

# Simulating Data

# What are simulation studies?

- Use of artificial data to investigate problems related to:
  - study design
  - analytic approach
  - implications of research findings

# Why learn to simulate data

- Very useful tool in the toolbox
  - study design, what-if analysis, utility analysis
- Gives you a new (and better) understanding of statistics

# Random Number Generators (RNG)

- Functions that generate (pseudo-) random numbers
- Typically based on a probability distribution
- A mathematical function that quantifies the probability of some random variable,  $X$  given a set of parameters

# Random Number Generators (RNG)

- `rnorm` function
- draws  $n$  random samples from a normal distribution
  - with some mean and standard deviation
- has 3 arguments
  - a sample size `n`, a `mean` and `sd` of the distribution

```
1 y <- rnorm(n = 100, mean = 10, sd = 2.5)
2 m <- mean(y)
3 s <- sd(y)
4
5 print(c("sample mean" = m, "sample sd" = s))
```

sample mean	sample sd
9.454534	2.238740

# Random Number Generators (RNG)

- Give us a unique set of numbers each time
- Sometimes we want replicability
  - demonstrations
  - replicating results of a simulation
- “Set the seed”
  - `set.seed()`

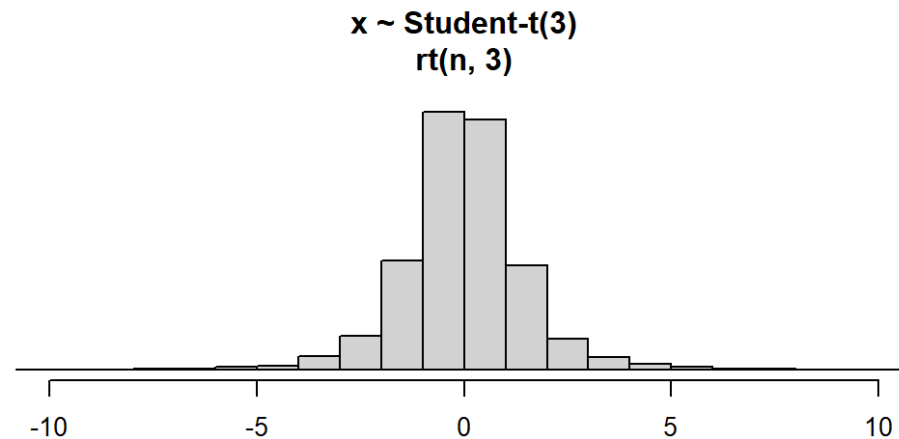
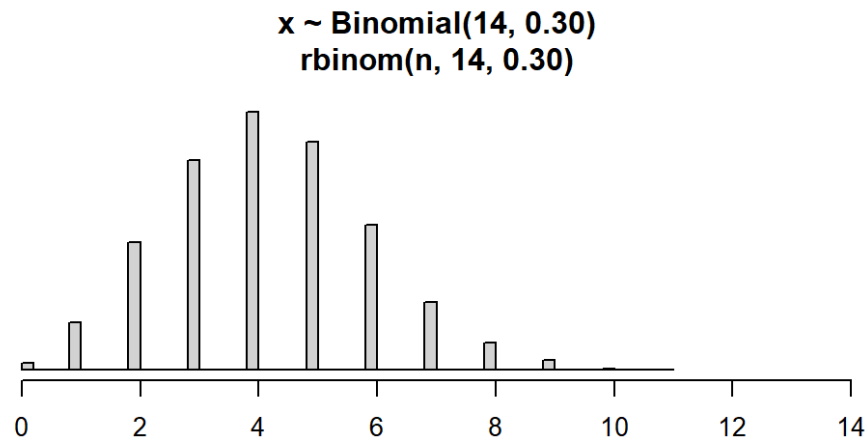
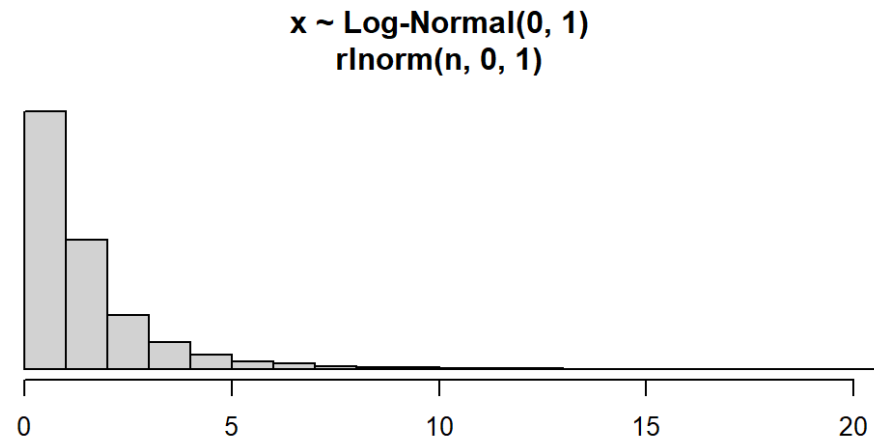
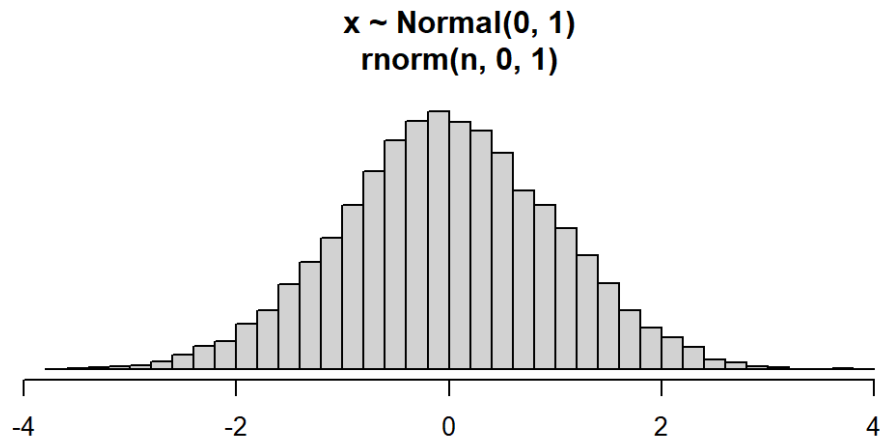
# Random Number Generators (RNG)

```
1  set.seed(1)
2  y1 <- rnorm(n = 100, mean = 0, sd = 1)
3  mean_1 <- mean(y1)
4
5  set.seed(1)
6  y2 <- rnorm(n = 100, mean = 0, sd = 1)
7  mean_2 <- mean(y2)
8
9  print(c("mean 1" = mean_1, "mean 2" = mean_2))
```

```
      mean 1      mean 2
0.1088874 0.1088874
```

