

*YOUNG-GEUN statistics*



# Statistical Computing

*R Lab*

**O RLY?**

*Young-geun Kim*



# R Lab for Statistical Computing

*Young-geun Kim*

*Department of Statistics, SKKU*

*dudrms33@g.skku.edu*

*19 May, 2019*



# Contents

<b>Welcome</b>	<b>5</b>
Statistical Computing . . . . .	5
<b>1 Methods for Generating Random Variables</b>	<b>7</b>
1.1 Introduction . . . . .	7
1.2 Pseudo-random Numbers . . . . .	7
1.3 The Inverse Transform Method . . . . .	9
1.4 The Acceptance-Rejection Method . . . . .	13
1.5 Transformation Methods . . . . .	19
1.6 Sums and Mixtures . . . . .	23
1.7 Multivariate Normal Random Vector . . . . .	25
1.8 Stochastic Processes . . . . .	29
<b>2 Monte Carlo Integration and Variance Reduction</b>	<b>37</b>
2.1 Monte Carlo Integration . . . . .	37
2.2 Variance and Efficiency . . . . .	41
2.3 Variance Reduction . . . . .	44
2.4 Importance Sampling . . . . .	54
<b>3 Monte Carlo Methods in Inference</b>	<b>61</b>
3.1 Parametric Bootstrap . . . . .	61
3.2 Monte Carlo Methods for Estimation . . . . .	62
3.3 Confidence interval . . . . .	66
3.4 Hypothesis tests . . . . .	69
3.5 Statistical Methods . . . . .	78
3.6 Bootstrap . . . . .	78



# Welcome

Statistical computing mainly treats useful simulation methods.

```
library(tidyverse)
```

`tidyverse` package family will be used in every chapter. Loading step is in `_common.R`, so it is not included in the text. Sometimes `data.table` library will be called for efficiency.

## Statistical Computing

We first look at *random generation* methods. Lots of simulation methods are built based on this random numbers.

### Sampling from a finite population

Generating random numbers is like sampling. From finite population, we can sample data with or without replacement. For example of sampling with replacement, we toss coins 10 times.

```
sample(0:1, size = 10, replace = TRUE)
#> [1] 1 0 0 1 0 1 1 0 1 1
```

Sampling without replacement: Choose some lottery numbers which consist of 1 to 100.

```
sample(1:100, size = 6, replace = FALSE)
#> [1] 61 83 50 74 34 35
```

### Random generators of common probability distributions

R provides some functions which generate random numbers following famous distributions. Although we will learn some skills generating these numbers in basis levels, these functions do the same thing more elegantly.

```
gg_curve(dbeta, from = 0, to = 1, args = list(shape1 = 3, shape2 = 2)) +
  geom_histogram(
    data = tibble(
      rand = rbeta(1000, 3, 2),
      idx = seq(0, 1, length.out = 1000)
    ),
    aes(x = rand, y = ..density..),
    position = "identity",
    bins = 30,
    alpha = .45,
    fill = gg_hcl(1)
  )
```

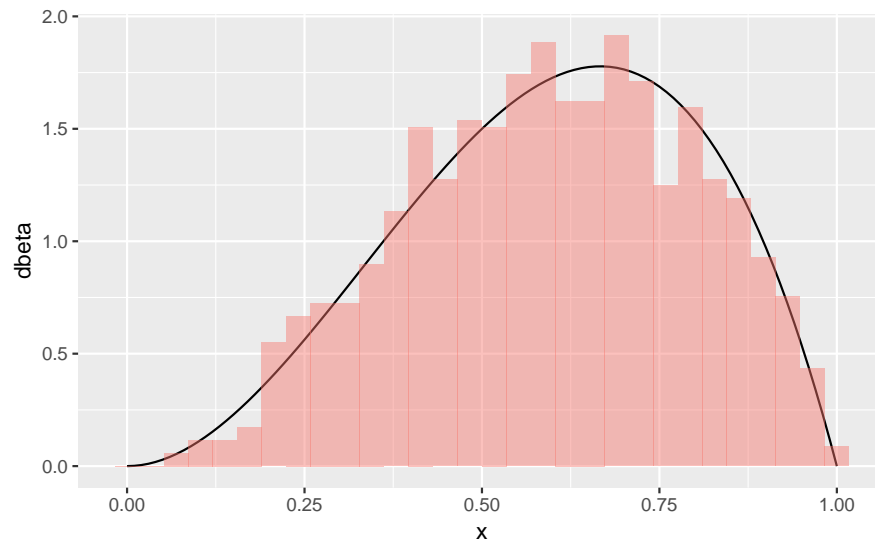


Figure 1: Beta(3,2) random numbers

Figure 1 shows that `rbeta()` function generate random numbers very well. Histogram is of the random number, and the curve is the true beta distribution.



# Chapter 1

## Methods for Generating Random Variables

### 1.1 Introduction

Most of the methods so-called *computational statistics* requires generation of random variables from specified probability distribution. In hand, we can spin wheels, roll a dice, or shuffle cards. The results are chosen randomly. However, we want the same things with computer. Here, **r**. As we know, computer cannot generate complete uniform random numbers. Instead, we generate **pseudo-random** numbers.

### 1.2 Pseudo-random Numbers

**Definition 1.1** (Pseudo-random numbers). Sequence of values generated deterministically which have all the appearances of being independent  $unif(0, 1)$  random variables, i.e.

$$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} unif(0, 1)$$

- behave *as if* following  $unif(0, 1)$
- typically generated from an *initial seed*

#### 1.2.1 Linear congruential generator

Then  $u_1, u_2, \dots, u_n \sim unif(0, 1)$

**Algorithm 1:** Linear congruential generator

```
input :  $a, c \in \mathbb{Z}_+$  and modulus  $m$ 
1 Initialize  $x_0$ ;
2 for  $i \leftarrow 1$  to  $n$  do
3    $x_i = (ax_{i-1} + c) \bmod m$ ;
4 end
5  $u_i = \frac{x_i}{m} \in (0, 1)$ ;
output:  $u_1, u_2, \dots, u_n \sim unif(0, 1)$ 
```

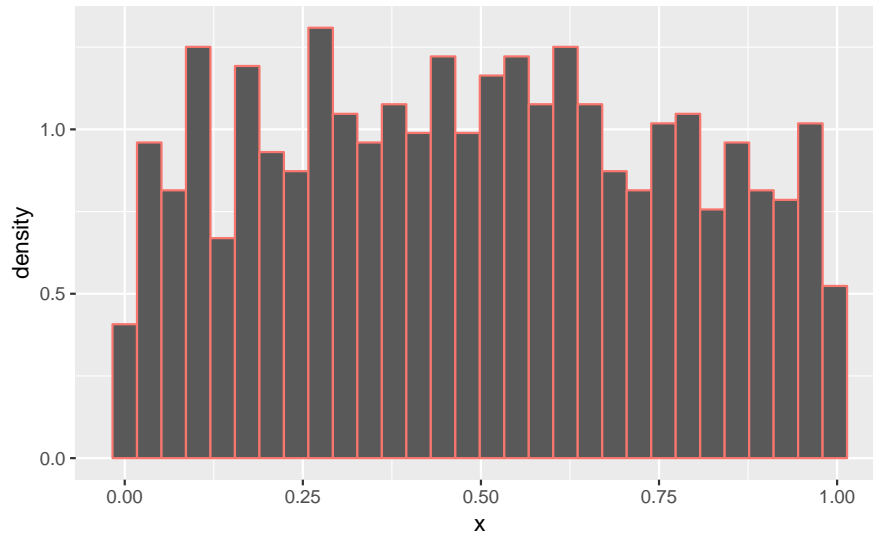
```
lcg <- function(n, seed, a, b, m) {
  x <- rep(seed, n + 1)
  for (i in 1:n) {
    x[i + 1] <- (a * x[i] + b) %% m
  }
}
```

```

  x[-1] / m
}

tibble(
  x = lcg(1000, 0, 1664525, 1013904223, 2^32)
) %>%
  ggplot(aes(x = x)) +
  geom_histogram(aes(y = ..density..), bins = 30, col = gg_hcl(1))

```



### 1.2.2 Multiplicative congruential generator

As we can expect from its name, this is congruential generator with  $c = 0$ .

#### Algorithm 2: Multiplicative congruential generator

**input** :  $a, \in \mathbb{Z}_+$  and modulus  $m$

- 1 Initialize  $x_0$ ;
- 2 **for**  $i \leftarrow 1$  **to**  $n$  **do**
- 3    $x_i = ax_{i-1} \bmod m$ ;
- 4 **end**
- 5  $u_i = \frac{x_i}{m} \in (0, 1)$ ;

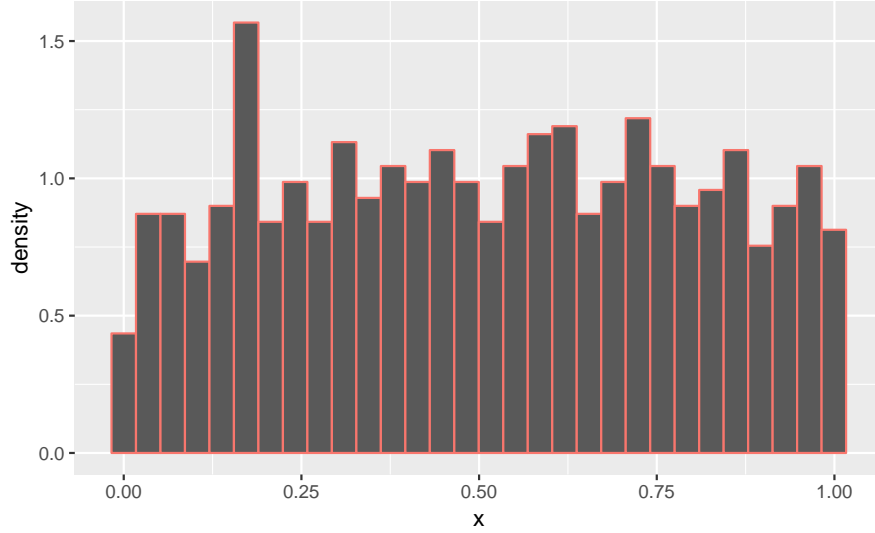
**output**:  $u_1, u_2, \dots, u_n \sim \text{unif}(0, 1)$

We just set  $b = 0$  in our `lcg()` function. The **seed must not be zero**.

```

tibble(
  x = lcg(1000, 5, 1664525, 0, 2^32)
) %>%
  ggplot(aes(x = x)) +
  geom_histogram(aes(y = ..density..), bins = 30, col = gg_hcl(1))

```



### 1.2.3 Cycle

Generate LCG  $n = 32$  with  $a = 1$ ,  $c = 1$ , and  $m = 16$  from the seed  $x_0 = 0$ .

```
lcg(32, 0, 1, 1, 16)
#> [1] 0.0625 0.1250 0.1875 0.2500 0.3125 0.3750 0.4375 0.5000 0.5625 0.6250
#> [11] 0.6875 0.7500 0.8125 0.8750 0.9375 0.0000 0.0625 0.1250 0.1875 0.2500
#> [21] 0.3125 0.3750 0.4375 0.5000 0.5625 0.6250 0.6875 0.7500 0.8125 0.8750
#> [31] 0.9375 0.0000
```

Observe that we have the cycle after  $m$ -th number. Against this problem, we give different seed from every  $(im + 1)$ th random number.

## 1.3 The Inverse Transform Method

**Definition 1.2** (Inverse of CDF). Since some cdf  $F_X$  is not strictly increasing, we define  $F_X^{-1}(y)$  for  $0 < y < 1$  by

$$F_X^{-1}(y) := \inf\{x : F_X(x) \geq y\}$$

Using this definition, we can get the following theorem.

**Theorem 1.1** (Probability Integral Transformation). *If  $X$  is a continuous random variable with cdf  $F(x)$ , then*

$$U \equiv F_X(X) \sim \text{unif}(0, 1)$$

*Probability Integral Transformation.* Let  $U \sim \text{unif}(0, 1)$ . Then

$$\begin{aligned} P(F_X^{-1}(U) \leq x) &= P(\inf\{t : F_X(t) = U\} \leq x) \\ &= P(U \leq F_X(x)) \\ &= F_U(F_X(x)) \\ &= F_X(x) \end{aligned}$$

□

Thus, to generate  $n$  random variables  $\sim F_X$ , we can use *uniform random numbers*.

**Algorithm 3:** Inverse transformation method

```

input : analytical form of  $F_X^{-1}$ 
1 for  $i \leftarrow 1$  to  $n$  do
2    $u_i \stackrel{iid}{\sim} \text{unif}(0, 1);$ 
3    $x_i = F_X^{-1}(u_i);$ 
4 end
output:  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} F_X$ 

```

Note that in R, vectorized operation would be better, i.e. generate `runif(n)` and plug it into given inverse cdf.

### 1.3.1 Continuous case

Denote that the *probability integral transformation* holds for a continuous variable. When generating continuous random variable, applying above algorithm might work.

**Example 1.1** (Exponential distribution). If  $X \sim \text{Exp}(\lambda)$ , then  $F_X(x) = 1 - e^{-\lambda x}$ . We can derive the inverse function of cdf

$$F_X^{-1}(u) = \frac{1}{\lambda} \ln(1 - u)$$

Note that

$$U \sim \text{unif}(0, 1) \Leftrightarrow 1 - U \sim \text{unif}(0, 1)$$

Then we just can use  $U$  instead of  $1 - U$ .

```

inv_exp <- function(n, lambda) {
  -log(runif(n)) / lambda
}

```

If we generate  $x_1, \dots, x_{500} \sim \text{Exp}(\lambda = 1)$ ,

```

gg_curve(dexp, from = 0, to = 10) +
  geom_histogram(
    data = tibble(x = inv_exp(500, lambda = 1)),
    aes(x = x, y = ..density..),
    bins = 30,
    fill = gg_hcl(1),
    alpha = .5
  )

```



Figure 1.1: Inverse Transformation: Exp(1)

### 1.3.2 Discrete case

**Algorithm 4:** Inverse transformation method in discrete case

```

input : analytical form of  $F_X$ 
1 for  $i \leftarrow 1$  to  $n$  do
2    $u_i \stackrel{iid}{\sim} \text{unif}(0, 1)$ ;
3   Take  $x_i$  s.t.  $F_X(x_{i-1}) < U \leq F_X(x_i)$ ;
4 end
output:  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} F_X$ 

```

Table 1.1: Example of a Discrete Random Variable

x	0.0	1.0	2.0	3.0	4.0
p	0.1	0.2	0.2	0.2	0.3

**Example 1.2** (Discrete Random Variable). Consider a discrete random variable  $X$  with a mass function as in Table 1.1.

i.e.



Figure 1.2: Probability Mass Function

Then we have the cdf

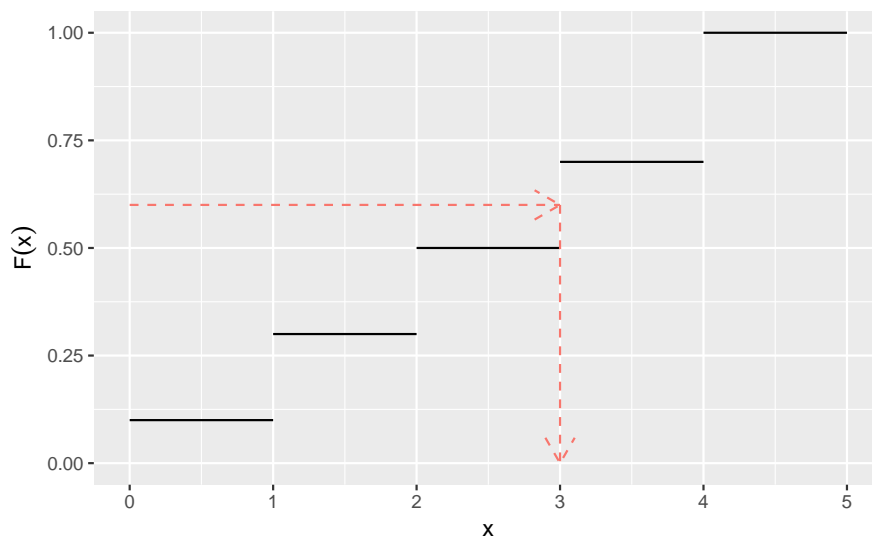


Figure 1.3: CDF of the Discrete Random Variable: Illustration for discrete case

Remembering the algorithm, we can implement `dplyr::case_when()` here.

```
rcustom <- function(n) {
  tibble(u = runif(n)) %>%
    mutate(
      x = case_when(
        u > 0 & u <= .1 ~ 0,
        u > .1 & u <= .3 ~ 1,
        u > .3 & u <= .5 ~ 2,
        u > .5 & u <= .7 ~ 3,
        TRUE ~ 4
      )
    )
}
```

```

) %>%
  select(x) %>%
  pull()
}

tibble(x = rcustom(100)) %>%
  count(x) %>%
  mutate(n = n / sum(n)) %>%
  bind_cols(px = pmf %>% select(p)) %>% # pmf table
  gather(-x, key = "key", value = "value") %>%
  ggplot(aes(x = x, fill = key)) +
  geom_bar(aes(y = value), stat = "identity", position = "dodge", width = .2) +
  scale_fill_discrete(
    name = "Compare",
    labels = c("InvTrans", expression(p(x)))
  ) +
  ylab("prob")

```

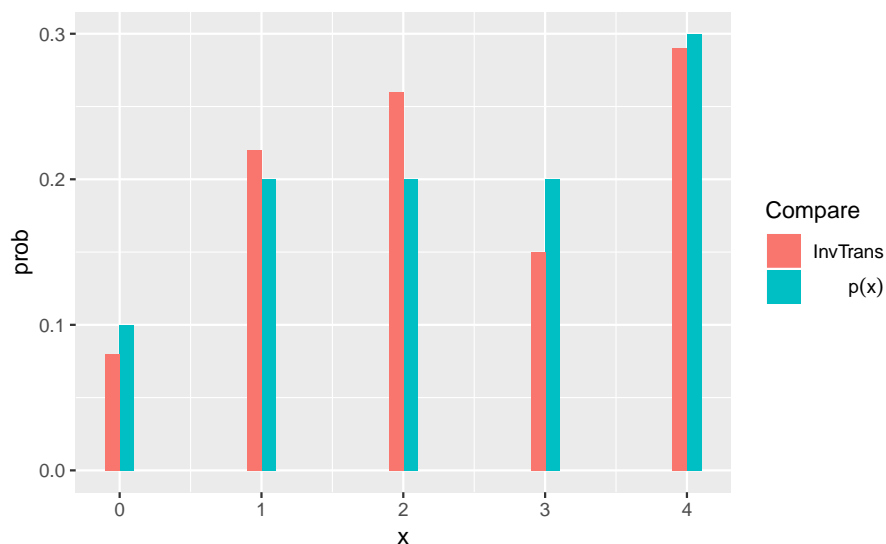


Figure 1.4: Generated discrete random numbers

See Figure 1.4. Comparing random numbers to true pmf, the result can be said okay.

