambda)$

## Binomial distribution

Throw a die once. The outcome is success if the side is 1, failure otherwise. This type of trial where the outcome is either success or failure with a certain probability is called a **Bernoulli trial**. The distribution of the number of successes  $x$  out of  $n$  independent Bernoulli trials ( $0 \leq x \leq n$ ) is the **binomial distribution**.

In R, we can use the `rbinom` function for binomial sampling. Confusingly, in R, the number of trials is `size` and the number of binomial samplings is `n`.

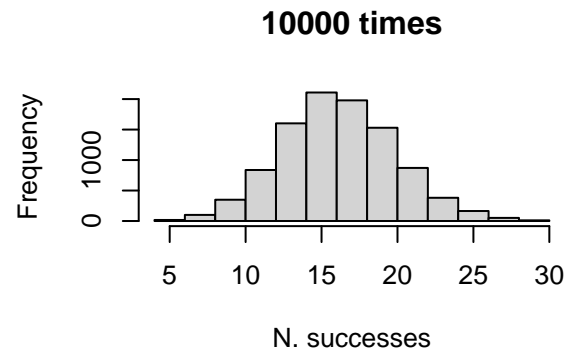
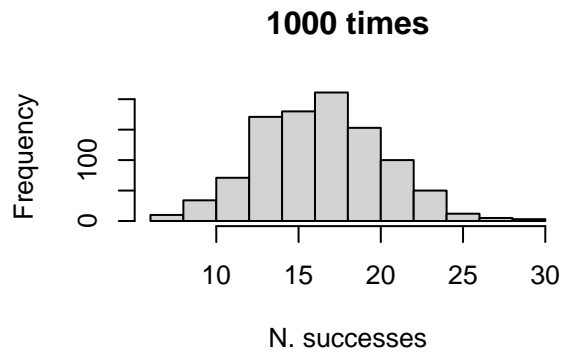
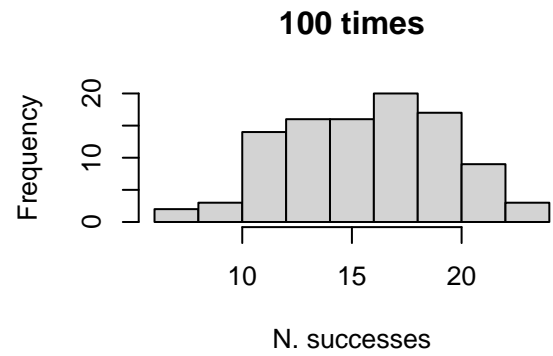
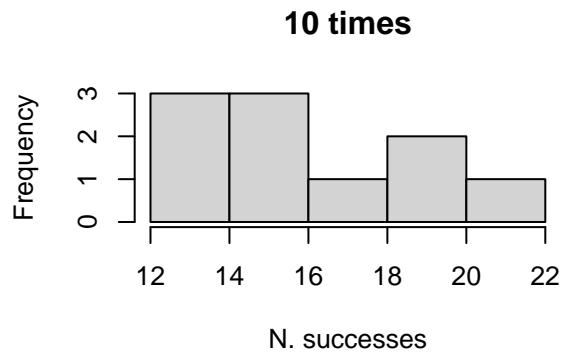
Let's throw a die 100 times and count the number of successes (e.g. the side of 1).

```
1 rbinom(n = 1, size = 100, prob = 1/6) # B(100, 1/6)
```

```
## [1] 18
```

Let's repeat this 10, 100, 1,000, 10,000 times to see the distribution of the number of successes.

```
1 par(mfrow=c(2,2))
2 for(n in c(10, 100, 1000, 10000)){
3     hist(rbinom(n = n, size = 100, prob = 1/6),
4         main = paste0(n, " times"),
5         xlab = "N. successes"
6     )
7 }
```