# Random Number Generators (RNG)

- Many types of distributions available in R





Distribution	R function	Common Uses
Normal	<code>rnorm</code>	noise/error, performance scores
Log-Normal	<code>rlnorm</code>	response times, anything that is strictly positive
Exponential	<code>rexp</code>	short response times, distances
Gamma	<code>rgamma</code>	response times, anything that is strictly positive

Distribution	R function	Common Uses
Binomial/Bernoulli	<code>rbinom</code>	number of successes out of $k$ attempts
Poisson	<code>rpois</code>	count data, number of events in a given time
Beta	<code>rbeta</code>	proportions or percentages
Multinomial	<code>sample</code>	categorical data, simulating group assignments

# A general framework for simulations in R

- Define your question or goal
  - power analysis, what-if analysis
- Define the generative model and generator
- Define the simulation and simulator
- Define the simulation conditions
  - e.g. sample sizes, test lengths, number of predictors
- Define the summary statistics

# Example 1: Power Analysis

# Define the goal

- What sample size do I need for a 1-sample  $t$ -test?
  - Population Mean:  $\mu = 0$
  - Power:  $1 - \beta = 0.80$  (20% false negative rate)
  - Observed Standard Deviation:  $s = 1$
  - Minimum Effect Size:  $d = \frac{\bar{x} - \mu}{s} = \frac{\bar{x} - 0}{1} = 0.20$

# Define the generative model

- How are the data going to be created?
  - What are the variables involved?
  - What distributions should they be sampled from?
  - Are the parameters fixed (i.e. chosen by you beforehand)?
  - Are the parameters randomly generated as well?
- What needs to be returned by the generator?
  - Just the data?
  - the parameters too?