### 1.3.3 Problems with inverse transformation

Examples 1.1 and 1.2. We could generate these random numbers because we aware of

1. analytical  $F_X$
2.  $F^{-1}$

In practice, however, not all distribution have analytical  $F$ . Numerical computing might be possible, but it is not efficient. There are other approaches.

## 1.4 The Acceptance-Rejection Method

Acceptance-rejection method does not require analytical form of cdf. What we need is our *target* density (or mass) function and *proposal* density (or mass) function. Target function is what we want to generate.

Proposal function is of any random variable that is *easy to generate random numbers*. From this approach, we can generate any distribution while computation is not efficient.

pdf or pmf	target or proposal
$f$	target
$g$	proposal - easy to generate random numbers

First of all,  $g$  should satisfy that

$$\text{spt}f \subseteq \text{spt}g$$

Next, for some (pre-specified)  $c > 0$

$$\forall x \in \text{spt}f : \frac{f(x)}{g(x)} \leq c$$

**Algorithm 5:** Acceptance-rejection algorithm

```

input : target  $f$ , proposal  $g$ , and  $c$ 
1 for  $i \leftarrow 1$  to  $n$  do
2    $Y \sim g(y)$ ;
3    $U \sim \text{unif}(0, 1) \perp\!\!\!\perp Y$ ;
4   if  $U \leq \frac{f(Y)}{cg(Y)}$  then
5     Accept  $x_i = Y$ ;
6   else
7     go to Line 2;
8   end
9 end
output:  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} f(x)$ 
```

### 1.4.1 Efficiency

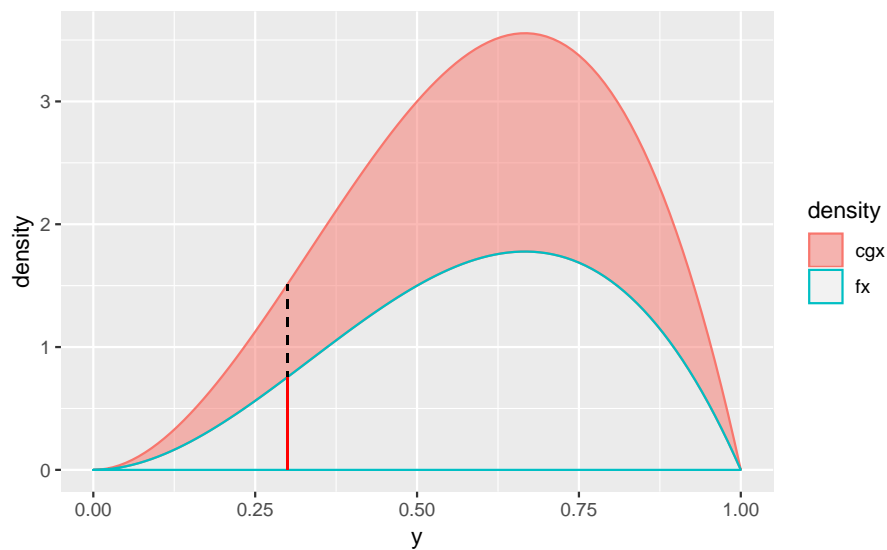


Figure 1.5: Property of AR method



See Figure 1.5. This illustrates the motivation of A-R method. Lower one is  $f(x)$  and the upper one is  $cg(x)$  which covers  $f$ . We can see that

$$0 < \frac{f(x)}{cg(x)} \leq 1$$

The algorithm takes random number from  $Y \sim g$  in each recursive step  $i$ , which is represented as a line in the figure. At this value, the algorithm accept  $Y$  as random number of  $f$  if

$$U \leq \frac{f(Y)}{cg(Y)}$$

Suppose that we choose a point at random on a line drawn in the figure 1.5. If we get the red line, we accept. Otherwise, we reject. In other words, the *colored area is where we reject the given value*. The smaller the area is, the more efficient the algorithm will be.

**Proposition 1.1** (Properties of A-R Method). *See Figure 1.5.*

$$1. \frac{f(Y)}{cg(Y)} \perp\!\!\!\perp U$$

$$2. 0 < \frac{f(x)}{cg(x)} \leq 1$$

3. Let  $N$  be the number of iterations needed to get an acceptance. Then

$$N \sim \text{Geo}(p) \quad \text{where } p \equiv P\left(U \leq \frac{f(Y)}{cg(Y)}\right)$$

and so

$$\begin{cases} P(N = n) = p(1-p)^{n-1} I_{\{1,2,\dots\}}(n) \\ E(N) = \text{average number of iterations} = \frac{1}{p} \end{cases}$$

$$4. X \sim Y \mid U \leq \frac{f(Y)}{cg(Y)}, \text{ i.e.}$$

$$P\left(Y \leq y \mid U \leq \frac{f(Y)}{cg(Y)}\right) = F_X(y)$$

*Remark* (Efficiency). Efficiency of the A-R method depends on  $p = P\left(U \leq \frac{f(Y)}{cg(Y)}\right)$ . In fact,

$$E(N) = \frac{1}{p} = c$$

The algorithm becomes efficient for small  $c$ .

*Proof.* Note that

$$P\left(U \leq \frac{f(y)}{cg(y)}, Y = y\right) = P\left(Y \leq \frac{g(y)}{cg(y)} \mid Y = y\right) P(Y = y)$$

Since  $U \sim \text{unif}(0, 1)$ ,  $P\left(Y \leq \frac{g(y)}{cg(y)} \mid Y = y\right) = \frac{f(y)}{cg(y)}$ .

By construction,  $P(Y = y) = g(y)$ .

It follows that

$$\begin{aligned}
p = P\left(U \leq \frac{f(y)}{cg(y)}\right) &= \int_{-\infty}^{\infty} P\left(U \leq \frac{f(y)}{cg(y)}, Y = y\right) dy \\
&= \int_{-\infty}^{\infty} \frac{f(y)}{cg(y)} g(y) dy \\
&= \frac{1}{c} \int_{-\infty}^{\infty} f(y) dy \\
&= \frac{1}{c}
\end{aligned}$$

Hence,

$$E(N) = \frac{1}{p} = c$$

We can say that the method is efficient when the acceptance rate  $p$  is large, i.e.  $c$  small. □

**Corollary 1.1** (Efficiency of A-R Method). *A-R method is efficient when*

*$g(\cdot)$  is close to  $f(\cdot)$  and*

*have small  $c$ .*

**Corollary 1.2** (Choosing  $c$ ). *To enhance the algorithm, we might choose  $c$  which satisfy*

$$c = \max \left\{ \frac{f(x)}{g(x)} : x \in \text{spt } f \right\}$$

## 1.4.2 Examples

**Example 1.3** (Beta(a,b)). Let  $X \sim \text{Beta}(a, b)$ . Then the pdf of  $X$  is given by

$$f(x) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1} I_{(0,1)}(x)$$

*Solution* (Generating Beta(a,b) with A-R method). Consider proposal density  $g(x) = I_{(0,1)}(x)$ , i.e.  $\text{unif}(0, 1)$ .

To determine the optimal  $c$  s.t.

$$c = \max \left\{ \frac{f(x)}{g(x)} : x \in (0, 1) \right\}$$

find the maximum of

$$\frac{f(x)}{g(x)} = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}$$

Solve

$$\begin{aligned}
\frac{d}{dx} \left( \frac{f(x)}{g(x)} \right) &= \frac{1}{B(a, b)} \left( (a-1)x^{a-2}(1-x)^{b-1} - (b-1)x^{a-1}(1-x)^{b-2} \right) \\
&= \frac{x^{a-2}(1-x)^{b-2}}{B(a, b)} \left( (a-1)(1-x) - (b-1)x \right) \\
&= \frac{x^{a-2}(1-x)^{b-2}}{B(a, b)} (a-1 - (a+b-2)x) = 0
\end{aligned}$$

It follows that

$$\frac{f(x)}{g(x)} \leq \frac{f(\frac{a-1}{a+b-2})}{g(\frac{a-1}{a+b-2})} = c$$

if  $\frac{a-1}{a+b-2} \neq 0, 1$

```
ar_beta <- function(n, a, b) {
  opt_x <- (a - 1) / (a + b - 2)
  opt_c <- dbeta(opt_x, shape1 = a, shape2 = b) / dunif(opt_x)
  X <- NULL
  N <- 0
  while (N <= n) {
    Y <- runif(n)
    U <- runif(n)
    X <- c(X, Y[U <= dbeta(Y, shape1 = a, shape2 = b) / opt_c])
    N <- length(X)
    if ( N > n ) X <- X[1:n]
  }
  X
}
```

Now we try to compare this A-R function to R `rbeta` function.

```
gen_beta <-
  tibble(
    ar_rand = ar_beta(1000, 3, 2),
    sam = rbeta(1000, 3, 2)
  ) %>%
  gather(key = "den", value = "value")

gg_curve(dbeta, from = 0, to = 1, args = list(shape1 = 3, shape2 = 2)) +
  geom_histogram(
    data = gen_beta,
    aes(x = value, y = ..density.., fill = den),
    position = "identity",
    bins = 30,
    alpha = .45
  ) +
  scale_fill_discrete(
    name = "random number",
    labels = c("AR", "rbeta")
  )
```



Figure 1.6: Beta(3,2) Random numbers from each function

In the Figure 1.6, the both histograms are very close to the true density curve. To see more statistically, we can draw a Q-Q plot.

```
gen_beta %>%
  ggplot(aes(sample = value)) +
  stat_qq_line(
    distribution = stats::qbeta,
    dparams = list(shape1 = 3, shape2 = 2),
    col = I("grey70"),
    size = 3.5
  ) +
  stat_qq(
    aes(colour = den),
    distribution = stats::qbeta,
    dparams = list(shape1 = 3, shape2 = 2)
  ) +
  scale_colour_discrete(
    name = "random number",
    labels = c("AR", "rbeta")
  )
)
```



Figure 1.7: Q-Q plot for Beta(3,2) random numbers

See Figure 1.7. We have got series of numbers that are stuck to the beta distribution line.

**Example 1.4** (A-R Method for Discrete case). A-R method can be also implemented to discrete case such as Example 1.2.

Table 1.3: Example of a Discrete Random Variable

x	0.0	1.0	2.0	3.0	4.0
p	0.1	0.2	0.2	0.2	0.3

*Solution* (Generating discrete random numbers using A-R methods). Consider proposal  $g(x) \sim \text{Discrete unif}(0, 1, 2, 3, 4)$ , i.e.

$$g(0) = g(1) = \dots = g(4) = 0.2$$

Then we set

$$c = \max \left\{ \frac{p(x)}{g(x)} : x = 0, \dots, 4 \right\} = \max \{0.5, 1, 1.5\} = 1.5$$

## 1.5 Transformation Methods

### 1.5.1 Continuous

**Proposition 1.2** (Transformation between continuous random variables). *Relation between random variables enables generating target numbers from the others.*

1.  $Z_1, \dots, Z_n \stackrel{iid}{\sim} N(0, 1) \Rightarrow \sum Z_i^2 \sim \chi^2(n)$
2.  $Y_1 \sim \chi^2(m) \perp\!\!\!\perp Y_2 \sim \chi^2(n) \Rightarrow \frac{Y_1/m}{Y_2/n} \sim F(m, n)$
3.  $Z \sim N(0, 1) \perp\!\!\!\perp Y \sim \chi^2(n) \Rightarrow \frac{Z}{\sqrt{Y/n}} \sim t(n)$
4.  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Exp}(\lambda) \Rightarrow \sum Y_i^2 \sim \text{Gamma}(n, \lambda)$

5.  $U \sim \text{unif}(0, 1) \Rightarrow (b - a)U + a \sim \text{unif}(a, b)$
6.  $U \sim \text{Gamma}(r, \lambda) \perp\!\!\!\perp V \sim \text{Gamma}(s, \lambda) \Rightarrow \frac{U}{U+V} \sim \text{Beta}(r, s)$
7.  $Z \sim N(0, 1) \Rightarrow \mu + \sigma Z \sim N(\mu, \sigma^2)$
8.  $Y \sim N(\mu, \sigma^2) \Rightarrow e^Y \sim \text{LogNormal}(\mu, \sigma^2)$

**Example 1.5** (Generating Beta(a, b) using rgamma). From Proposition 1.2, we can generate  $\text{Beta}(a, b)$  random numbers using  $\text{Gamma}(a, 1)$  and  $\text{Gamma}(b, 1)$ .

```
trans_beta <- function(n, shape1, shape2) {
  u <- rgamma(n, shape = shape1, rate = 1)
  v <- rgamma(n, shape = shape2, rate = 1)
  u / (u + v)
}
```

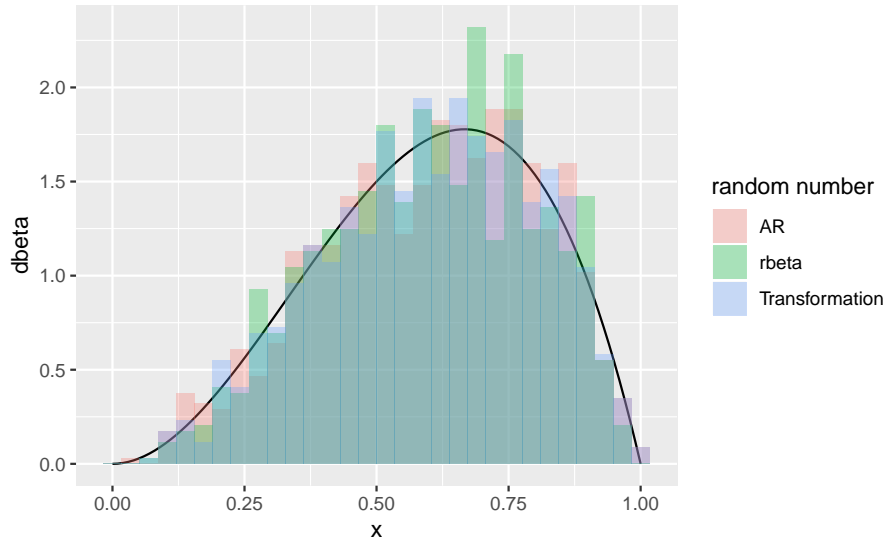


Figure 1.8: Beta(3,2) Random numbers from each function, including transformation method

### 1.5.2 Box-Muller transformation

Denote that Gaussian cdf has no closed form of  $F_X^{-1}$ . Using polar coordinates, we can generate Normal random numbers.

**Theorem 1.2** (Box-Muller transformation). Let  $U_1, U_2 \stackrel{iid}{\sim} \text{unif}(0, 1)$ . Then

$$\begin{cases} Z_1 = \sqrt{-2 \ln U_2} \cos(2\pi U_1) \\ Z_2 = \sqrt{-2 \ln U_2} \sin(2\pi U_1) \end{cases}$$

*Proof.* Write

$$(Z_1, Z_2)^T \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

Then the joint pdf is given by

$$f_{Z_1, Z_2}(x_1, x_2) = \frac{1}{2\pi} \exp\left(-\frac{x_1^2 + x_2^2}{2}\right)$$

Consider polar coordinate transformation  $(R, \theta)$ :  $x_1 = R \cos \theta$  and  $x_2 = R \sin \theta$ .