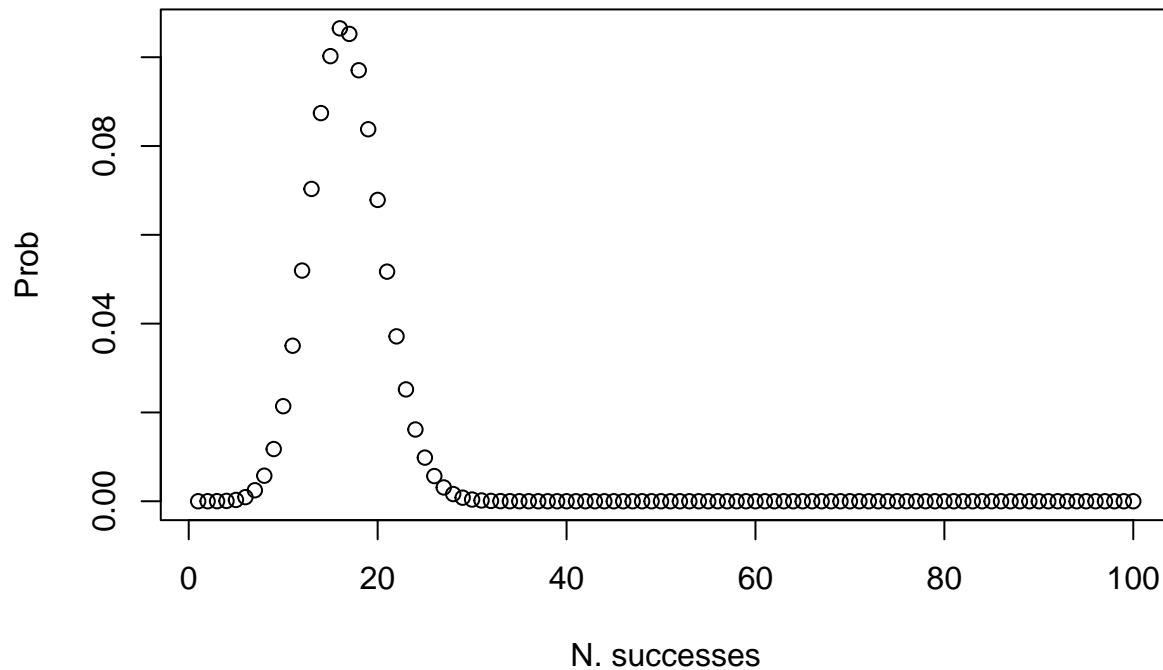
This distribution will approach the probability mass function of the binomial distribution.

```

1 x <- 1:100
2 plot(x, dbinom(x = x, size = 100, prob = 1/6),
3     main = "B(100, 1/6)",
4     ylab = "Prob",
5     xlab = "N. successes"
6 )

```

## B(100, 1/6)



## Geometric distribution

Let's keep throwing a die until the side is 1. How many times does it take?

Naively, we can loop binomial sampling until it succeeds.

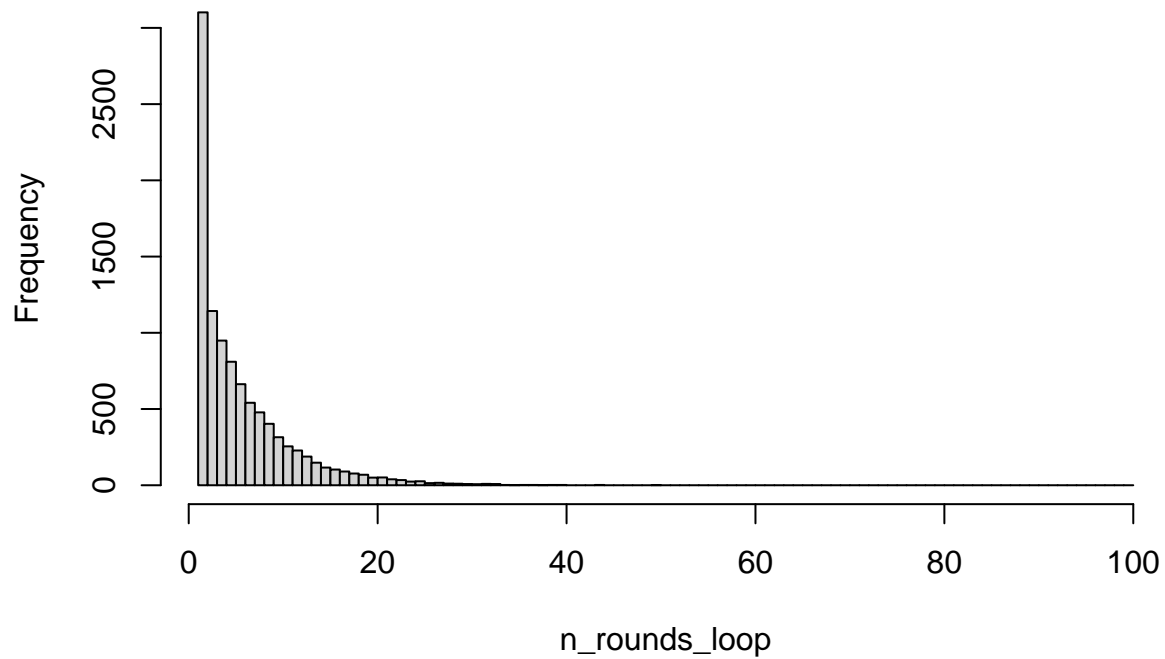
```
1 get_nrounds <- function(prob){
2   c <- 0
3   n <- 0
4   while(n == 0){
5     c <- c + 1
6     n <- rbinom(n = 1, size = 1, prob = prob)
7   }
8   return(c)
9 }
10
11
12 get_nrounds(prob = 1/6)
```

```
## [1] 3
```

