

Simulation

One of the great advantages of using a statistical programming language like R is its vast collection of tools for simulating random numbers.

This lesson assumes familiarity with a few common probability distributions, but these topics will only be discussed with respect to random number generation. Even if you have no prior experience with these concepts, you should be able to complete the lesson and understand the main ideas.

The first function we'll use to generate random numbers is `sample()`. Use `?sample` to pull up the documentation.

```
?sample
```

Let's simulate rolling four six-sided dice: `sample(1:6, 4, replace = TRUE)`.

```
sample(1:6, 4, replace = TRUE)
## [1] 3 2 5 6
```

Now repeat the command to see how your result differs. (The probability of rolling the exact same result is $(1/6)^4 = 0.00077$, which is pretty small!)

```
sample(1:6, 4, replace = TRUE)
## [1] 5 4 6 4
```

`sample(1:6, 4, replace = TRUE)` instructs R to randomly select four numbers between 1 and 6, WITH replacement. Sampling with replacement simply means that each number is “replaced” after it is selected, so that the same number can show up more than once. This is what we want here, since what you roll on one die shouldn't affect what you roll on any of the others.

Now sample 10 numbers between 1 and 20, WITHOUT replacement. To sample without replacement, simply leave off the ‘replace’ argument.

```
sample(1:20, 10)
## [1] 16 11 8 1 12 7 3 14 2 9
```

Since the last command sampled without replacement, no number appears more than once in the output.

`LETTERS` is a predefined variable in R containing a vector of all 26 letters of the English alphabet. Take a look at it now.

```
LETTERS
## [1] "A" "B" "C" "D" "E" "F" "G" "H" "I" "J" "K" "L" "M" "N" "O" "P" "Q"
## [18] "R" "S" "T" "U" "V" "W" "X" "Y" "Z"
```

The `sample()` function can also be used to permute, or rearrange, the elements of a vector. For example, try `sample(LETTERS)` to permute all 26 letters of the English alphabet.

```
sample(LETTERS)
## [1] "F" "A" "H" "J" "O" "Q" "V" "N" "U" "L" "S" "E" "G" "Y" "K" "P" "T"
## [18] "W" "M" "C" "B" "R" "Z" "I" "X" "D"
```

This is identical to taking a sample of size 26 from LETTERS, without replacement. When the 'size' argument to sample() is not specified, R takes a sample equal in size to the vector from which you are sampling.

Now, suppose we want to simulate 100 flips of an unfair two-sided coin. This particular coin has a 0.3 probability of landing 'tails' and a 0.7 probability of landing 'heads'.

Let the value 0 represent tails and the value 1 represent heads. Use sample() to draw a sample of size 100 from the vector c(0,1), with replacement. Since the coin is unfair, we must attach specific probabilities to the values 0 (tails) and 1 (heads) with a fourth argument, prob = c(0.3, 0.7). Assign the result to a new variable called flips.

```
flips <- sample(c(0,1), 100, replace = TRUE, prob = c(0.3, 0.7))
```

View the contents of the flips variable.

```
flips
##      [1] 1 0 1 0 1 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 0 1 1 0 1 1 1 1 1 0 1 1 1 1 1 1
##      [36] 0 1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 0 1 1 0 0 1 1 1 1 0 0 0 1 0 1 1 1 1 1
##      [71] 1 1 1 0 1 1 0 0 0 1 1 1 0 1 1 1 0 1 0 1 0 1 0 1 1 0 0 1 1 1
```

Since we set the probability of landing heads on any given flip to be 0.7, we'd expect approximately 70 of our coin flips to have the value 1. Count the actual number of 1s contained in flips using the sum() function.

```
sum(flips)
```

```
## [1] 70
```

A coin flip is a binary outcome (0 or 1) and we are performing 100 independent trials (coin flips), so we can use rbinom() to simulate a binomial random variable. Pull up the documentation for rbinom() using ?rbinom.

```
?rbinom
```

Each probability distribution in R has an r*** function (for "random"), a d*** function (for "density"), a p*** (for "probability"), and q*** (for "quantile"). We are most interested in the r*** functions in this lesson, but I encourage you to explore the others on your own.

A binomial random variable represents the number of 'successes' (heads) in a given number of independent 'trials' (coin flips). Therefore, we can generate a single random variable that represents the number of heads in 100 flips of our unfair coin using rbinom(1, size = 100, prob = 0.7). Note that you only specify the probability of 'success' (heads) and NOT the probability of 'failure' (tails). Try it now.

```
rbinom(1, size = 100, prob = 0.7)
```

```
## [1] 73
```

Equivalently, if we want to see all of the 0s and 1s, we can request 100 observations, each of size 1, with success probability of 0.7. Give it a try, assigning the result to a new variable called flips2.

```
flips2 <- rbinom(100, size = 1, prob = 0.7)
```

View the contents of flips2.

```
flips2
##      [1] 1 0 1 0 0 1 1 0 0 0 0 1 1 1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 0 1 1 0 1 1
##      [36] 1 1 1 0 1 0 1 1 0 1 0 0 1 0 1 1 1 1 1 1 0 1 1 0 1 1 0 1 0 1 1 1 1 0
```

```
## [71] 0 1 1 0 1 0 0 0 1 1 0 0 1 1 1 1 1 1 1 1 0 0 1 0 1 1 0 1
```