Since it is also random vector,

$$\begin{aligned} f_{R,\theta}(r, \theta) &= f_{Z_1, Z_2}(x_1, x_2) |J| \\ &= \frac{1}{2\pi} \exp\left(-\frac{x_1^2 + x_2^2}{2}\right) \left| \begin{array}{cc} \frac{\partial x_1}{\partial r} & \frac{\partial x_1}{\partial \theta} \\ \frac{\partial x_2}{\partial r} & \frac{\partial x_2}{\partial \theta} \end{array} \right| \\ &= \frac{1}{2\pi} \exp\left(-\frac{r^2}{2}\right) \left| \begin{array}{cc} \frac{\partial x_1}{\partial r} & \frac{\partial x_1}{\partial \theta} \\ \frac{\partial x_2}{\partial r} & \frac{\partial x_2}{\partial \theta} \end{array} \right| \\ &= \frac{r}{2\pi} \exp\left(-\frac{r^2}{2}\right) \end{aligned}$$

Then each marginal density function can be computed as

$$\begin{aligned} f_\theta(\theta) &= \int_0^\infty \frac{r}{2\pi} \exp\left(-\frac{r^2}{2}\right) dr \\ &= \frac{1}{2\pi} I_{(0, 2\pi)}(\theta) \\ &\stackrel{d}{=} \text{unif}(0, 2\pi) \end{aligned}$$

$$\begin{aligned} f_R(r) &= \int_0^\theta \frac{r}{2\pi} \exp\left(-\frac{r^2}{2}\right) d\theta \\ &= r \exp\left(-\frac{r^2}{2}\right) I_{(0, \infty)}(r) \end{aligned}$$

Thus,

$$f_{R,\theta} = f_\theta f_R \Rightarrow R \perp\!\!\!\perp \theta$$

It follows from inverse transformation theorem that

$$Z_1 = R \cos \theta = \sqrt{-2 \ln U_2} \cos(2\pi U_1)$$

and that

$$Z_2 = R \sin \theta = \sqrt{-2 \ln U_2} \sin(2\pi U_1)$$

where  $U_1, U_2 \stackrel{iid}{\sim} \text{unif}(0, 1)$

□

**Algorithm 6:** Box-Muller transformation

```

1 for  $i \leftarrow 1$  to  $n$  do
2    $U_1, U_2 \stackrel{iid}{\sim} \text{unif}(0, 1);$ 
3    $z_{2i-1} = \sqrt{-2 \ln U_2} \cos(2\pi U_1);$ 
4    $z_{2i} = \sqrt{-2 \ln U_2} \sin(2\pi U_1);$ 
5 end
output:  $z_1, \dots, z_n \stackrel{iid}{\sim} N(0, 1)$ 
```

```

bmnorm <- function(n, mean = 0, sd = 1) {
  n_bm <- ceiling(n / 2)
  tibble(
    theta = runif(n = n_bm, max = 2 * pi),
    R = sqrt(-2 * log(runif(n_bm)))
  ) %>%
  mutate(
    x1 = R * cos(theta),
    x2 = R * sin(theta)
  ) %>%
  gather(x1, x2, key = "key", value = "value") %>%
  mutate(value = mean + sd * value) %>%
  select(value) %>%
  pull()
}

gg_curve(dnorm, from = 0, to = 6, args = list(mean = 3, sd = 1)) +
  geom_histogram(
    data = tibble(x = bmnorm(1000, mean = 3, sd = 1)),
    aes(x = x, y = ..density..),
    bins = 30,
    fill = gg_hcl(1),
    alpha = .5
  )

```

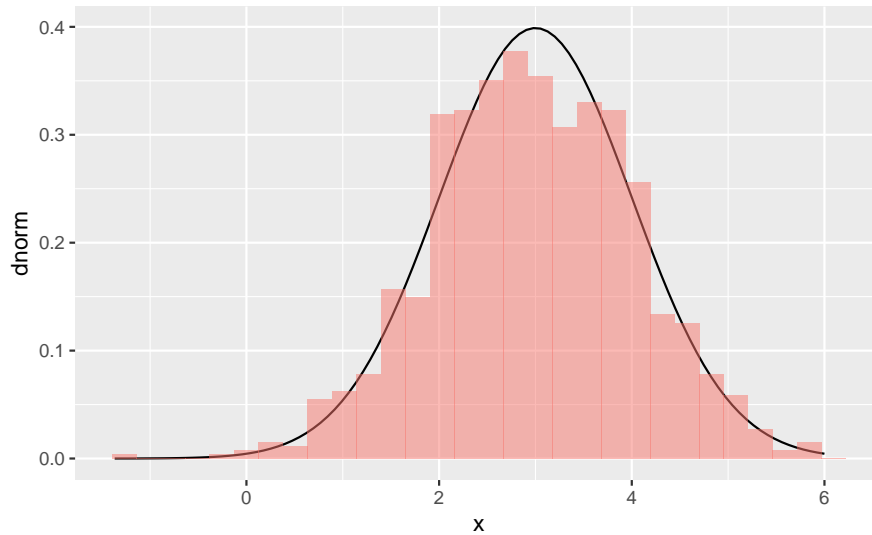


Figure 1.9: Normal random numbers by Box-Muller transformation

### 1.5.3 Discrete

**Proposition 1.3** (Transformation between discrete random variables). *Relation between random variables enables generating target numbers from the others.*

1.  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Bernoulli}(p) \Rightarrow \sum Y_i^2 \sim B(n, p)$
2.  $U \sim \text{unif}(0, 1) \Rightarrow X_i = \lfloor mU \rfloor + 1$
3.  $X = \text{the number of events occurring in 1 unit of time} \sim \text{Poisson}(\lambda)$



**Proposition 1.4** (Bernoulli process). Let  $X_1, X_2, \dots \stackrel{iid}{\sim} \text{Bernoulli}(p)$ .

1.  $N = \text{the number of trials until we see a success, i.e. } X_N = 1 \Rightarrow N \sim \text{Geo}(p)$

2.  $Y_1, \dots, Y_r \stackrel{iid}{\sim} \text{Geo}(p) \Rightarrow \sum_{i=1}^r Y_i = \text{the number of trials until we see } r \text{ successes} \sim \text{NegBin}(r, p)$

**Proposition 1.5** (Count process). Let  $Y_1, Y_2, \dots \stackrel{iid}{\sim} \text{Exp}(\lambda)$  be interarrival times. Then

$$X = \max\{n : \sum Y_i \leq 1\} = \text{the number of events occurring in 1 unit of time} \sim \text{Poisson}(\lambda)$$

## 1.6 Sums and Mixtures

### 1.6.1 Convolutions

**Definition 1.3** (Convolution). Let  $X_1, \dots, X_n$  be independent and identically distributed and let  $S = X_1 + \dots + X_n$ . Then the distribution of  $S$  is called the  $n$ -fold convolution of  $X$  and denoted by  $F_X^{*(n)}$ .

In the last chapter, we have already seen a bunch of random variables that can be generated by summing the other.

**Example 1.6** (Chisquare). Let  $Z_1, \dots, Z_n \stackrel{iid}{\sim} N(0, 1)$ . We know from Proposition 1.2 that

$$V = \sum_{i=1}^n Z_i^2 \sim \chi^2(n)$$

Building a  $n \times \text{df}$  matrix can be a good strategy here. After that, `rowSums` or `colSums` ends the generation work.

```
conv_chisq <- function(n, df) {
  X <-
    matrix(rnorm(n * df), nrow = n, ncol = df)^2
  rowSums(X)
}

gg_curve(dchisq, from = 0, to = 15, args = list(df = 5)) +
  geom_histogram(
    data = tibble(x = conv_chisq(1000, df = 5)),
    aes(x = x, y = ..density..),
    bins = 30,
    fill = gg_hcl(1),
    alpha = .5
  )
```

Figure 1.10:  $\chi^2$  random numbers from Normal sums

### 1.6.2 Mixtures

**Definition 1.4** (Discrete mixture). A random variable  $X$  is a discrete mixture if the distribution of  $X$  is a weighted sum

$$F_X(x) = \sum \theta_i F_{X_i}(x)$$

where constants  $\theta_i$  are called the mixing weights or mixing probabilities.

**Definition 1.5** (Continuous mixture). A random variable  $X$  is a continuous mixture if the distribution of  $X$  is a weighted sum

$$F_X(x) = \int_{-\infty}^{\infty} F_{X|Y=y}(x) f_Y(y) dy$$

**Example 1.7** (Mixture of several Normal distributions). Generate a random sample of size 1000 from a normal location mixture with components of the mixture  $N(0, 1)$  and  $N(3, 1)$ , i.e.

$$F_X = p_1 F_{X_1} + (1 - p_1) F_{X_2}$$

To combine samples easily, we use `foreach` library.

```
library(foreach)
```

As in A-R method, Bernoulli splitting would be used.

$$\begin{cases} F_{X_1} & U > p_1 \\ F_{X_2} & \text{otherwise} \end{cases}$$

```
mix_norm <- function(n, p1, mean1, sd1, mean2, sd2) {
  x1 <- rnorm(n, mean = mean1, sd = sd1)
  x2 <- rnorm(n, mean = mean2, sd = sd2)
  k <- as.integer(runif(n) <= p1)
  k * x1 + (1 - k) * x2
}
```

Try various  $p_1$ , from 0.1 to 1. We would loop and combine by `dplyr::bind_rows()`. Reason for binding is to plot.

```
mixture <-
  foreach(p1 = 0:10 / 10, .combine = bind_rows) %do% {
    tibble(
      value = mix_norm(n = 1000, p1 = p1, mean1 = 0, sd1 = 1, mean2 = 3, sd2 = 1),
      key = rep(p1, 1000)
    )
  }
```

Output is long data format. So we can easily draw a line for each group (`key`).

```
mixture %>%
  ggplot(aes(x = value, colour = factor(key))) +
  stat_density(geom = "line", position = "identity") +
  scale_colour_discrete(
    name = expression(p[1]),
    labels = 0:10 / 10
  ) +
  xlab("x")
```

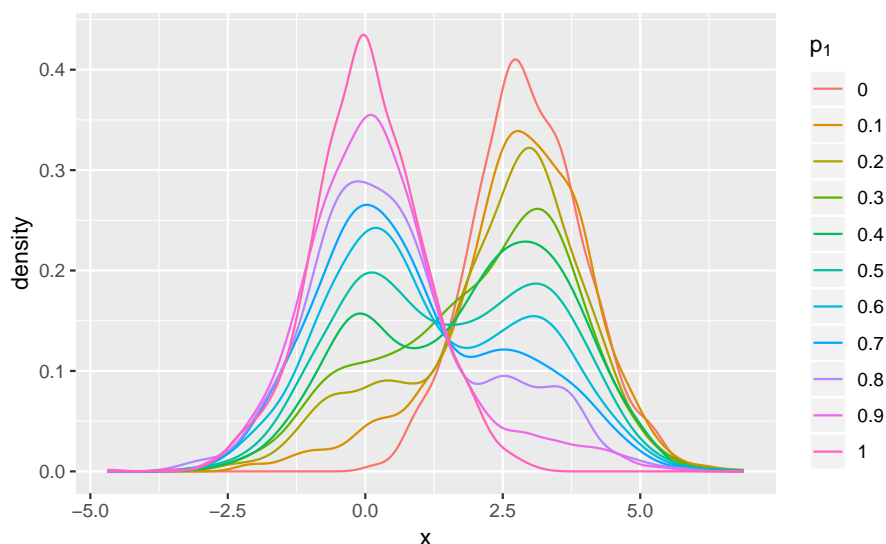


Figure 1.11: Mixture normal random number for each mixing probability

As  $p_1$  becomes larger to 1, the distribution becomes  $N(0, 1)$ . On the contrary,  $p_1 = 0$  results in  $N(3, 1)$ .

## 1.7 Multivariate Normal Random Vector

**Definition 1.6** (Multivariate normal random vector). A random vector  $\mathbf{X} = (X_1, \dots, X_p)^T$  follows multivariate normal distribution if

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{p}{2}|\Sigma|}} \exp \left[ -\frac{1}{2}(\mathbf{x}\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}\boldsymbol{\mu}) \right]$$

*Remark.* Let  $\mathbf{Z} \sim MVN(\mathbf{0}, I)$ . Then

$$\Sigma^{\frac{1}{2}}\mathbf{Z} + \boldsymbol{\mu} \sim MVN(\boldsymbol{\mu}, \Sigma) \quad (1.1)$$

From this remark, we get to generate *standard normal random vector*.

### 1.7.1 Spectral decomposition method

Note that covariance matrix is symmetric.

**Theorem 1.3** (Spectral decomposition). *Suppose that  $\Sigma$  is symmetric. Then*

$$\Sigma = P\Lambda P^T$$

where  $(\mathbf{v}_j, \lambda_j)$  corresponding eigenvector-eigenvalue

$$\begin{cases} P = [\mathbf{v}_1 & \cdots & \mathbf{v}_p] \in \mathbb{R}^{p \times p} \text{ orthogonal} \\ \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p) \end{cases}$$

**Corollary 1.3.** *Suppose that  $\Sigma$  is symmetric. Then*

$$\Sigma^{\frac{1}{2}} = P\Lambda^{\frac{1}{2}}P^T$$

where  $\Lambda^{\frac{1}{2}} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_p})$

`eigen()` performs spectral decomposition. `$values` has eigenvalues and `$vectors` has eigenvectors. We first generate matrix that consists of standard normal random vector:

$$\begin{bmatrix} Z_{11} & Z_{12} & \cdots & Z_{1p} \\ Z_{21} & Z_{22} & \cdots & Z_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ Z_{n1} & Z_{n2} & \cdots & Z_{np} \end{bmatrix}$$

Denote that each observation is row. To use Equation (1.1), we should multiply  $\Sigma^{\frac{1}{2}}$  behind this matrix, not in front of.  $\boldsymbol{\mu}$  matrix should be also made to matrix, in form of

$$\begin{bmatrix} \mu_{11} & \mu_{12} & \cdots & \mu_{1p} \\ \mu_{21} & \mu_{22} & \cdots & \mu_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \mu_{n1} & \mu_{n2} & \cdots & \mu_{np} \end{bmatrix} \in \mathbb{R}^{n \times p}$$

```
rmvn_eigen <- function(n, mu, sig) {
  d <- length(mu)
  ev <- eigen(sig, symmetric = TRUE)
  lambda <- ev$values
  P <- ev$vectors
  sig2 <- P %*% diag(sqrt(lambda)) %*% t(P)
  Z <- matrix(rnorm(n * d), nrow = n, ncol = d)
  X <- Z %*% sig2 + matrix(mu, nrow = n, ncol = d, byrow = TRUE)
  colnames(X) <- paste0("x", 1:d)
  X %>% tbl_df()
}
```

```
# mean vector -----
mu <- c(0, 1, 2)
# symmetric matrix -----
sig <- matrix(numeric(9), nrow = 3, ncol = 3)
```

```
diag(sig) <- rep(1, 3)
sig[lower.tri(sig)] <- c(-.5, .5, -.5) * 2
sig <- (sig + t(sig)) / 2
```

Generate

$$\mathbf{X}_i \sim MVN\left((0, 1, 2), \begin{bmatrix} 1 & -0.5 & 0.5 \\ -0.5 & 1 & -0.5 \\ 0.5 & -0.5 & 1 \end{bmatrix}\right)$$

```
(mvn3 <- rmvn_eigen(1000, mu = mu, sig = sig))
#> # A tibble: 1,000 x 3
#>       x1      x2      x3
#>   <dbl>   <dbl> <dbl>
#> 1 -0.168  1.41  1.80
#> 2  1.39  -0.00942 2.40
#> 3 -0.710  1.30  1.37
#> 4  0.0314 2.04  1.80
#> 5  0.177  0.568  1.71
#> 6 -0.960  1.23  1.61
#> 7 -1.01  1.28  0.106
#> 8  0.272  0.0842 2.12
#> 9  0.148  1.63  2.53
#> 10 -1.24  1.53  1.28
#> # ... with 990 more rows
```

```
mvn3 %>%
  GGally::ggpairs(
    lower = list(continuous = GGally::wrap(gg_scatter, size = 1))
  )
```

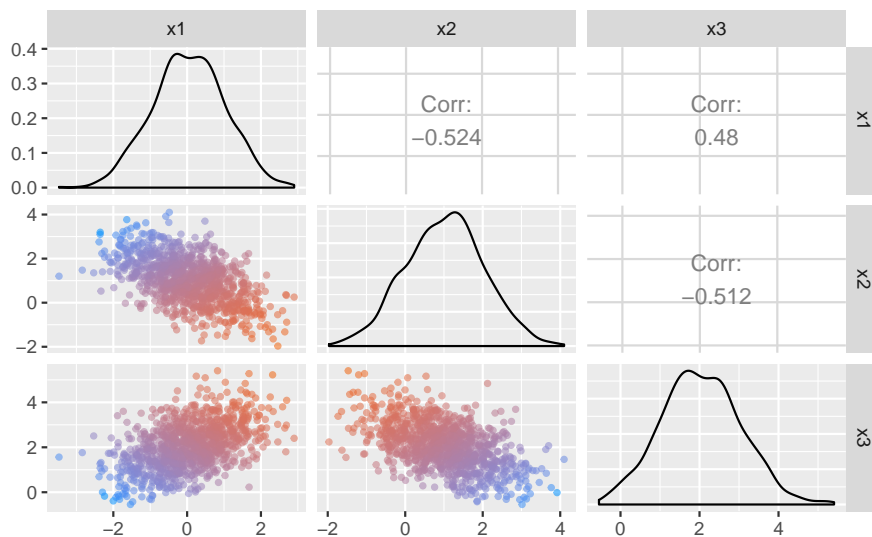


Figure 1.12: Multivariate normal random vector - spectral decomposition method

### 1.7.2 Singular value decomposition

SVD can be said to be a kind of generalization of spectral decomposition. This method can be used for any matrix, i.e. non-symmetric matrix. For  $\Sigma$ , SVD and spectral decomposition is equivalent. However, SVD does not account for symmetric property, so this method is less efficient compared to spectral decomposition.

```
rmvn_svd <- function(n, mu, sig) {
  d <- length(mu)
  S <- svd(sig)
  sig2 <- S$u %*% diag(sqrt(S$d)) %*% t(S$v)
  Z <- matrix(rnorm(n * d), nrow = n, ncol = d)
  X <- Z %*% sig2 + matrix(mu, nrow = n, ncol = d, byrow = TRUE)
  colnames(X) <- paste0("x", 1:d)
  X %>% tbl_df()
}
```

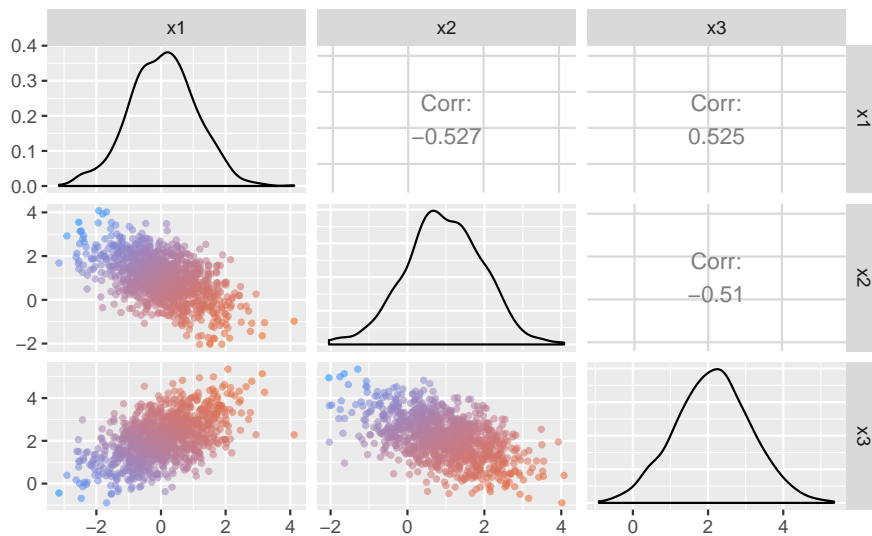


Figure 1.13: Multivariate normal random vector - svd

### 1.7.3 Choleski decomposition

**Theorem 1.4** (Cholesky decomposition). *Suppose that  $\Sigma$  is symmetric and positive definite. Then*

$$\Sigma = Q^T Q$$

where  $Q$  is an upper triangular matrix.