# Define the generative model

$$\mu = 0$$

$$s = 1$$

$$\mu_x = \frac{d - \mu}{s} = 0.20$$

$$x \sim \text{Normal}(\mu_x, s)$$

# Define the generative model

$$\mu = 0$$

$$s = 1$$

$$\mu_x = \frac{d - \mu}{s} = 0.20$$

$$x \sim \text{Normal}(\mu_x, s)$$

```
1  n <- 100
2
3  d <- 0.20
4  mu <- 0
5  s <- 1
6  mu_x <- (d - mu) / s
7
8  x <- rnorm(n, mu_x, s)
9
10 data <- data.frame(x = x)
```



# Define the generative model

```
1 generator <- function(n_persons) {  
2   d <- 0.20  
3   mu <- 0  
4   s <- 1  
5   mu_x <- (d - mu) / s  
6  
7   x <- rnorm(n_persons, mu_x, s)  
8  
9   data <- data.frame(x = x)  
10  
11   return(data)  
12 }
```

# Define simulation conditions

- What conditions should affect our simulation?
  - sample size, test-length, effect sizes, careless responder %, etc.
- Sample sizes: from 5 to 50
- Simulation iterations: 1,000

```
1 simulation_conditions <- data.frame(n_persons = 5:500)
2
3 n_iters <- 1000
```

# Define the simulation

- What sort of analysis are you investigating?
- t-test

```
1 simulator <- function(conditions) {  
2   n_persons <- conditions[["n_persons"]]  
3   data <- generator(n_persons)  
4  
5   fit <- t.test(x = data$x, mu=0)  
6   p_val <- fit$p.value  
7  
8   output <- data.frame(n_persons=n_persons, p_val=p_val)  
9  
10  return(output)  
11 }
```

# Define the summary statistics

- Simulations generate a lot of data
- We are typically interested in facts about the data and parameters rather than the data and parameters themselves

# Common summary statistics

- Mean Square Error (MSE) of Prediction

- lower = better prediction

- $$\text{MSE} = \frac{1}{n} \sum (b - \hat{b})^2$$

- Mean Bias

- how far is the estimate from the true value on average?

- $$\text{Bias} = \frac{1}{n} \sum b - \hat{b}$$

# Common summary statistics

- Coverage
  - do 95% confidence intervals contain the data generating parameter?
  - $\frac{1}{n} \sum_i^n \delta(-1.98 \cdot \text{SE} < b < 1.98 \cdot \text{SE})$
- Power:  $1 - \beta \approx \frac{1}{n} \sum_i^n \delta(p < \alpha)$