Let's visualise how this waiting time is distributed

```
1 n_rounds_loop <- sapply(1:10000,
2   function(x){
3     get_nrounds(prob = 1/6)
4   })
5
6 hist(n_rounds_loop, breaks = seq(1,100))
```

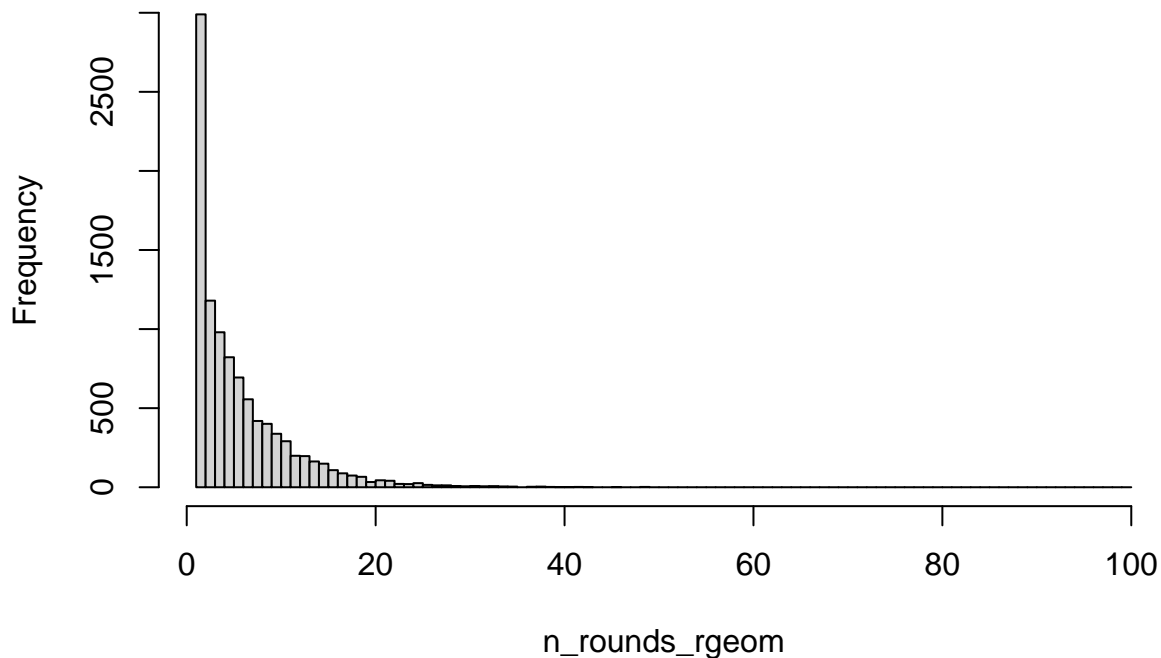
## Histogram of n\_rounds\_loop



The number of failures before the first success is known to follow a **geometric distribution**. So, instead of looping binomial/Bernoulli trials, we can directly obtain the number of throws from a geometric distribution using `rgeom` function.

```
1 n_rounds_rgeom <- rgeom(n = 10000, prob = 1/6) + 1 # + 1 for the first success
2 hist(n_rounds_rgeom, breaks = seq(1,100))
```

## Histogram of n\_rounds\_rgeom



Using a geometric sampling instead of looping is more efficient especially when the probability of success is very low. Let's measure the time it takes to obtain the number of throws until the first success 100 times using the loop.

```
1 system.time(  
2     sapply(1:100,  
3           function(x){  
4               get_nrounds(prob = 1e-6)  
5           })  
6 )  
7
```

```
##    user  system elapsed  
##  3.009   0.007   3.017
```

And if we use geometric sampling

```
1 system.time(  
2     rgeom(n = 10000, prob = 1/6) + 1 # + 1 for the first success  
3 )
```

```
##    user  system elapsed  
##  0.001   0.000   0.000
```

This is a good example that you can still implement your simulation without knowing probability distributions, yet knowing distributions makes your simulation faster and more efficient!