**Corollary 1.4.** *Suppose that  $\Sigma$  is symmetric and positive definite. For cholesky decomposition 1.4, define*

$$\Sigma^{\frac{1}{2}} = Q$$

`chol()` computes cholesky decomposition. In R, it gives upper triangular  $Q$ . Since some statements cholesky decomposition by  $\Sigma = LL^T$  with lower triangular matrix, try not to confuse.

```
rmvn_chol <- function(n, mu, sig) {
  d <- length(mu)
  sig2 <- chol(sig)
  Z <- matrix(rnorm(n * d), nrow = n, ncol = d)
```

```

X <- Z %*% sig2 + matrix(mu, nrow = n, ncol = d, byrow = TRUE)
colnames(X) <- paste0("x", 1:d)
X %>% tbl_df()
}

```

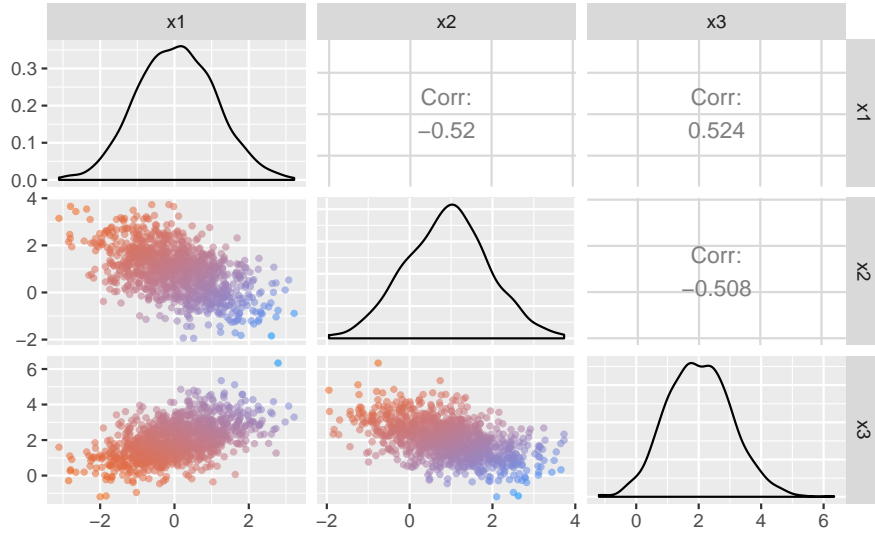


Figure 1.14: Multivariate normal random vector - cholesky decomposition

## 1.8 Stochastic Processes

**Definition 1.7** (Stochastic process). A stochastic process is a collection  $\{X(t) : t \in T\}$  of random variables indexed by the set  $T$ . The index set  $T$  could be discrete or continuous.

A State space is called the set of possible values that  $X(t)$  can take.

**Definition 1.8** (Discrete Time Markov Chain).  $\{X_n : n = 0, 1, 2, \dots\}$  is a Discrete time markov chain on  $S$  if and only if

1.  $S$  is at most countable
2. Markov property  $P(X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) = P(X_{n+1} = j \mid X_n = i) = P_{ij}$

If  $P_{ij}$  is fixed, then  $\{X_n\}$  is called time homogeneous. Otherwise, it is called nonhomogeneous.

**Definition 1.9** (Random walk model). Let  $\{Y_n : n \in \mathbb{N}\}$  be an IID process on  $S$  s.t.

$$P(Y_n = k) = p_k$$

Define

$$S_n := \begin{cases} 0 & n = 0 \\ S_0 + Y_1 + \dots + Y_n & n \in \mathbb{N} \end{cases}$$

### 1.8.1 Gambler's ruin model

**Definition 1.10** (Gambler's ruin model). Let  $\{Y_n : n \in \mathbb{N}\}$  be a process on  $\{-1, 1\}$  s.t.

$$P(Y_n = 1) = p, \quad P(Y_n = -1) = 1 - p$$

Define

$$X_n := \begin{cases} a & n = 0 \\ a + Y_1 + \cdots + Y_n & n \in \mathbb{N} \end{cases}$$

**Example 1.8** (Gambling with coin). Suppose that A and B each start with a stake of \$10, and bet \$1 on consecutive coin flips. The game ends when either one of the players has all the money. Let  $S_n$  be the fortune of player A at time  $n$ . Then  $\{S_n, n \geq 0\}$  is a symmetric random walk with absorbing barriers at 0 and 20. Simulate a realization of the process  $\{S_n, n \geq 0\}$  and plot  $S_n$  vs the time index from time 0 until a barrier is reached.

Here we have

$$P(Y_n = 1) = P(Y_n = -1) = \frac{1}{2}$$

$$S_n := \begin{cases} 10 & n = 0 \\ 10 + Y_1 + \cdots + Y_n & n \in \mathbb{N} \end{cases}$$

```
gambling <- function(begin = 10, betting = 1, prob = .5) {
  N <- begin * 2
  sa <- begin
  record <- tibble(a = begin, b = begin)
  while(all(record > 0)) {
    sa <- ifelse(runif(1) <= prob, sa + betting, sa - betting)
    record <-
      record %>%
      bind_rows(c(a = sa, b = N - sa))
    if (sa == N) break()
  }
  record %>%
    mutate(idx = 1:n()) %>%
    select(idx, a, b)
}
```

```
gambling(begin = 10, betting = 1, prob = .5) %>%
  gather(-idx, key = "player", value = "fortune") %>%
  ggplot(aes(x = idx, y = fortune, colour = player)) +
  geom_path() +
  geom_point(alpha = .5, size = 1) +
  geom_hline(yintercept = c(0, 20), col = I("grey")) +
  labs(
    x = "Betting",
    y = "Fortune"
  )
```



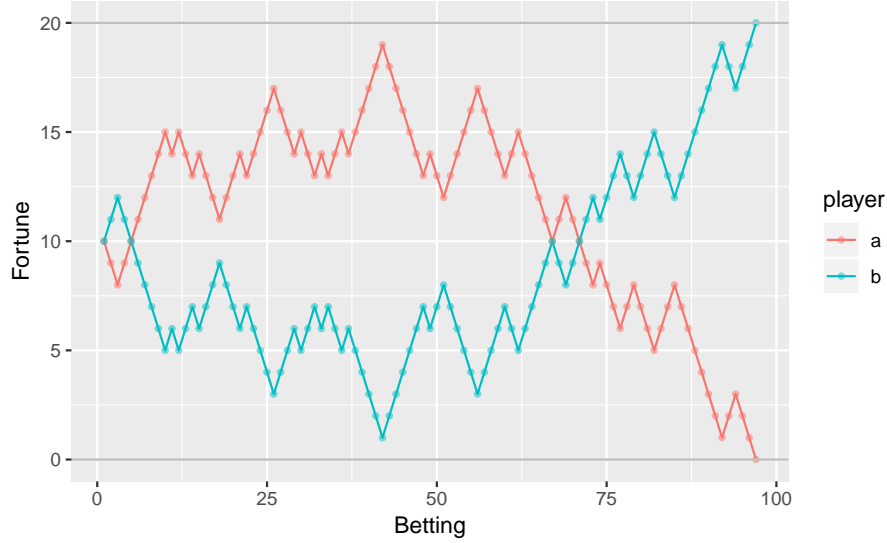


Figure 1.15: Sample path of Gambler's ruin model

In Figure 1.15, we can see the result of process with  $p = \frac{1}{2}$ . In fact, this process with probability 0.5 is also called *symmetric random walk*.

### 1.8.2 Homogeneous poisson process

**Definition 1.11** (Count process). A stochastic process  $\{N(t) : t \geq 0\}$  where  $N(t)$  is total number of events that occur by time  $t$  is called counting process.

1.  $N(t) \geq 0$
2.  $N(t) \in \mathbb{Z}$
3.  $s \leq t \Rightarrow N(s) \leq N(t)$
4. For  $s < t$ ,  $N(t) - N(s)$  = the number of events that occur in  $(s, t]$

*Poisson process* is one of this counting process.

**Definition 1.12** (Poisson process). The counting process  $\{N(t), t \geq 0\}$  is said to be a Poisson process with rate  $\lambda > 0$

1.  $N(0) = 0$
2.  $N(t) \perp\!\!\!\perp N(t+s) - N(t)$
3. Distribution of  $N(t+s) - N(t)$  is the same for all values of  $t$
4.  $\lim_{h \rightarrow 0} \frac{P(N(h)=1)}{h} = \lambda$
5.  $\lim_{h \rightarrow 0} \frac{P(N(h) \geq 2)}{h} = 0$

*Remark.*

$$\{N(t), t \geq 0\} \sim PP(\lambda) \Rightarrow N(t) \sim \text{Poisson}(\lambda t)$$

We can generate Poisson process using this relationship. However, it is slow. Thus, we find another way.

**Theorem 1.5** (Campbell's Theorem). Let  $\{N(t), t \geq 0\} \sim PP(\lambda)$ , let  $X_t$  be the interarrival time, and let  $S_t$  be the waiting time until  $t$ -th event, i.e.  $S_t := \sum_{i=1}^t X_i$ . Then

$$S_1, S_2, \dots, S_n \mid N(t) = n \stackrel{d}{=} (U_{(1)}, U_{(2)}, \dots, U_{(n)})$$

where  $U_i \sim \text{unif}(0, t)$ .

This Campbell's theorem gives solution to the PP generation.

**Algorithm 7:** Fast algorithm for Poisson Process

**input :** end time  $T$

- 1 Generate  $N \sim \text{Poisson}(\lambda T)$ ;
- 2 For  $N$ , generate  $U_1, \dots, U_N \stackrel{iid}{\sim} \text{unif}(0, T)$ ;
- 3 Sort  $U_1, \dots, U_N$  in ascending order, i.e.  $\{U_{(1)}, \dots, U_{(N)}\}$ ;
- 4 Set  $S_1 = U_{(1)}, \dots, S_N = U_{(N)}$ ;

**output:**  $\begin{bmatrix} 1 & S_1 \\ 2 & S_2 \\ \dots & \dots \\ N & S_n \end{bmatrix}$   
 $\Rightarrow \{N(S_n)\} \sim PP(\lambda) \text{ on } [0, T]$

```
rpp <- function(lambda, t0) {
  N <- rpois(1, lambda = lambda * t0)
  tibble(sn = runif(N) * t0) %>%
    arrange(sn) %>% # arrival time
    mutate(pp = 1:n()) # N(sn) ~ PP
}
```

```
rpp(lambda = 1, t0 = 50) %>%
  mutate(
    true_mean = sn, # sn * lambda
    true_sd = sqrt(sn) # sn * lambda
  ) %>%
  ggplot(aes(x = sn)) +
  geom_ribbon(
    aes(ymin = true_mean - true_sd, ymax = true_mean + true_sd),
    fill = "grey70",
    alpha = .5
  ) +
  geom_line(aes(y = true_mean), col = I("white"), size = 2) +
  geom_path(aes(y = pp)) +
  labs(
    x = "t",
    y = expression(N(t))
  )
```

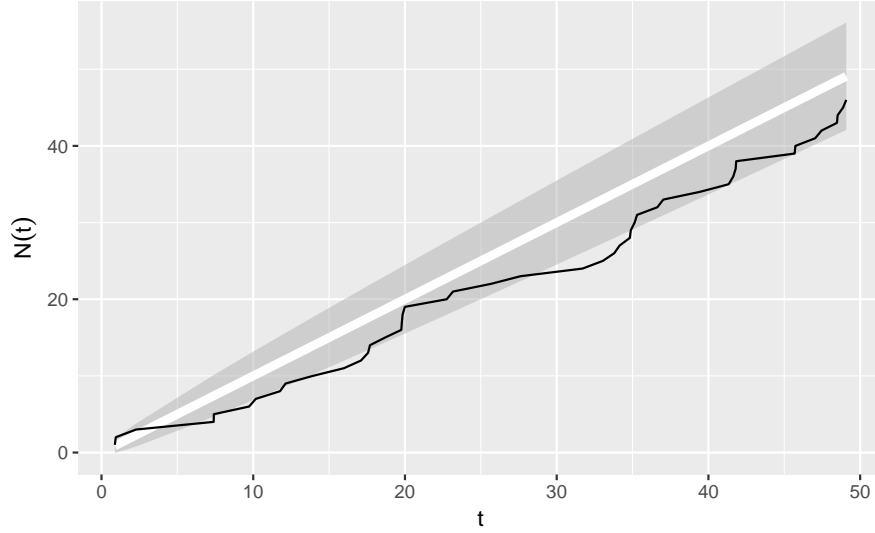


Figure 1.16: Sample path of Poisson process

### 1.8.3 Nonhomogeneous poisson process

Last section, what we have seen was homogeneous PP whose distribution does not depend on  $t$ . On the other hand, this condition can be broken.

**Definition 1.13** (Nonhomogeneous Poisson Process). The counting process  $\{N(t), t \geq 0\}$  is said to be a Nonhomogeneous Poisson process with rate  $\lambda(t) > 0$  if the third condition does not hold.

1.  $N(0) = 0$
2.  $N(t) \perp\!\!\!\perp N(t+s) - N(t)$
3.  $\lim_{h \rightarrow 0} \frac{P(N(h)=1)}{h} = \lambda$
4.  $\lim_{h \rightarrow 0} \frac{P(N(h) \geq 2)}{h} = 0$

$\lambda(t)$  is called intensity at time  $t$ .

$m(t) := \int_0^t \lambda(s) ds, t \geq 0$  is called mean-value function.

We can generate this NPP by Bernoulli splitting, which is called *thinning approach*.