Now use `sum()` to count the number of 1s (heads) in `flips2`. It should be close to 70!

```
sum(flips2)
```

```
## [1] 66
```

Similar to `rbinom()`, we can use R to simulate random numbers from many other probability distributions. Pull up the documentation for `rnorm()` now.

```
?rnorm
```

The standard normal distribution has mean 0 and standard deviation 1. As you can see under the 'Usage' section in the documentation, the default values for the 'mean' and 'sd' arguments to `rnorm()` are 0 and 1, respectively. Thus, `rnorm(10)` will generate 10 random numbers from a standard normal distribution. Give it a try.

```
rnorm(10)
```

```
## [1] -0.07952488 -1.81224950 -0.14665427 1.74160612 -1.41498614
```

```
## [6] -0.67057891 1.58949270 2.34070911 -0.03885349 -0.21897641
```

Now do the same, except with a mean of 100 and a standard deviation of 25.

```
rnorm(10, 100, 25)
```

```
## [1] 99.96306 97.58117 105.05747 72.06460 108.29377 95.72120 61.67664
```

```
## [8] 108.57520 94.74809 93.84602
```

Finally, what if we want to simulate 100 *groups* of random numbers, each containing 5 values generated from a Poisson distribution with mean 10? Let's start with one group of 5 numbers, then I'll show you how to repeat the operation 100 times in a convenient and compact way.

Generate 5 random values from a Poisson distribution with mean 10. Check out the documentation for `rpois()` if you need help.

```
rpois(5, 10)
```

```
## [1] 8 5 16 11 13
```

Now use `replicate(100, rpois(5, 10))` to perform this operation 100 times. Store the result in a new variable called `my_pois`.

```
my_pois <- replicate(100, rpois(5, 10))
```

Take a look at the contents of `my_pois`.

```
my_pois
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13]
## [1,]   10   10   10    8   11    6   10    8    9    4   10   13    9
## [2,]    3   12   10    9    7   13    9    7   12   13   10    9   11
## [3,]   12    8   11    6    5   15   13   11   15   10   10   16   10
## [4,]   13   11   11   10    7   10   10    4   11   13   10   14    9
## [5,]   11    9    4   11   12    9   10    8    8   12    7   11   10
```

##	[,14]	[,15]	[,16]	[,17]	[,18]	[,19]	[,20]	[,21]	[,22]	[,23]	[,24]
## [1,]	15	5	7	8	10	17	11	14	5	8	7
## [2,]	11	6	13	7	15	13	13	10	8	8	11
## [3,]	7	12	7	8	11	8	8	11	10	4	8
## [4,]	15	11	8	13	8	10	6	8	12	15	6
## [5,]	4	10	8	13	9	6	11	9	8	11	10
##	[,25]	[,26]	[,27]	[,28]	[,29]	[,30]	[,31]	[,32]	[,33]	[,34]	[,35]
## [1,]	6	10	4	10	10	8	16	12	9	8	14
## [2,]	11	8	10	9	10	9	6	17	5	12	10
## [3,]	19	10	9	9	14	7	9	10	10	6	10
## [4,]	15	17	11	11	10	12	11	17	7	13	16
## [5,]	14	9	16	8	13	10	10	7	9	13	11
##	[,36]	[,37]	[,38]	[,39]	[,40]	[,41]	[,42]	[,43]	[,44]	[,45]	[,46]
## [1,]	7	10	15	11	11	9	14	10	11	5	8
## [2,]	12	10	11	11	7	7	13	13	11	13	9
## [3,]	8	8	11	8	17	9	16	9	7	14	5
## [4,]	12	11	5	5	9	10	5	13	13	15	11
## [5,]	7	10	6	8	17	6	11	9	11	4	16
##	[,47]	[,48]	[,49]	[,50]	[,51]	[,52]	[,53]	[,54]	[,55]	[,56]	[,57]
## [1,]	9	9	10	7	12	11	6	12	11	8	10
## [2,]	8	10	10	11	5	9	12	7	13	13	12
## [3,]	8	13	12	13	12	11	13	8	10	8	10
## [4,]	12	13	7	8	8	11	7	3	14	7	5
## [5,]	7	12	12	9	11	11	15	6	10	8	14
##	[,58]	[,59]	[,60]	[,61]	[,62]	[,63]	[,64]	[,65]	[,66]	[,67]	[,68]
## [1,]	11	12	9	11	15	13	7	14	9	13	9
## [2,]	5	11	13	10	10	13	12	14	11	6	7
## [3,]	11	14	11	6	10	6	17	12	8	12	7
## [4,]	14	7	8	13	13	10	12	6	3	8	11
## [5,]	7	13	12	15	12	10	7	7	14	4	8
##	[,69]	[,70]	[,71]	[,72]	[,73]	[,74]	[,75]	[,76]	[,77]	[,78]	[,79]
## [1,]	10	10	15	6	11	9	11	13	8	14	12
## [2,]	10	12	8	6	4	7	15	7	8	10	9
## [3,]	8	14	8	13	8	12	13	6	8	11	9
## [4,]	8	9	6	11	6	13	8	11	14	11	7
## [5,]	11	5	8	5	9	10	13	9	9	6	13
##	[,80]	[,81]	[,82]	[,83]	[,84]	[,85]	[,86]	[,87]	[,88]	[,89]	[,90]
## [1,]	15	5	11	11	13	11	11	14	12	8	8

```
## [2,] 10 7 10 13 15 8 7 8 9 8 12
## [3,] 12 6 18 10 12 9 7 13 9 10 11
## [4,] 11 13 14 11 8 9 8 17 6 4 11
## [5,] 10 6 6 11 10 12 11 11 12 11 14
##      [,91] [,92] [,93] [,94] [,95] [,96] [,97] [,98] [,99] [,100]
## [1,] 2 12 13 9 15 13 11 8 9 10
## [2,] 9 6 15 9 8 8 3 17 10 10
## [3,] 6 14 10 14 11 7 13 11 8 8
## [4,] 12 8 8 15 9 7 6 8 5 18
## [5,] 13 4 4 10 8 12 10 9 10 8
```