## Poisson distribution

Consider taking 1 mL from a cell suspension at 5,000,000 cells/L (5,000 cell/mL on average). Instead of 1 time of pipetting of 1 mL, let's pipette  $n$  times,  $1/n$  mL each. When  $n$  is very large, the probability that multiple cells are in the taken  $1/n$  mL is so low that it contains either one cell at a probability of  $5000/n$

or no cells at a probability of  $1 - 5000/n$ . By considering of having one cell as success, the total number of cells in the  $n$  portions (thus 1 mL) should follow the binomial distribution  $B(n, p)$ , where  $p = 5000/n$ . When  $n \rightarrow \infty$ , the binomial distribution approaches a **Poisson** distribution. The Poisson distribution has one parameter, the rate  $\lambda$ , which depicts the number of events per unit of time or space (in this case, volume). In this example, the number of cells in a sampled suspension of 1 mL should follow a Poisson distribution with  $\lambda = 5000$ .

Let's take a 1 mL of cell suspension from 5,000 cells/mL.

```
1 rpois(n = 1, lambda = 5000)
```

```
## [1] 5077
```

A nice thing about Poisson distribution is that you can define the unit of time/space as you like as long as the rate parameter is linearly scaled. In the example above, if we take 10 mL instead of 1 mL, the number of sampled cells should now follow another Poisson distribution with  $\lambda = 5000 \times 10 = 50,000$ .

Let's think of the number of mutations (mutation rate of  $1 \times 10^{-9}$  per year per bp) between two homologous sequences of DNA of 1,000,000 bp, where TMRCA is 500,000 years ago. The total branch length between the two sequences is 100,000 years, and mutation occurs at a rate of  $1,000,000[\text{years}] \times 1,000,000[\text{bp}] \times 10^{-9}[\text{year}^{-1}\text{bp}^{-1}] = 1,000$ . So the number of mutations follows a Poisson distribution with  $\lambda = 1,000$ .

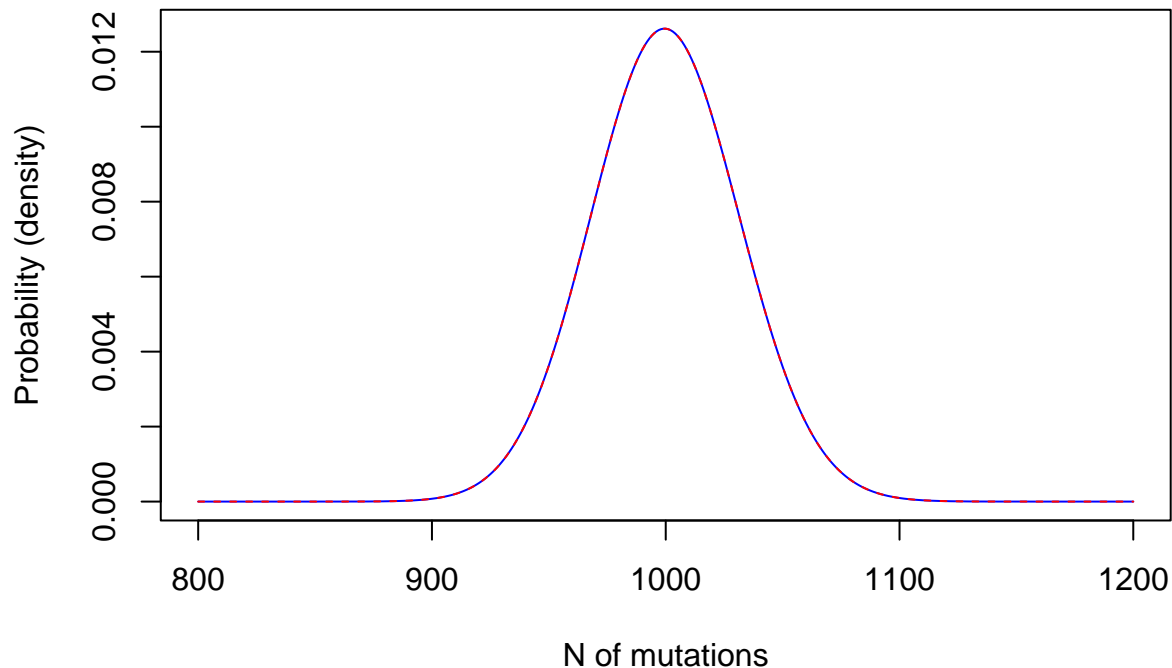
```
1 rpois(n = 1, lambda = 1000)
```

```
## [1] 980
```