**Lemma 1.1** (Thinning approach). *Choosing  $\lambda$  s.t.  $\forall t \leq T \lambda(t) \leq \lambda$ . If an event of  $PP(\lambda)$  counted with*

$$p(t) = \frac{\lambda(t)}{\lambda}$$

*then the process follows  $NPP(\lambda(t))$  on  $[0, T]$*

**Algorithm 8:** Thinning algorithm

```

1 Set  $t = 0, N = 0$ ;
2 repeat
3   Generate  $Y \sim \text{Exp}(\lambda)$ ;
4   Set  $t \leftarrow t + Y$ ;
5   Generate  $U \sim \text{unif}(0, 1)$ ;
6   if  $U \leq \frac{\lambda(t)}{\lambda}$  then
7     Set  $N \leftarrow N + 1$ ;
8     Set  $S(N) \leftarrow t$ ;
9 until  $t > T$ ;

```

```

npp <- function(lambda, t0, intensity) {
  t <- 0
  N <- 0
  S <- tibble(t = t, N = N)
  while (t <= t0) {
    t <- t + rexp(1, lambda)
    U <- runif(1)
    if (U <= intensity(t) / lambda) {
      N <- N + 1
      S <-
        S %>%
        bind_rows(tibble(t = t, N = N))
    }
  }
  S %>% slice(-1)
}

```

```

intensity <- function(x) {
  tt <- x %% 480
  rate <-
    case_when(
      tt >= 0 && tt <= 120 ~ .5,
      tt > 120 && tt <= 240 ~ 1,
      tt > 240 && tt <= 360 ~ 2,
      tt > 360 && tt <= 480 ~ 1.5
    )
  rate
}
#-----
npp(2, 480, intensity = intensity) %>%
  ggplot(aes(x = t, y = N)) +
  geom_path()

```

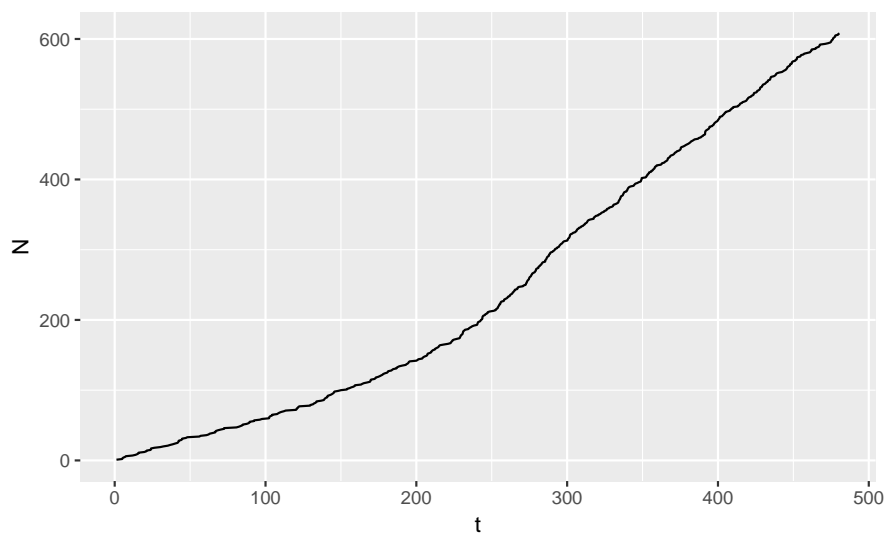


Figure 1.17: Sample path of nonhomogeneous poisson process



## Chapter 2

# Monte Carlo Integration and Variance Reduction

### 2.1 Monte Carlo Integration

Consider integration problem of a integrable function  $g(x)$ . We want to compute

$$\theta \equiv \int_a^b g(x)dx$$

For instance, standard normal cdf.

**Example 2.1** (Standard normal cdf). Compute values for

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt$$

It might be impossible to compute this integral with hand. So we implement *simulation* concept here, based on the following theorems.

**Theorem 2.1** (Weak Law of Large Numbers). Suppose that  $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2 < \infty)$ . Then

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu$$

Let  $g$  be a measurable function. Then

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \xrightarrow{p} g(\mu)$$

**Theorem 2.2** (Strong Law of Large Numbers). Suppose that  $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \sigma^2 < \infty)$ . Then

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \mu$$

Let  $g$  be a measurable function. Then

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \xrightarrow{a.s.} g(\mu)$$

### 2.1.1 Simple Monte Carlo estimator

**Theorem 2.3** (Monte Carlo Integration). *Consider integration (2.1). This can be approximated via appropriate pdf  $f(x)$  by*

$$\hat{\theta}_M = \frac{1}{N} \sum_{i=1}^N g(X_i)$$

Suppose that we have a distribution  $f(x) = I_{sptg}(x)$ , i.e. *uniform distribution*. Let  $sptg = (a, b)$ .

$$\begin{aligned} \theta &\equiv \int_{sptg} g(x) dx \\ &= \int_a^b g(x) dx \\ &= \int_0^1 g(a + (b-a)t)(b-a) dt \\ &\equiv \int_0^1 h(t) dt \\ &= \int_0^1 h(t) I_{(a,b)}(t) dt \\ &= E[h(U)] \quad U \sim \text{unif}(0, 1) \end{aligned} \tag{2.1}$$

By the Strong law of large numbers 2.2,

$$\frac{1}{n} \sum_{i=1}^n h(U_i) \xrightarrow{a.s.} E[h(U)] = \theta$$

where  $U \sim \text{unif}(0, 1)$ . Thus, what we have to do here are two things.

1. representing  $g$  as  $h$ .
2. generating lots of  $U_i$

Go back to Example 2.1.

*Solution.* Case 1:  $x > 0$

Since  $\Phi(x)$  is symmetry,

$$\Phi(0) = \frac{1}{2}$$

Fix  $x > 0$ .

$$\begin{aligned} \int_0^x \exp\left(-\frac{t^2}{2}\right) dt &= \int_0^x x \exp\left(-\frac{t^2}{2}\right) \frac{I_{(0,x)}(t)}{x} dt \\ &\approx \frac{1}{N} \sum_{i=1}^N x \exp\left(-\frac{U_i^2}{2}\right) \end{aligned}$$



with  $U_1, \dots, U_N \stackrel{iid}{\sim} \text{unif}(0, x)$ .

Case 2:  $x \leq 0$

Recall that  $\Phi(x)$  is symmetry.

Hence,

$$\hat{\Phi}(x) = \begin{cases} \frac{1}{\sqrt{2\pi}} \frac{1}{N} \sum_{i=1}^N x \exp\left(-\frac{U_i^2}{2}\right) + \frac{1}{2} \equiv \hat{\theta}(x) & x \geq 0 \\ 1 - \hat{\theta}(-x) & x < 0 \end{cases}$$

```
phihat <- function(x, y) {
  yi <- abs(y)
  theta <- mean(yi * exp(-x^2 / 2)) / sqrt(2 * pi) + .5
  ifelse(y >= 0, theta, 1 - theta)
}
```

Then compute  $\hat{\Phi}(x)$  for various  $x$  values.

```
phi_simul <- foreach(y = seq(.1, 2.5, length.out = 10), .combine = bind_rows) %do% {
  tibble(
    x = y,
    phi = pnorm(y),
    Phihat =
      tibble(x = runif(10000, max = y)) %>%
      summarise(cdf = phihat(x, y = y)) %>%
      pull()
  )
}
```

Table 2.1: Simple MC estimates of Normal cdf for each  $x$

x	pnorm	mc
0.100	0.540	0.540
0.367	0.643	0.643
0.633	0.737	0.737
0.900	0.816	0.816
1.167	0.878	0.878
1.433	0.924	0.923
1.700	0.955	0.958
1.967	0.975	0.976
2.233	0.987	0.987
2.500	0.994	0.990

### 2.1.2 Hit-or-Miss Monte Carlo

Hit-or-Miss approach is another way to evaluate integrals.

**Example 2.2** (Estimation of  $\pi$ ). Consider a circle in  $\mathbb{R}$  coordinate.

$$x^2 + y^2 = 1$$

Since  $y = \sqrt{1 - x^2}$ ,

$$\int_0^1 \sqrt{1-t^2} dt = \frac{\pi}{4} \quad (2.2)$$

By estimating Equation (2.2), we can estimate  $\pi$ , i.e.

$$\pi = 4 \int_0^1 \sqrt{1-t^2} dt$$

Simple MC integration can also be used.

$$\begin{aligned} \int_0^1 \sqrt{1-t^2} dt &= \int_0^1 \sqrt{1-t^2} I_{(0,1)}(t) dt \\ &\approx \frac{1}{N} \sum_{i=1}^N \sqrt{1-U_i^2} \end{aligned}$$

```
circ <- function(x) {
  4 * sqrt(1 - x^2)
}

tibble(x = runif(10000)) %>%
  summarise(mc_pi = mean(circ(x)))
#> # A tibble: 1 x 1
#>   mc_pi
#>   <dbl>
#> 1  3.14
```

On the other way, hit-or-miss MC method applies geometric probability.

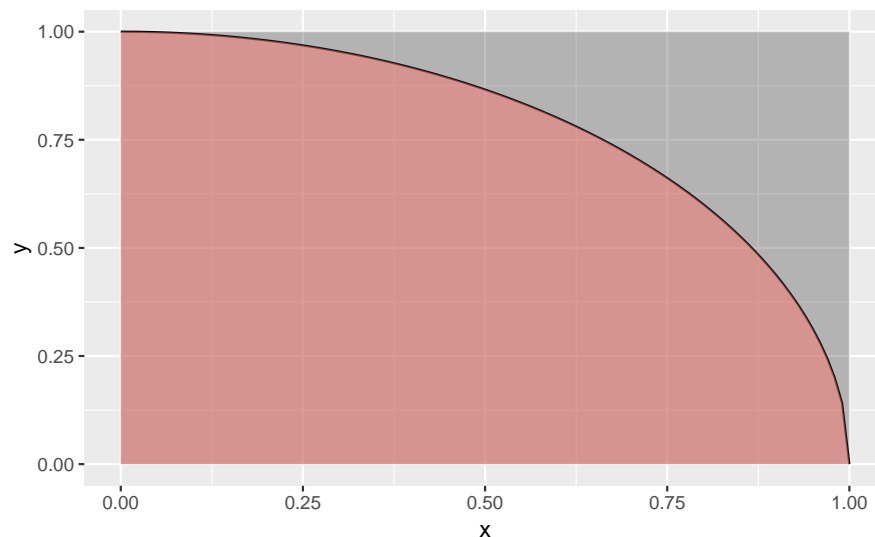


Figure 2.1: Hit-or-Miss

See Figure 2.1. From each coordinate, generate

- $X_i \stackrel{iid}{\sim} \text{unif}(0,1)$

- $Y_i \stackrel{iid}{\sim} \text{unif}(0, 1)$

Then the proportion of  $Y_i \leq \sqrt{1 - X_i^2}$  estimates  $\frac{\pi}{4}$ .

```
tibble(x = runif(10000), y = runif(10000)) %>%
  summarise(hitormiss = mean(y <= sqrt(1 - x^2)) * 4)
#> # A tibble: 1 x 1
#>   hitormiss
#>   <dbl>
#> 1      3.15
```

## 2.2 Variance and Efficiency

We have seen two approaches doing the same task. Now we want to *evaluate them*. Denote that simple Monte Carlo integration 2.3 is estimating the *expected value of some random variable*. Proportion, which approximates probability is expected value of identity function.

The common statistic that can evaluate estimators expected value might be their variances.

### 2.2.1 Variance

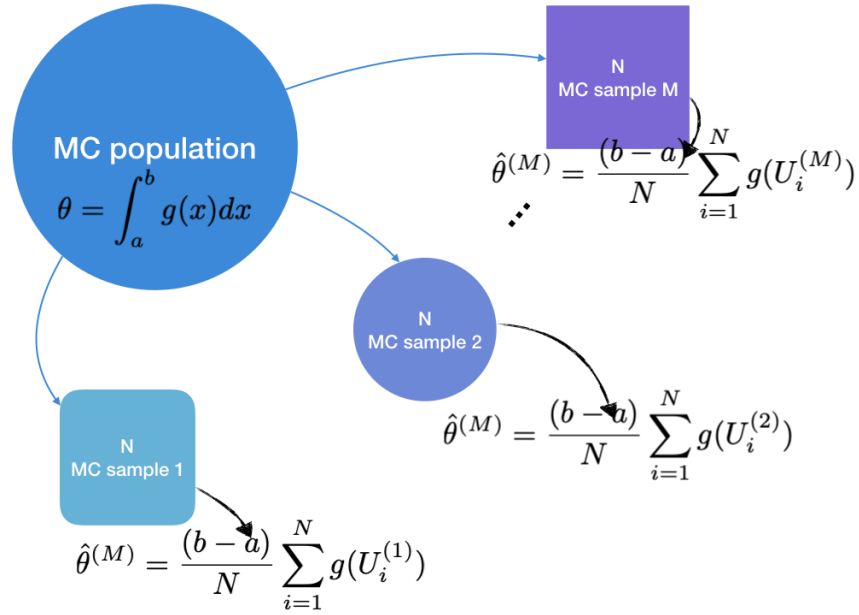
Note that variance of sample mean is  $\text{Var}(\overline{g(X)}) = \frac{\text{Var}(g(X))}{N}$ . This property is one of estimating variance of  $\hat{\theta}$ .

$$\widehat{\text{Var}}(\hat{\theta}) = \frac{1}{N} \left( \frac{1}{N} \sum_{i=1}^N (g(X_i) - \overline{g(X)}) \right) = \frac{1}{N^2} \sum_{i=1}^N (g(X_i) - \overline{g(X)}) \quad (2.3)$$

For example,

```
tibble(x = runif(10000)) %>%
  summarise(mc_pi = var(circ(x)) / 10000)
#> # A tibble: 1 x 1
#>   mc_pi
#>   <dbl>
#> 1 0.0000795
```

However, this *variance of sample mean* is used in situation when we are in sample limitation situation. We do not have to stick to this. Now, Generating samples as many as we want is possible. So we try another approach: *parametric bootstrap*.

Figure 2.2: Empirical distribution of  $\hat{\theta}$ 

See Figure 2.2. If we estimate  $E[g(U \sim \text{unif}(a, b))]$ , we can get  $\theta$ . Generate  $M$  samples  $\{U_1^{(j)}, \dots, U_N^{(j)}\}, j = 1, \dots, M$  from this  $U \sim \text{unif}(a, b)$ . In each sample, calculate MC estimates  $\hat{\theta}^{(j)}$ . Now we have  $M$  MC estimates  $\hat{\theta}$ . This gives empirical distribution of  $\hat{\theta}$ . By *drawing a histogram*, we can see the outline.

**Algorithm 9:** Variance of  $\hat{\theta}$ 

```

input :  $\theta = \int_a^b g(x)dx$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   | Generate  $U_1^{(m)}, \dots, U_N^{(m)} \stackrel{iid}{\sim} \text{unif}(a, b)$ ;
3   | Compute  $\hat{\theta}^{(j)} = \frac{(b-a)}{N} \sum g(U_i^{(j)})$ ;
4 end
5  $\bar{\hat{\theta}} = \frac{1}{M} \sum \hat{\theta}^{(j)}$ ;
6  $\widehat{Var}(\hat{\theta}) = \frac{1}{M-1} \sum (\hat{\theta}^{(j)} - \bar{\hat{\theta}})^2$ ;
output:  $\widehat{Var}(\hat{\theta})$ 

```

Since we have to generate large size of data, `data.table` package will be used.

```
library(data.table)
```

Group operation can be used. Additional column (`sam`) would indicate group, and for each group MC operation would be processed. The following is the function generating `data.table` before group operation.

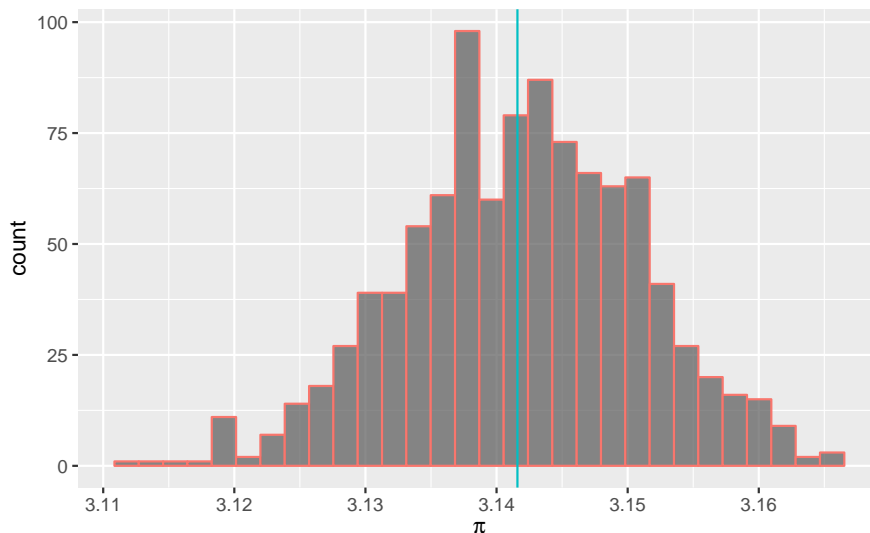
```

mc_data <- function(rand, N = 10000, M = 1000, char = "s", ...) {
  data.table(
    u = rand(n = N * M, ...),
    sam = gl(M, N, labels = paste0(char, 1:M))
  )
}

```

```
pi_mc <-
  mc_data(runif)[,
    .(mc_pi = mean(circ(u))),
    keyby = sam]

pi_mc %>%
  ggplot(aes(x = mc_pi)) +
  geom_histogram(bins = 30, col = gg_hcl(1), alpha = .7) +
  xlab(expression(pi)) +
  geom_vline(xintercept = pi, col = gg_hcl(2)[2])
```

Figure 2.3: Empirical distribution of  $\hat{\pi}$  by simple MC

As in Algorithm 9, we can compute the variance as below.

```
(mc_var <-
  pi_mc[,
    .(mc_variance = var(mc_pi))])
#>   mc_variance
#> 1:   8.09e-05
```

On the other hand, we need to generate two sets of random numbers for hit-or-miss MC.

```
pi_hit <-
  mc_data(runif)[
    , u2 := runif(10000 * 1000)
  ][,
    .(hitormiss = mean(u2 <= sqrt(1 - u^2)) * 4),
    keyby = sam]

pi_mc[pi_hit] %>%
  melt(id.vars = "sam", variable.name = "hat") %>%
  ggplot(aes(x = value, fill = hat)) +
  geom_histogram(bins = 30, alpha = .5, position = "identity") +
  xlab(expression(pi)) +
  geom_vline(xintercept = pi, col = I("red")) +
  scale_fill_discrete(
```

```

name = "MC",
labels = c("Simple", "Hit-or-Miss")
)

```

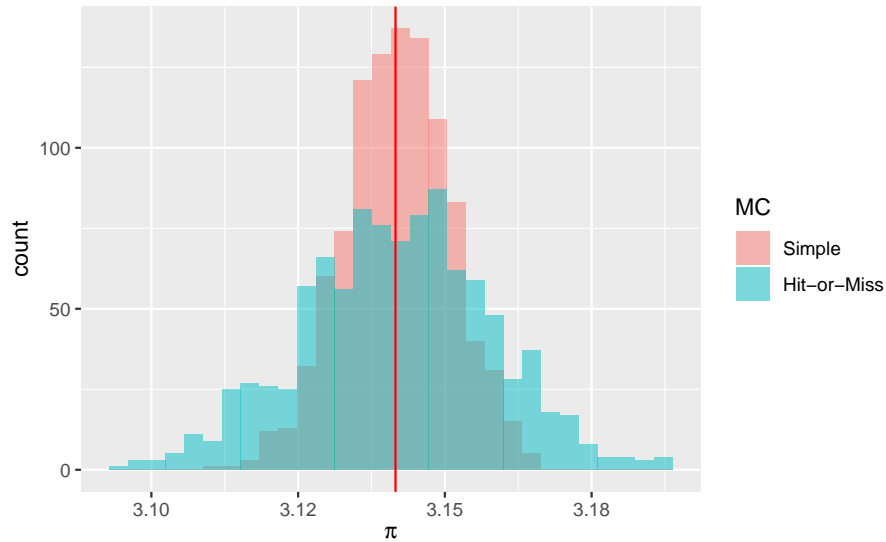


Figure 2.4: Simple MC and Hit-or-Miss MC

```

(hit_var <-
  pi_hit[,
    .(hit_variance = var(hitormiss))])
#> hit_variance
#> 1: 0.000258

```

### 2.2.2 Efficiency

See Figure 2.4. It is obvious that Hit-or-Miss estimate produces larger variance than simple MC.