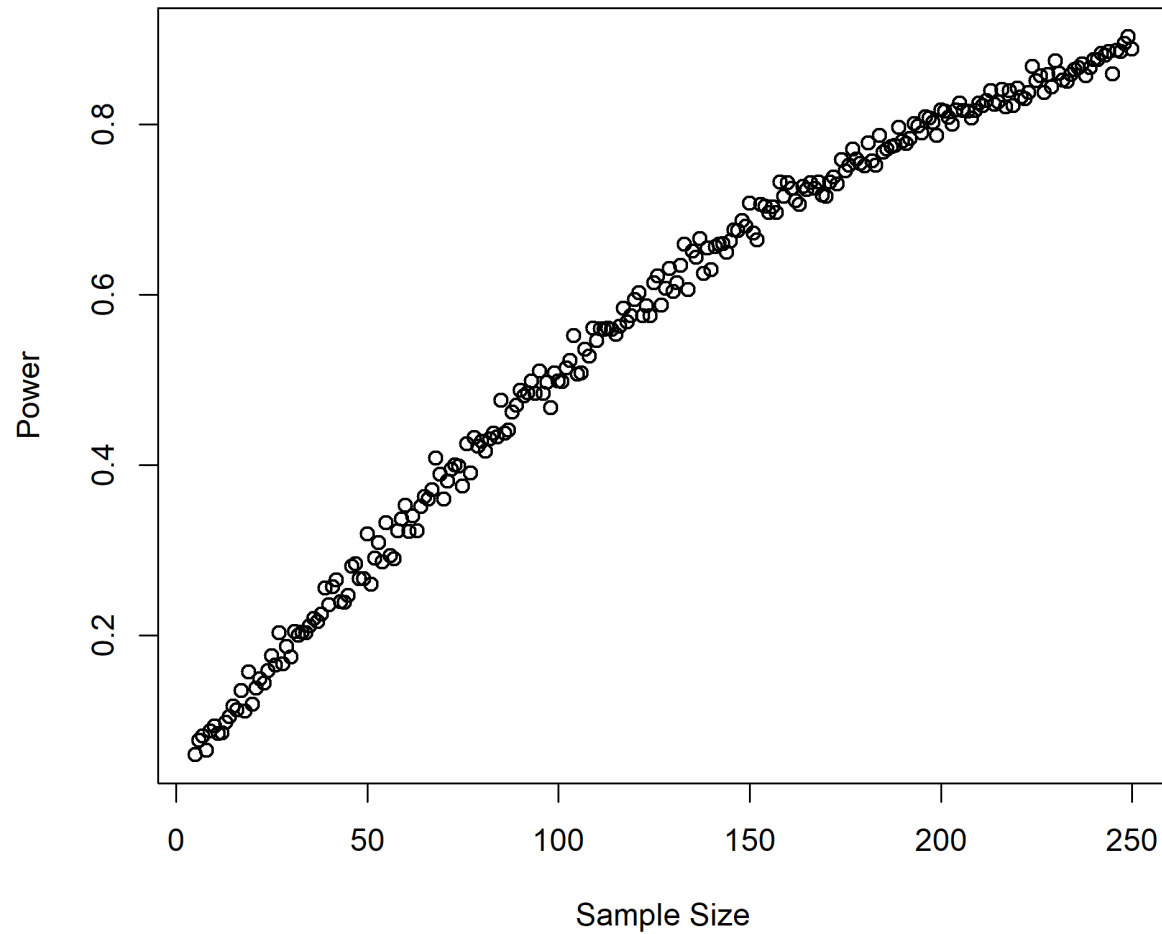
```
1 calculate_power <- function(p_values, alpha=0.05) {  
2   mean(p_values < alpha)  
3 }
```