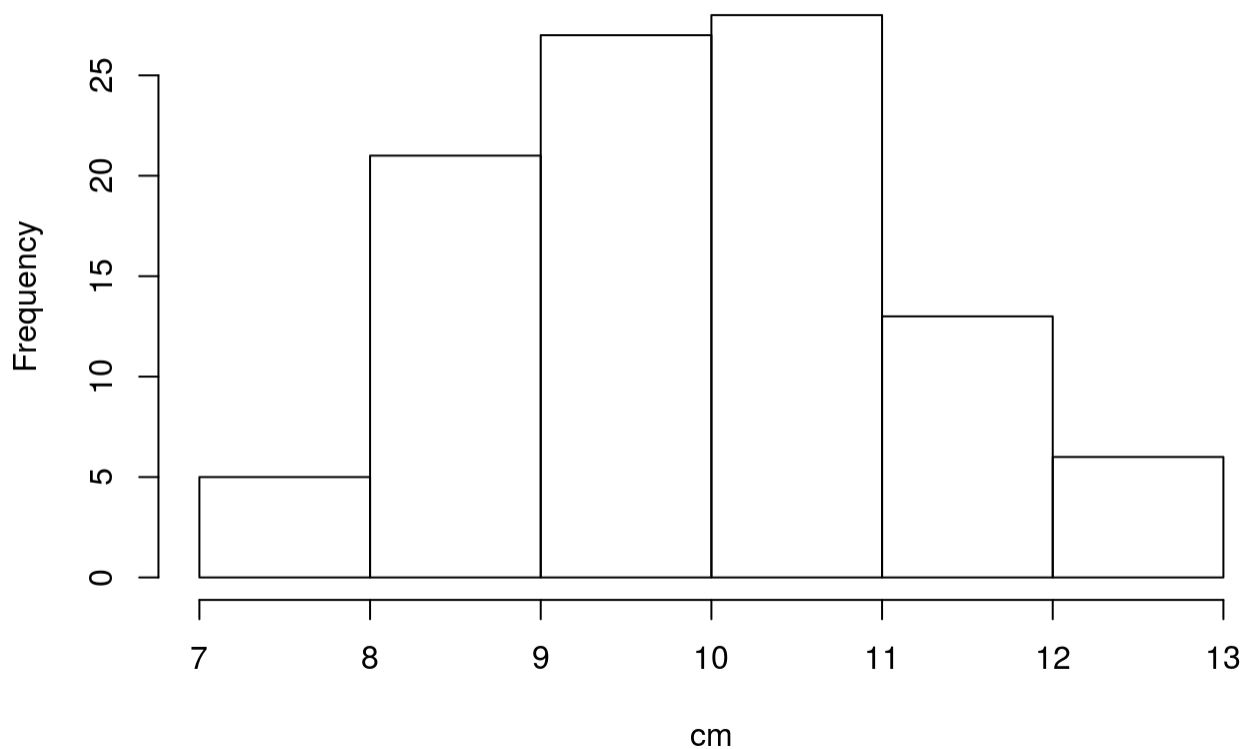
`replicate()` created a matrix, each column of which contains 5 random numbers generated from a Poisson distribution with mean 10. Now we can find the mean of each column in `my_pois` using the `colMeans()` function. Store the result in a variable called `cm`.

```
cm <- colMeans(my_pois)
```

And let's take a look at the distribution of our column means by plotting a histogram with `hist(cm)`.

```
hist(cm)
```

Histogram of cm



Looks like our column means are almost normally distributed, right? That's the Central Limit Theorem at work, but that's a lesson for another day!

All of the standard probability distributions are built into R, including exponential (`rexp()`), chi-squared (`rchisq()`), gamma (`rgamma()`), Well, you see the pattern.

Simulation is practically a field of its own and we've only skimmed the surface of what's possible. I encourage you to explore these and other functions further on your own.