**Definition 2.1** (Efficiency). Let  $\hat{\theta}_1$  and  $\hat{\theta}_2$  be two estimators for  $\theta$ . Then  $\hat{\theta}_1$  is more efficient than  $\hat{\theta}_2$  if

$$\frac{\text{Var}(\hat{\theta}_1)}{\text{Var}(\hat{\theta}_2)} < 1$$

In other words, if  $\hat{\theta}_1$  has smaller variance than  $\hat{\theta}_2$ , then  $\hat{\theta}_1$  is said to be efficient, which is preferable.

Table 2.2: Simple MC versus Hit-or-Miss

SimpleMC	Hit-or-Miss	SimpleMCEfficiency
0	0	TRUE

## 2.3 Variance Reduction

Consider Equation (2.3) based on  $\text{Var}(\hat{\theta}) = \frac{\sigma^2}{N}$ . This variance can always be reduced by adding  $N$ . But we want to reduce variance less computationally.

### 2.3.1 Antithetic Variables

Consider correlated random variables  $U_1$  and  $U_2$ . Then we have

$$\text{Var}\left(\frac{U_1 + U_2}{2}\right) = \frac{1}{4}\left(\text{Var}(U_1) + \text{Var}(U_2) + 2\text{Cov}(U_1, U_2)\right)$$

See the last term  $\text{Cov}(U_1, U_2)$ . If we generate  $U_{i1}$  and  $U_{i2}$  negatively correlated, we can get reduced variance than previous i.i.d. sample

$$\text{Var}\left(\frac{U_1 + U_2}{2}\right) = \frac{1}{4}\left(\text{Var}(U_1) + \text{Var}(U_2)\right)$$

**Lemma 2.1.**  $U$  and  $1 - U$  are identically distributed, but negatively correlated.

1.  $U \sim \text{unif}(0, 1) \Leftrightarrow 1 - U \sim \text{unif}(0, 1)$
2.  $\text{Corr}(U, 1 - U) = -1$

This is well-known property of uniform distribution. Instead of generating  $N$  uniform numbers, try  $\frac{N}{2} U_i$  and make corresponding  $\frac{N}{2} 1 - U_i$ . This sequence becomes negatively correlated, so we can reduce the variance as mentioned.

When can we replace previous numbers with these *antithetic variables*? We usually plug-in the numbers in some function  $h$  to get Monte carlo integration. The thing is, our target is  $h$ , not  $U$ .  $h(U)$  and  $h(1 - U)$  should *still be negatively correlated*. Hence,  $h$  should be *monotonic function*.

**Corollary 2.1.** If  $g = g(X_1, \dots, X_n)$  is monotone, then

$$Y = g(F_X^{-1}(U_1), \dots, F_X^{-1}(U_n))$$

and

$$Y' = g(F_X^{-1}(1 - U_1), \dots, F_X^{-1}(1 - U_n))$$

are negatively correlated.

**Algorithm 10:** Variance of  $\hat{\theta}$  using antithetic variables

<pre> <b>input</b> : <math>h</math> : monotonic 1 <b>for</b> <math>m \leftarrow 1</math> <b>to</b> <math>M</math> <b>do</b> 2   Generate <math>U_{1,1}^{(m)}, \dots, U_{\frac{N}{2},1}^{(m)} \stackrel{iid}{\sim} \text{unif}(0, 1)</math>; 3   Set <math>U_{i,2}^{(m)} := 1 - U_{i,1}^{(m)} \stackrel{iid}{\sim} \text{unif}(0, 1)</math>; 4   <math>\{U_i^{(m)}\}_1^N = \{U_{1,1}^{(m)}, \dots, U_{\frac{N}{2},2}^{(m)}\}</math>; 5   <math>\hat{\theta}^{(j)} = \frac{1}{N} \sum h(U_i^{(j)})</math>; 6 <b>end</b> 7 <math>\bar{\hat{\theta}} = \frac{1}{M} \sum \hat{\theta}^{(j)}</math>; 8 <math>\widehat{\text{Var}}(\hat{\theta}) = \frac{1}{M-1} \sum (\hat{\theta}^{(j)} - \bar{\hat{\theta}})^2</math>; <b>output</b>: <math>\widehat{\text{Var}}(\hat{\theta})</math> </pre>
--

Check again Example 2.1. We have try to calculate

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt$$

using simple monte carlo. To make the support  $(0, 1)$ , let  $y = \frac{t}{x}$  be a change of variable. Then

$$\begin{aligned} \int_0^x \exp\left(-\frac{t^2}{2}\right) dt &= \int_0^1 x \exp\left(-\frac{(xy)^2}{2}\right) dy \\ &\approx \frac{1}{N} \sum_{i=1}^N x \exp\left(-\frac{(xU_i)^2}{2}\right) \end{aligned}$$

```
phiunif <- function(x, y) {
  yi <- abs(y)
  theta <- mean(yi * exp(-(yi * x)^2 / 2)) / sqrt(2 * pi) + .5
  ifelse(y >= 0, theta, 1 - theta)
}
```

Consider  $\Phi(2)$ .

```
phi2 <-
  mc_data(runif)[,
    .(p2 = phiunif(u, y = 2)),
    keyby = sam]
```

Now apply antithetic variables.

```
phi2_anti <-
  mc_data(runif, N = 10000 / 2)[,
    u2 := 1 - u] %>%
  melt(id.vars = "sam", value.name = "U") %>%
  .[,
    .(anti_p2 = phiunif(U, y = 2)),
    keyby = sam]

phi2[phi2_anti] %>%
  melt(id.vars = "sam", variable.name = "hat") %>%
  ggplot(aes(x = value, fill = hat)) +
  geom_histogram(bins = 30, alpha = .5, position = "identity") +
  xlab(expression(pi)) +
  geom_vline(xintercept = pnorm(2), col = I("red")) +
  scale_fill_discrete(
    name = "MC",
    labels = c("Simple", "Antithetic")
  )
```



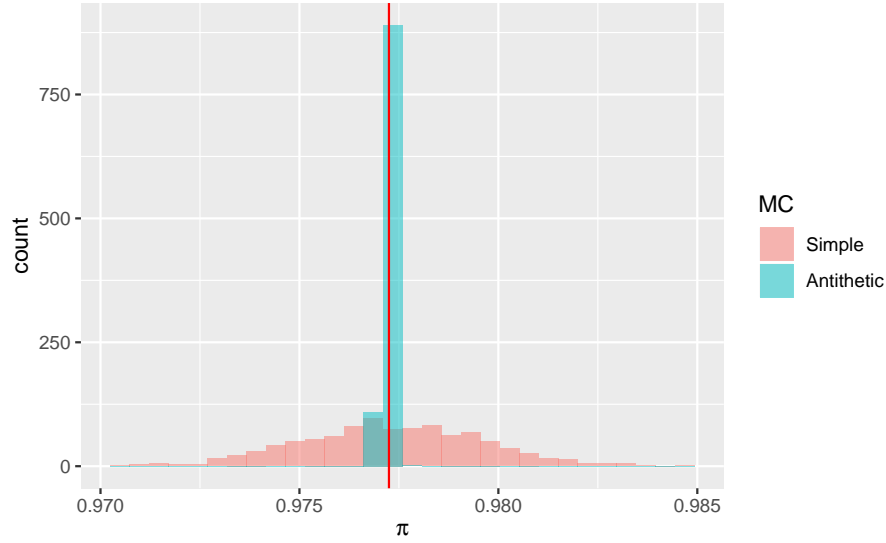


Figure 2.5: Use of antithetic variables

Obviously, variance has been reduced.

```
phi2[phi2_anti] %>%
  melt(id.vars = "sam", variable.name = "hat") %>%
  .[,
    .(variance = var(value)),
    by = hat]
#>      hat variance
#> 1:    p2 5.2e-06
#> 2: anti_p2 1.5e-08
```

### 2.3.2 Control Variates

Recall that we are trying to estimate  $\theta = Eg(X)$  here in MC integration. Consider other output random variable. Suppose that  $\mu_f \equiv Ef(Y)$  is known. It is obvious that

$$\hat{\theta}_c = g(X) + c(f(Y) - \mu_f) \quad (2.4)$$

is an unbiased estimator for  $\theta$  for any  $c \in \mathbb{R}$ . Then we have

$$Var\hat{\theta}_c = Var g(X) + c^2 Var f(X) + 2c Cov(g(X), f(X)) \quad (2.5)$$

Recall that our goal is to minimize this  $Var\hat{\theta}_c$ . What value of  $c$  is to be determined? Note that Equation (2.5) is quadratic function of  $c$ .

$$\begin{aligned} Var\hat{\theta}_c &= Var f(X)c^2 + 2c Cov(g(X), f(X)) + Var g(X) \\ &= Var f(X) \left( c + \frac{Cov(g(X), f(X))}{Var f(X)} \right)^2 + Var g(X) - \frac{Cov(g(X), f(X))^2}{Var f(X)} \end{aligned} \quad (2.6)$$

From Equation (2.6), the variance is minimized at

$$c^* = -\frac{\text{Cov}(g(X), f(X))}{\text{Var}f(X)} \quad (2.7)$$

with minimum variance

$$\text{Var}\hat{\theta}_{c^*} = \text{Var}g(X) - \frac{\text{Cov}(g(X), f(X))^2}{\text{Var}f(X)}$$

By this, we can reduce the variance of estimation as much as possible (using  $f(X)$ ). Here,  $f(X)$  is called a *control variate* for  $g(X)$ .

**Algorithm 11:** Variance of  $\hat{\theta}$  using control variables

```

input :  $g$ , control variate  $f$  with mean  $\mu_f$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   Generate  $U_1^{(m)}, \dots, U_N^{(m)} \stackrel{iid}{\sim} \text{unif}$ ;
3   Set  $g = g(U_i)$  and  $f = f(U_i)$ ;
4   Compute  $\hat{c}^{*(m)} = -\frac{\widehat{\text{Cov}}(g, f)}{\widehat{\text{Var}}(f)}$ ;
5    $\hat{\theta}_{c^*}^{(j)} = g + \hat{c}^{*(m)}(f - \mu_f)$ ;
6 end
7  $\bar{\hat{\theta}} = \frac{1}{M} \sum \hat{\theta}_{c^*}^{(j)}$ ;
8  $\widehat{\text{Var}}(\hat{\theta}) = \frac{1}{M-1} \sum (\hat{\theta}_{c^*}^{(j)} - \bar{\hat{\theta}})^2$ ;
output:  $\widehat{\text{Var}}(\hat{\theta})$ 
```

**Example 2.3** (Variance reduction by control variate). Apply each simple MC, antithetic variate, and control variate to

$$\int_0^1 e^x dx$$

```

N <- 100
M <- 1000
```

Denote that the true value is

$$\int_0^1 e^x dx = e - 1 = 1.718$$

We might compare each estimate to this.

*Solution* (Simple MC). We only need  $U \sim \text{unif}(0, 1)$ .

$$\begin{aligned}
\theta &= \int_0^1 e^x dx \\
&= \int_0^1 e^x I_{(0,1)}(x) dx \\
&\approx \frac{1}{N} \sum_{i=1}^N e^{u_i}, \quad u_i \stackrel{iid}{\sim} \text{unif}(0, 1)
\end{aligned} \quad (2.8)$$

```
theta_sim <-
  mc_data(runif, N = N, M = M)[,
    .(mc = mean(exp(u))),
    keyby = sam]
```

```
#>      sam  mc
#>   1:   s1 1.70
#>   2:   s2 1.71
#>   3:   s3 1.76
#>   4:   s4 1.71
#>   5:   s5 1.73
#> ---
#>  996: s996 1.79
#>  997: s997 1.82
#>  998: s998 1.68
#>  999: s999 1.73
#> 1000: s1000 1.76
```

*Solution* (Antithetic variate). For  $N' = \frac{N}{2}$ ,

Consider  $u_1, \dots, u_{N'} \stackrel{iid}{\sim} \text{unif}(0, 1)$  and  $1 - u_1, \dots, 1 - u_{N'} \stackrel{iid}{\sim} \text{unif}(0, 1)$ .

See Equation (2.8). Then we can compute antithetic estimator by

$$\begin{aligned}\hat{\theta}_A &= \frac{1}{N} \sum_{i=1}^{N/2} (e^{u_i} + e^{1-u_i}) \\ &= \frac{1}{N/2} \sum_{i=1}^{N/2} \left( \frac{e^{u_i} + e^{1-u_i}}{2} \right) \\ &= \text{sample mean}\end{aligned}$$

```
theta_anti <-
  mc_data(runif, N = N / 2, M = M)[,
    u2 := 1 - u][,
    .(anti = mean((exp(u) + exp(u2)) / 2)),
    keyby = sam]
```

Now look at the results of the two.

```
theta_sim[theta_anti] %>%
  melt(id.vars = "sam", variable.name = "simul", value.name = "integral") %>%
  ggplot(aes(x = integral, fill = simul)) +
  geom_histogram(bins = 30, position = "identity", alpha = .5) +
  xlab(expression(theta)) +
  geom_vline(xintercept = exp(1) - 1, col = I("red")) +
  scale_fill_discrete(
    name = "MC",
    labels = c("Simple", "Antithetic")
  )
```

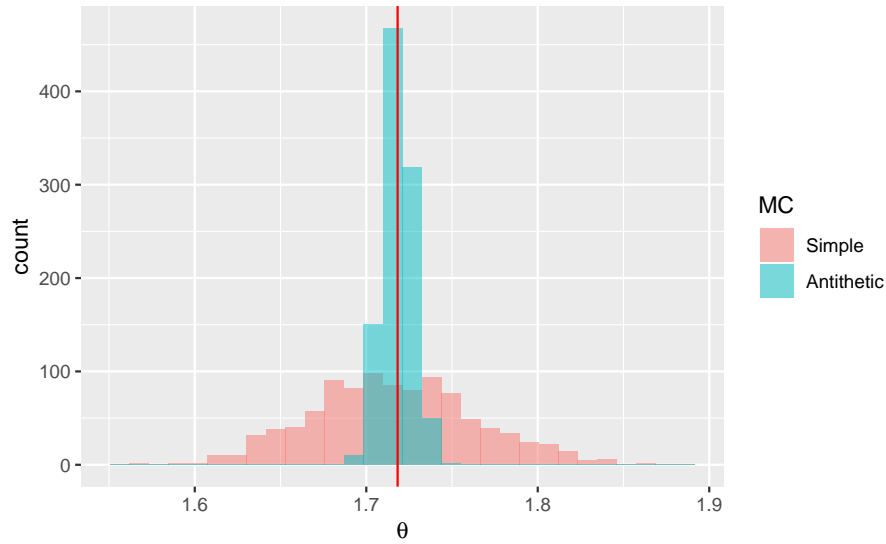


Figure 2.6: Antithetic variate estimator

It is clear that antithetic variate have reduced variance.

*Solution* (Control variate). Consider

$$g(U) = e^U$$

and

$$f(U) = U$$

with  $U \sim \text{unif}(0, 1)$ .

Note that

$$E(U) = \frac{1}{2}$$

Then

$$\hat{\theta}_C = e^U + c \left( U - \frac{1}{2} \right)$$

is an unbiased estimator of  $\theta = \int_0^1 e^x dx$ .