# Roll it into a simulation

- Create a for-loop

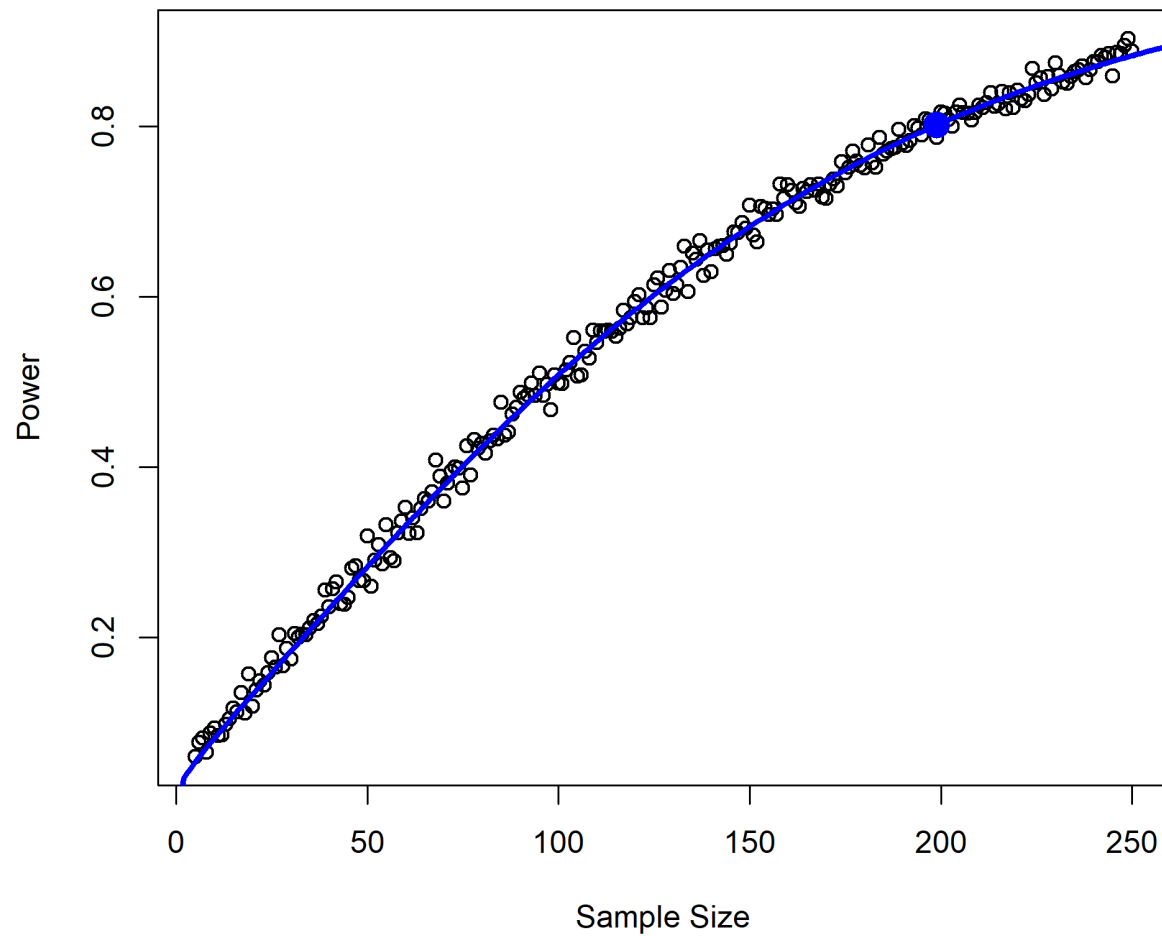
```
1 # create an empty data.frame to store simulated p-values
2 results <- data.frame()
3 n_conditions <- nrow(simulation_conditions)
4
5 # loop over conditions
6 for (c in 1:n_conditions) {
7   conditions <- simulation_conditions[c,,drop=FALSE]
8   # for some number of iterations
9   for (i in 1:n_iters) {
10     new_result <- simulator(conditions)
11
12     results <- rbind(results, new_result)
13
14     # print progress every 100 iterations
15     if (i %% 100 == 0) {
16       cat("Sample Size: ", conditions$n_persons, "; Iteration: ", i
17     }
18   }
19 }
```

# Review Results





# Review Results



# Caveats

- Using RNGs introduces some variability
  - solutions are not exact
  - need lots of iterations
- Can be computationally intensive
  - takes time and computing power
  - running simulations across cores in parallel helps
- Answers are only as thorough as your simulation conditions
  - Do you need to simulate all possible scenarios?

# Resources

- Hallgren, K. A. (2013). Conducting Simulation Studies in the R Programming Environment. *Tutorials in Quantitative Methods for Psychology*, 9(2), 43–60.  
<https://doi.org/10.20982/tqmp.09.2.p043>
- Morris, T. P., White, I. R., & Crowther, M. J. (2019). Using simulation studies to evaluate statistical methods. *Statistics in Medicine*, 38(11), 2074–2102.  
<https://doi.org/10.1002/sim.8086>

<https://distribution-explorer.github.io/index.html>

<https://github.com/tylerjamesryan/Simulation>