We could think of the same problem with binomial sampling, where mutation can occur at a very low probability ( $10^{-9}$  per year per bp) with many trials ( $1,000,000 [\text{bp}] \times 1,000,000 [\text{years}] = 1,000,000,000,000$  trials) We can see how close both ways are:

```
1 plot(800:1200, dpois(800:1200, lambda = 1000),
2     type = 'l',
3     col = "blue",
4     xlab = "N of mutations",
5     ylab = "Probability (density)"
6 )
7 lines(800:1200, dbinom(800:1200, 1e12, 1e-9), type = 'l', col = "red", lty = 2)
```





## Exponential distribution

Continuing on the example of taking cells, let's keep sucking the solution until we get the first cell. The random amount of suspension to take until we get the first cell follows a continuous distribution called **exponential distribution**, which is the limit of geometric distribution where  $\lambda = np$  and  $n \rightarrow \infty$ . As is the Poisson distribution, the exponential distribution has one parameter, the rate  $\lambda$ .

Let's sample the amount of suspension at 5,000 cells/mL until you get the first cell. Confusingly, `lambda` in `rgeom()` is the same rate parameter as `rate` in `rexp()`.

```
1 rexp(1, 5000)
```

```
## [1] 0.0003215086
```

According to [Wikipedia](#), the lifetime probability of a fatal lightning strike is 1/60000. So the time it takes for you to die of a lightning hit can be modelled with an exponential distribution. Let's sample one such event.

```
1 print(paste(rexp(1, 1/60000) * 70, "years"))
```

```
## [1] "2802819.4011189 years"
```

## Other distributions used in the course

### Hypergeometric distribution

The hypergeometric distribution describes the probability of  $k$  successes (random draws for which the object drawn has a specified feature) in  $n$  draws, **without replacement**, from a finite population of size  $N$  that contains exactly  $K$  objects with that feature, wherein each draw is either a success or a failure. This is similar to binomial distribution, but it is sampling without replacement.

The hypergeometric distribution is often seen in enrichment analysis. For example, in the context of gene ontology analysis,  $N$  is the number of all genes,  $K$  is the number of all genes of a focal gene ontology term,  $n$  is the number of significant genes,  $k$  is the number of significant genes of a focal gene ontology.

Let's sample a number of genes out of 100 significant genes that are of one GO term, when there are 2000 genes of this GO term in 10000 genes throughout the genome.

```
1 rhyper(1,  
2     m = 2000, # number of genes of a GO type in the genome  
3     n = 10000 - 2000, # number of genes that are not the GO type in the genome  
4     k = 100 # number of significant genes  
5 )
```

```
## [1] 20
```

Let's compute the p-value of observing 50 out of 100 significant genes that are of this GO term.

```
1 phyper(50, # number of significant genes of the GO type  
2     m = 2000, # number of genes of the GO type in the genome  
3     n = 10000 - 2000, # number of genes that are not the GO type in the genome  
4     k = 100, # number of significant genes  
5     lower.tail = F  
6 )
```

```
## [1] 3.839329e-12
```

## Normal distribution

According to the central limit theorem, the average of many observations of some random variable is a random variable, whose distribution converges to a normal distribution as the number of samples increases. Quantities that are expected to be the sum of many independent processes, such as quantitative genetic traits and errors, can be modelled to follow a normal distribution.

Let's sample a normal random variable with mean of 175 and variance of standard deviation of 6.

```
1 rnorm(1,  
2     mean = 175,  
3     sd = 6  
4 )
```

```
## [1] 175.2156
```