To reduce variance, we need to set  $c$  to be

$$c^* = - \frac{\text{Cov}(e^U, U)}{\text{Var}(U)}$$

Since we do not know the exact number, we estimate this from each Monte Carlo sample.

```
theta_con <-
  mc_data(runif, N = N, M = M)[,
    chat := - cov(exp(u), u) / var(u),
```

```

      by = sam[,
        .(con = mean(exp(u) + chat * (u - 1 / 2))),
        keyby = sam]

thetahat <- theta_sim[theta_anti][theta_con]

thetahat %>%
  melt(id.vars = "sam", variable.name = "simul", value.name = "integral") %>%
  ggplot(aes(x = integral, fill = simul)) +
  geom_histogram(bins = 30, position = "identity", alpha = .5) +
  xlab(expression(theta)) +
  geom_vline(xintercept = exp(1) - 1, col = I("red")) +
  scale_fill_discrete(
    name = "MC",
    labels = c("Simple", "Antithetic", "Control")
  )

```

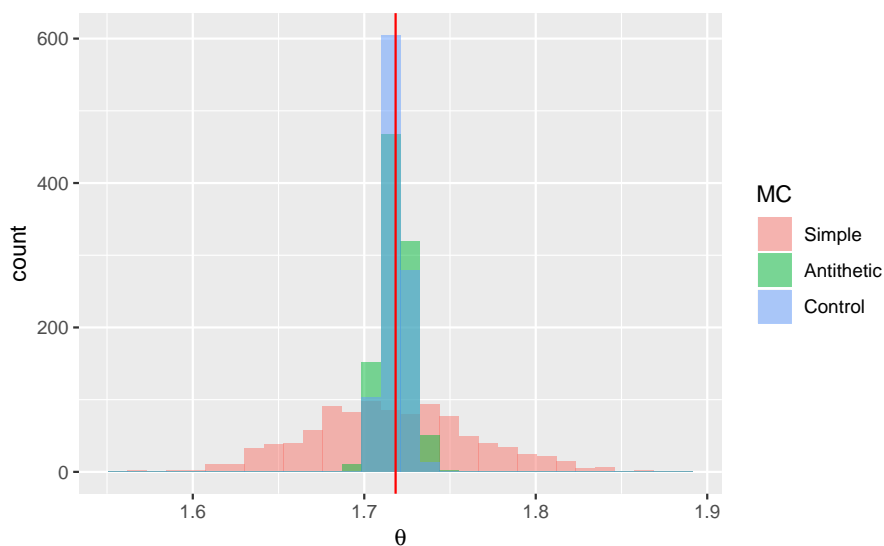


Figure 2.7: Use of Control variable

It looks like control variate have less variance, but what is more important is that both methods successfully have reduced it.

```

thetahat[,
  lapply(.SD, sd),
  ,.SDcols = -"sam"]
#>      mc      anti      con
#> 1: 0.049 0.00872 0.00643

```

### 2.3.3 Antithetic variate as control variate

Both antithetic variate and control variate reduce variance using covariance between two random variables. Actually, *antithetic variate is a special case of control variate*. See Equation (2.4).

**Lemma 2.2.** *Control variate estimator is a linear combination of unbiased estimators of  $\theta$ .*

Consider any two unbiased estimator  $\hat{\theta}_1$  and  $\hat{\theta}_2$  for  $\theta = Eg(X)$ . Build control variate as following.

$$\hat{\theta}_c = c\hat{\theta}_1 + (1-c)\hat{\theta}_2$$

It is obvious that  $\hat{\theta}_c$  is also unbiased of  $\theta$  for every  $c \in \mathbb{R}$ .

$$Var(\hat{\theta}_c) = Var(\hat{\theta}_2) + c^2 Var(\hat{\theta}_1 - \hat{\theta}_2) + 2c Cov(\hat{\theta}_2, \hat{\theta}_1 - \hat{\theta}_2) \quad (2.9)$$

Let  $\hat{\theta}_1$  and  $\hat{\theta}_2$  be antithetic variate choice. Recall that *antithetic variate* give that for  $\hat{\theta}_1$  and  $\hat{\theta}_2$ ,

$$\hat{\theta}_1, \hat{\theta}_2 \sim IID, \quad Corr(\hat{\theta}_1, \hat{\theta}_2) = -1$$

It follows that

$$Cov(\hat{\theta}_1, \hat{\theta}_2) = -Var(\hat{\theta}_1)$$

and that

$$Var(\hat{\theta}_c) = (4c^2 - 4c + 1)Var(\hat{\theta}_1)$$

Hence, it leads to choosing optimal

$$\hat{\theta}_{c^*} = \frac{\hat{\theta}_1 + \hat{\theta}_2}{2}$$

which we have been used in antithetic variate.

### 2.3.4 Several control variates

To summarize, control variate try to reduce variance by combining unbiased estimatos of the target parameter. We have used one variate  $f(X)$ . It might be possible to extend to multiple variates, so to speak,  $f_1(X), \dots, f_k(X)$ . Thanks to the linearity of expectation,

$$\hat{\theta}_c = g(X) + \sum_{i=1}^k c_i (f_i(X) - \mu_i)$$

is also unbiased estimator, where  $\mu_i = Ef_i(X)$ . How to get each  $c_i^*$ ? Rather than using variance and covariance, we can *fitting linear regression*.

### 2.3.5 Control variates and regression

See Equation (2.4) and Equation (2.7). It can be found that we were estimating linear regression coefficient as LSE.

**Lemma 2.3** (Least squares estimator). *Consider  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ . Then*

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (Y_i - \bar{Y})} = \frac{\widehat{Cov}(X, Y)}{\widehat{Var}(Y)}$$

Control variate estimator  $\hat{\theta}_c = g(X) + c(f(Y) - \mu_f)$  can be expressed in regression model as

$$Eg(X) = \beta_0 + \beta_1 Ef(X)$$

Then

$$\hat{\beta}_1 = \text{LSE of } g(X) \text{ on } f(X) = \frac{\widehat{Cov}(g(X), f(X))}{\widehat{Var}(f(X))} = -\hat{c}^* \quad (2.10)$$

Note that

$$\hat{\beta}_0 = \overline{g(X)} + \hat{c}^* \overline{f(X)}$$

This matches to  $\hat{\theta}_{c^*}$  in previous section.

$$\hat{\beta}_0 + \hat{\beta}_1 \mu_f = \overline{g(X)} + \hat{c}^* (\overline{f(X)} - \mu_f) = \hat{\theta}_{c^*} \quad (2.11)$$

Also, we can get the error variance estimate

$$\hat{\sigma}^2 = \widehat{Var}(X + \hat{c}^* Y) = MSE$$

and

$$\widehat{Var}\hat{\theta}_c^* = \frac{\hat{\sigma}^2}{N}$$

From Example 2.3, we can change the code computing  $c^* = \text{cov}(\exp(u), u) / \text{var}(u)$  to `lm(exp(u) ~ u)$coef[2]`.

```
mc_data(runif, N = N, M = M) %>%
  .[,
    chat := lm(exp(u) ~ u)$coef[2],
    by = sam] %>%
  .[,
    .(con = mean(exp(u) + chat * (u - 1 / 2))),
    by = sam]
#>      sam  con
#>   1:   s1 1.82
#>   2:   s2 1.93
#>   3:   s3 1.62
#>   4:   s4 1.62
#>   5:   s5 1.76
#> ---
#> 996: s996 1.46
#> 997: s997 1.72
#> 998: s998 1.61
#> 999: s999 1.59
#> 1000: s1000 1.71
```

In fact, we can use Equation (2.11) directly: `predict(lm, newdata = data.frame(u = mean(u)))`.

**Algorithm 12:** Control variables and regression

**input** :  $g$ , control variate  $f$  with mean  $\mu_f$   
**1 for**  $m \leftarrow 1$  **to**  $M$  **do**  
**2**   Generate  $U_1^{(m)}, \dots, U_N^{(m)} \stackrel{iid}{\sim} unif$ ;  
**3**   Set  $g = g(U_i)$  and  $f = f(U_i)$ ;  
**4**   Regression  $g \sim f$ ;  
**5**   Predict the regression at  $\bar{U}^{(m)}$ . It is  $\hat{\theta}_{c^*}^{(j)}$ ;  
**6 end**  
**7**    $\hat{\theta} = \frac{1}{M} \sum \hat{\theta}_{c^*}^{(j)}$ ;  
**8**    $\widehat{Var}(\hat{\theta}) = \frac{1}{M-1} \sum (\hat{\theta}_{c^*}^{(j)} - \hat{\theta})^2$ ;  
**output:**  $\widehat{Var}(\hat{\theta})$

```
mc_data(runif, N = N, M = M) [,
                                     .(con = predict(lm(exp(u) ~ u, data = .SD),
                                                         newdata = data.table(u = 1 / 2))),
                                     by = sam]

#>      sam con
#> 1:    s1 1.73
#> 2:    s2 1.71
#> 3:    s3 1.71
#> 4:    s4 1.72
#> 5:    s5 1.72
#> ---
#> 996: s996 1.72
#> 997: s997 1.72
#> 998: s998 1.72
#> 999: s999 1.72
#> 1000: s1000 1.71
```

Now, how to deal with multiple control variates?

$$X = \beta_0 + \sum_{i=1}^k \beta_i Y_i + \epsilon$$

Using *multiple linear regression* model, we can choose optimal  $c^*$  and estimate control variate estimate.

## 2.4 Importance Sampling

Simple MC computes

$$\int_A g(x)f(x)dx = Eg(X) = \theta$$

for some density function  $f$ . This method uses random number from  $f$  itself so that

$$\int_A g(x)f(x)dx \approx \frac{1}{N} \sum_{i=1}^N g(X_i)$$



where  $X_1, \dots, X_N \stackrel{iid}{\sim} f$ . This is why MC integration is called *direct sampling*. Sometimes, however, we face unknown distribution. In this case, generating from  $f$  directly is not easy. Even we can, it can be inefficient. The solution is *indirect method*: draw a sample from another pdf  $h$ . This is called **importance sampling**.

### 2.4.1 Importance sampling

Consider MC integration as before.

$$\int_A g(x)f(x)dx = Eg(X) = \theta$$

How about uniform random number set with simple MC as before? However, uniform random numbers does not apply to unbounded intervals. When the target function is not that uniform, especially, generating numbers uniformly can be inefficient.

$$\begin{aligned} E_f g(X) &= \int_A g(x)f(x)dx \\ &= \int_A g(x) \frac{f(x)}{\phi(x)} \phi(x)dx, \quad \phi : \text{density on } A \\ &= E_\phi \frac{g(X)f(X)}{\phi(X)} \\ &\approx \frac{1}{N} \sum_{i=1}^N \frac{g(X_i)f(X_i)}{\phi(X_i)}, \quad X_i \stackrel{iid}{\sim} \phi \end{aligned} \tag{2.12}$$

Here,  $\phi$  is called the *envelope* or the *importance sampling function*. This is just simple arithmetic, so it is possible to choose any density  $\phi$ . However, we should take good one. Typically, one should select  $\phi$  so that

$$\phi(x) \approx |g(x)|f(x) \quad \text{on } A \tag{2.13}$$

with finite variance.

**Example 2.4** (Choice of importance function). Obtain MC estimate of

$$\int_0^1 \frac{e^{-x}}{1+x^2} dx$$

by importance sampling.

```
g_target <- function(x) {
  exp(-x - log(1 + x^2)) * (x > 0) * (x < 1)
}
```

Consider candidate envelopes

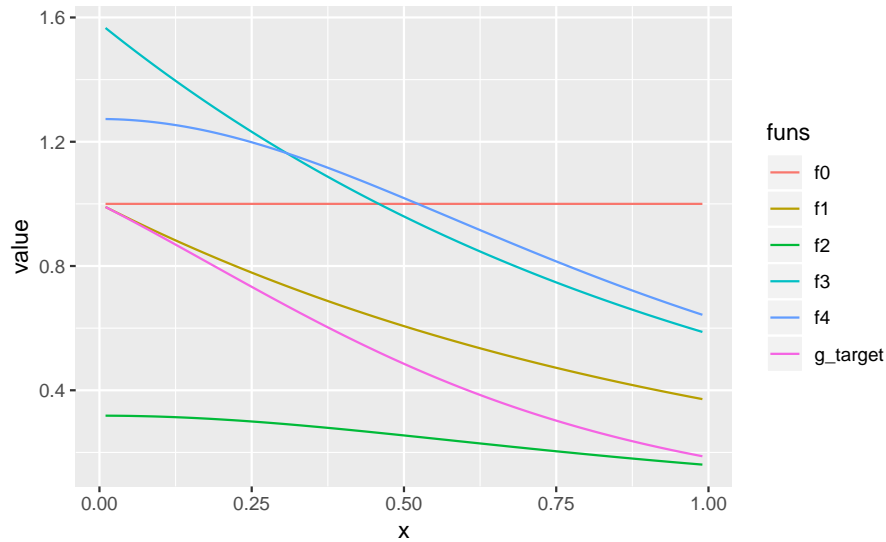
$$\left\{ \begin{array}{ll} \phi_0(x) = 1, & 0 < x < 1 \\ \phi_1(x) = e^{-x}, & 0 < x < \infty \\ \phi_2(x) = \frac{1}{\pi(1+x^2)}, & x \in \mathbb{R} \\ \phi_3(x) = \frac{e^{-x}}{1-e^{-1}}, & 0 < x < 1 \\ \phi_4(x) = \frac{4}{\pi(1+x^2)} & 0 < x < 1 \end{array} \right.$$

```

f0 <- function(x) {
  (x > 0) * (x < 1)
}
#-----
f1 <- function(x) {
  exp(-x) * (x > 0)
}
#-----
f2 <- function(x) {
  1 / (pi * (1 + x^2))
}
#-----
f3 <- function(x) {
  exp(-x) / (1 - exp(-1)) * (x > 0) * (x < 1)
}
#-----
f4 <- function(x) {
  4 / (pi * (1 + x^2)) * (x > 0) * (x < 1)
}

tibble(x = seq(.01, .99, by = .01)) %>%
  mutate_all(.funs = list(~g_target, ~f0, ~f1, ~f2, ~f3, ~f4)) %>%
  gather(~x, key = "funs", value = "value") %>%
  ggplot(aes(x = x, y = value, colour = funs)) +
  geom_path()

```

Figure 2.8: Importance functions  $\phi_0, \dots, \phi_4$ 

Each importance function is drawn in Figure 2.8.  $f_1$  shows similar patterns to  $g$ .

### 2.4.2 Variance in importance sampling

From Equation (2.12),

$$\theta = \int_A g(x) dx = \int_A \frac{g(x)}{\phi(x)} \phi(x) dx = E \left[ \frac{g(X)}{\phi(X)} \right] \approx \frac{1}{N} \sum \frac{g(X_i)}{\phi(X_i)}$$

where  $X_1, \dots, X_N \stackrel{iid}{\sim} \phi$ . Then

$$\begin{aligned} \text{Var} \hat{\theta} &= E \hat{\theta}^2 - (E \hat{\theta})^2 \\ &= \int_A \left( \frac{g(x)}{\phi(x)} \right)^2 \phi(x) dx - \theta^2 \\ &= \int_A \frac{g(x)^2}{\phi(x)} dx - \theta^2 \end{aligned}$$

Hence, the minimum variance

$$\left( \int_A |g(x)| dx \right)^2 - \theta^2$$

is obtained when

$$\phi(x) = \frac{|g(x)|}{\int_A |g(x)| dx}$$

But we do not know the value of denominator. It might be hard to get the exact function giving the minimum variance, but choosing  $\phi$  close to the shape of  $|g|$  would produce good result. To check our criterion (2.13) more clearly, compute  $\frac{g}{\phi_i}$ .

```
tibble(x = seq(.01, .99, by = .01)) %>%
  mutate_all(.funs = list(~g_target, ~f0, ~f1, ~f2, ~f3, ~f4)) %>%
  gather(-x, -g_target, key = "funs", value = "value") %>%
  mutate(value = g_target / value) %>%
  ggplot(aes(x = x, y = value, colour = funs)) +
  geom_path()
```

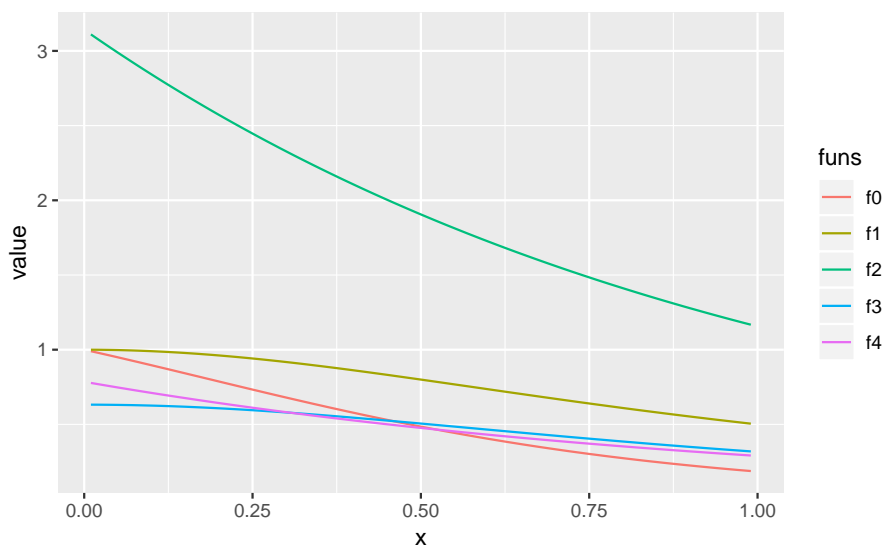


Figure 2.9: Ratio  $\frac{g}{\phi_i}$

What is the closest to 1?  $f_1$ , of course. Would this function produce the best result, i.e. variance?

```

theta_imp0 <-
  mc_data(runif, N = 100)[,
    .(phi0 = mean(g_target(u) / f0(u))),
    keyby = sam]

theta_imp1 <-
  mc_data(rexp, N = 100, rate = 1)[,
    .(phi1 = mean(g_target(u) / f1(u))),
    keyby = sam]

rf2 <- function(n) {
  x <- rcauchy(n)
  x[(x > 1) | (x < 0)] <- 2 # catch overflow errors in g
  x
}
theta_imp2 <-
  mc_data(rf2, N = 100)[,
    .(phi2 = mean(g_target(u) / f2(u))),
    keyby = sam]

rf3 <- function(n) {
  u <- runif(n)
  x <- -log(1 - u * (1 - exp(-1))) # inverse transformation method
  x
}
#-----
theta_imp3 <-
  mc_data(rf3, N = 100)[,
    .(phi3 = mean(g_target(u) / f3(u))),
    keyby = sam]

rf4 <- function(n) {
  u <- runif(n)
  tan(pi * u / 4) # inverse transformation method
}
#-----
theta_imp4 <-
  mc_data(rf4, N = 100)[,
    .(phi4 = mean(g_target(u) / f4(u))),
    keyby = sam]

theta_imp <- theta_imp0[theta_imp1][theta_imp2][theta_imp3][theta_imp4]

theta_imp %>%
  melt(id.vars = "sam", variable.name = "imp_fun", value.name = "integral") %>%
  ggplot(aes(x = integral, fill = imp_fun)) +
  geom_histogram(bins = 30, position = "identity", alpha = .5) +
  xlab(expression(theta))

```

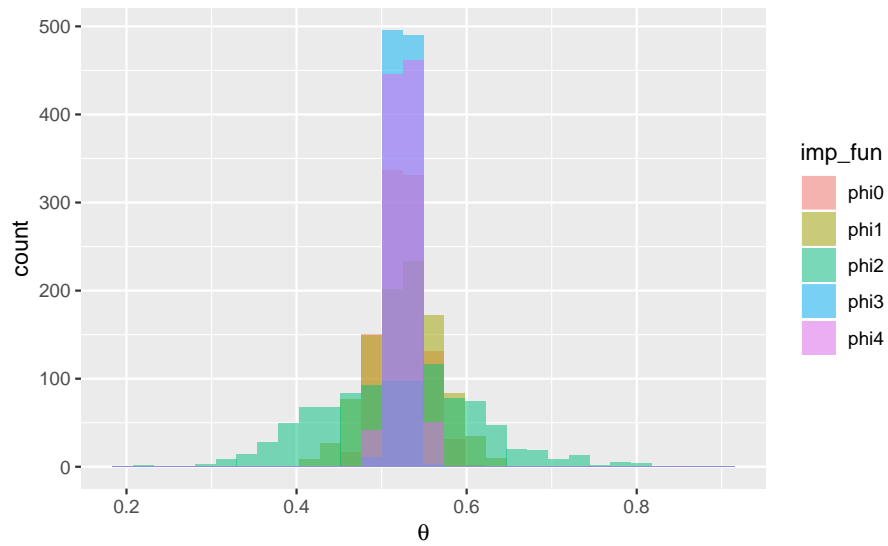


Figure 2.10: Empirical distribution of each importance sampling

```
theta_imp[,
  lapply(.SD, sd),
  , .SDcols = -"sam"]
#>      phi0  phi1  phi2  phi3  phi4
#> 1: 0.0246 0.042 0.0941 0.00977 0.0144
```

$f_3$  and possibly  $f_4$  yields the lowest variance. What happened to  $f_1$ ? Its support is  $(0, \infty)$ , so many values would be generated outside of  $(0, 1)$ . This results in many zeros in the sum of  $\frac{q}{f}$ .



## Chapter 3

# Monte Carlo Methods in Inference

### 3.1 Parametric Bootstrap

In this setting, we know distribution of  $X$ . We can freely generate from this distribution.

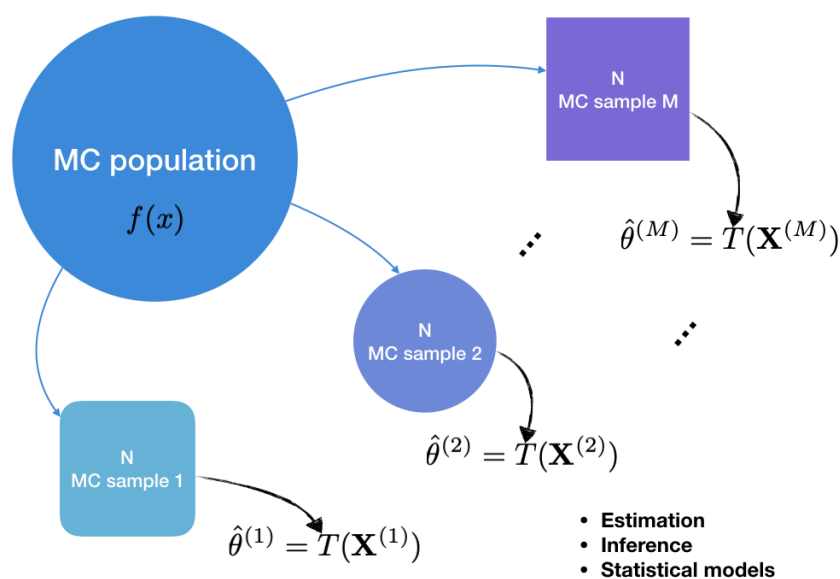


Figure 3.1: Parametric bootstrap

See Figure 3.1. From the “true” distribution, we can generate multiple samples. From each sample estimator can be computed. Then we can check these multiple estimates. Multiple estimates are close to motivation of estimator, so it helps exploring statistical inference with simple steps.

```
mc_data <- function(rand, N = 10000, M = 1000, char = "s", ...) {  
  data.table(  
    x = rand(n = N * M, ...),  
    sam = gl(M, N, labels = paste0(char, 1:M))  
  )  
}
```

## 3.2 Monte Carlo Methods for Estimation

**Example 3.1** (Any quantity of interest). Suppose that  $X_1, X_2 \stackrel{iid}{\sim} N(0, 1)$ . We want to estimate

$$\theta = E|X_1 - X_2|$$

### 3.2.1 Empirical distribution

**Algorithm 13:** Empirical distribution of  $\hat{\theta}$

```

input : distribution  $f$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   | Generate  $(X_1^{(m)}, X_2^{(m)}) \stackrel{iid}{\sim} N(0, 1)$ ;
3   | Compute  $\hat{\theta}^{(m)} = |X_1^{(m)} - X_2^{(m)}|$ ;
4 end
5 Draw a histogram;
output:  $\bar{\theta} = \frac{1}{M} \sum_{m=1}^M \hat{\theta}_m^{(m)}, \{\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(M)}\}$ 

```

```

basicmc <-
  mc_data(rnorm, N = 2)[,
    xname := gl(2, 1, length = 2000, labels = c("x1", "x2"))] %>%
  dcast(sam ~ xname, value.var = "x") %>%
  .[,
    .(that = mean(abs(x1 - x2))),
    by = sam]

```

```

basicmc[,
  .(est = mean(that))]
#>      est
#> 1: 1.1

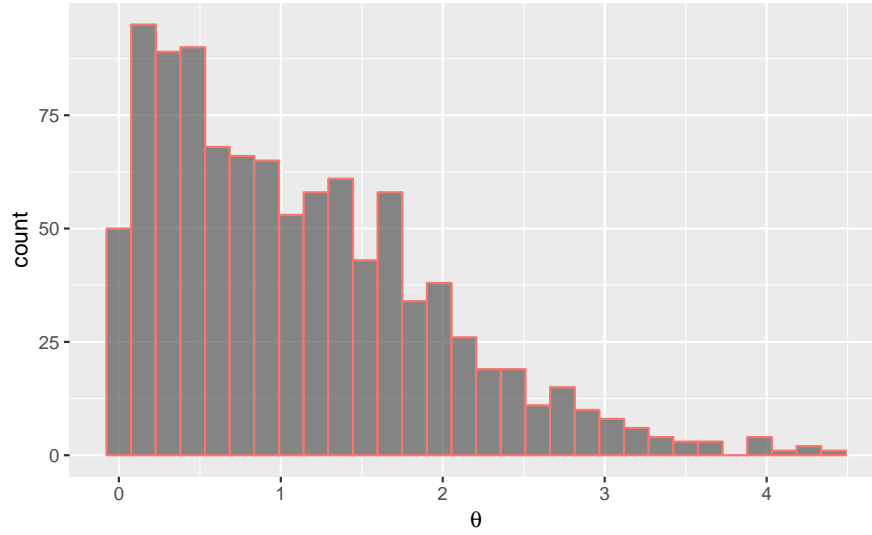
```

```

basicmc %>%
  ggplot(aes(x = that)) +
  geom_histogram(bins = 30, col = gg_hcl(1), alpha = .7) +
  xlab(expression(theta))

```



Figure 3.2: Empirical distribution of  $\hat{\theta}$  for  $|X_1 - X_2|$ 

### 3.2.2 Standard error

In Algorithm 13, we can get standard error by just calculating standard deviation of

$$\{\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(M)}\}$$

#### Algorithm 14: Standard error of $\hat{\theta}$

```

input : distribution  $f$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   Generate  $(X_1^{(m)}, X_2^{(m)}) \stackrel{iid}{\sim} N(0, 1)$ ;
3   Compute  $\hat{\theta}^{(m)} = |X_1^{(m)} - X_2^{(m)}|$ ;
4 end
5  $\bar{\hat{\theta}} = \frac{1}{M} \sum_{m=1}^M \hat{\theta}_m^{(m)}$ ;
6  $\widehat{SE}(\hat{\theta}) = \sqrt{\frac{1}{M-1} \sum_{m=1}^M (\hat{\theta}^{(m)} - \bar{\hat{\theta}})^2}$ ;
output:  $\widehat{SE}(\hat{\theta})$ 

```

```

basicmc[,
      .(se = sd(that))]
#>      se
#> 1: 0.844

```

### 3.2.3 Mean squared error

$MSE$  is used when comparing several estimators.

**Definition 3.1** (Mean squared error).

$$MSE(\hat{\theta}) := E(\hat{\theta} - \theta)^2$$

To know  $MSE$ , however, we should compute expectation. Some of them might be complicated even though we know true distribution. As the last chapter, we can apply Monte carlo method.

**Example 3.2** (MSE of a trimmed mean). Suppose that  $X_1, \dots, X_n \stackrel{iid}{\sim} N(2, 1)$ . Consider three estimators for  $\mu = 2$ .

1. mean  $\bar{X}$
2. median  $\tilde{X}$
3.  $k$ th trimmed mean  $\bar{X}_{[-k]}$

**Algorithm 15:** MSE of mean, median, and  $k$ th trimmed mean

```

input : distribution  $f$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   Generate  $(X_1^{(m)}, \dots, X_N^{(m)}) \stackrel{iid}{\sim} N(2, 1)$ ;
3   Sort  $(X_1^{(m)}, \dots, X_N^{(m)})$  in increasing order, i.e.  $(X_{(1)}^{(m)}, \dots, X_{(N)}^{(m)})$ ;
4   Mean  $\bar{X}^{(m)} = \frac{1}{N} \sum_{i=1}^N X_i^{(m)}$ ;
5   Median  $\tilde{X}^{(m)} = \begin{cases} X_{\frac{N}{2}+1}^{(m)} & N \text{ odd} \\ \frac{X_{\frac{N}{2}}^{(m)} + X_{\frac{N}{2}+1}^{(m)}}{2} & N \text{ even} \end{cases}$ ;
6    $k$ th trimmed mean  $\bar{X}_{[-k]}^{(m)} = \frac{1}{N-2k} \sum_{i=k+1}^{n-k} X_{(i)}^{(m)}$ 
7 end
8  $\widehat{MSE}(\bar{X}) = \frac{1}{M} \sum_{m=1}^M (\bar{X}^{(m)} - 2)^2$ ;
9  $\widehat{MSE}(\tilde{X}) = \frac{1}{M} \sum_{m=1}^M (\tilde{X}^{(m)} - 2)^2$ ;
10  $\widehat{MSE}(\bar{X}_{[-k]}) = \frac{1}{M} \sum_{m=1}^M (\bar{X}_{[-k]}^{(m)} - 2)^2$ ;
output:  $\widehat{MSE}(\bar{X})$ ,  $\widehat{MSE}(\tilde{X})$ , and  $\widehat{MSE}(\bar{X}_{[-k]})$ 

```

```

trim <- function(x, k = 1) {
  n <- length(x)
  x <- sort(x)
  sum(x[(k + 1):(n - k)]) / (n - 2 * k)
}
#-----
mu_list <- function(x, k) {
  list(mean = mean(x), median = median(x), trim = trim(x, k))
}

```

Try  $k = 1$ .

```

(trim_mc <-
  mc_data(rnorm, mean = 2, sd = 1)[,
    unlist(lapply(.SD, mu_list, k = 1)) %>% as.list,
    by = sam])
#>      sam x.mean x.median x.trim
#> 1:    s1  2.02    2.02    2.02
#> 2:    s2  2.00    2.00    2.00
#> 3:    s3  2.00    2.01    2.00
#> 4:    s4  1.99    1.98    1.99
#> 5:    s5  2.00    1.99    2.00
#> ---
#> 996: s996  2.02    2.02    2.02

```

```
#> 997: s997 2.00 1.99 2.00
#> 998: s998 1.99 1.99 1.99
#> 999: s999 1.99 1.99 1.99
#> 1000: s1000 2.00 2.01 2.00
```

```
trim_mc %>%
  melt(id.vars = "sam", variable.name = "hat") %>%
  ggplot(aes(x = value, fill = hat)) +
  geom_histogram(bins = 30, alpha = .3, position = "identity") +
  xlab(expression(mu)) +
  geom_vline(xintercept = 2, col = I("red")) +
  scale_fill_discrete(
    name = "Estimates",
    labels = c("Mean", "Median", "Trimmed")
  )
)
```

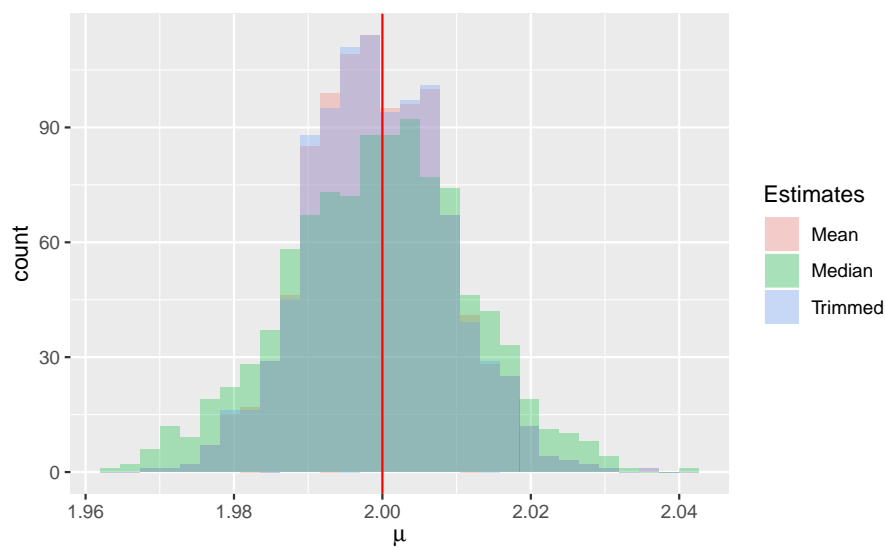


Figure 3.3: Empirical distribution of each estimator for  $\mu = 2$

Here, median shows the largest standard error.

```
trim_mc[,
  lapply(.SD, sd),
  .SDcols = c("sam")]
#>   x.mean x.median x.trim
#> 1: 0.00942 0.0123 0.00942
```

Now try various  $k$  for trimmed mean.

```
mse_list <- function(x, k) {
  list(mse = mean((x - 2)^2), se = sd(x))
}
#-----
trim_mse <-
  mc_data(rnorm, mean = 2, sd = 1)[,
    lapply(.SD, function(x) {
      sapply(0:9, function(k) {
        trim(x = x, k = k)
      })
    })
  ]
```

```

    })
  }) %>%
    unlist() %>%
    as.list(),
  by = sam][,
    lapply(.SD, mse_list) %>%
    unlist() %>%
    as.list(),
    .SDcols = -"sam"]

trim_mse %>%
  transpose() %>%
  .[,
    `:=`(
      k = rep(0:9, each = 2),
      hat = gl(2, k = 1, length = 2 * 10, labels = c("mse", "se"))
    )] %>%
  dcast(k ~ hat, value.var = "V1")
#>      k      mse      se
#>  1: 0 9.83e-05 0.00992
#>  2: 1 9.83e-05 0.00992
#>  3: 2 9.83e-05 0.00992
#>  4: 3 9.83e-05 0.00992
#>  5: 4 9.82e-05 0.00992
#>  6: 5 9.83e-05 0.00992
#>  7: 6 9.83e-05 0.00992
#>  8: 7 9.83e-05 0.00992
#>  9: 8 9.83e-05 0.00992
#> 10: 9 9.83e-05 0.00992

```

### 3.3 Confidence interval

Remember the meaning of 95% confidence interval. *If we have 100 samples and construct confidence interval in each sample, 95 intervals would include true parameter.* In this Monte Carlo setting, we know true population distribution, so we can generate multiple samples. Thus, we can reproduce this confidence interval situation.

#### 3.3.1 Empirical confidence interval

See one of histograms of Figure 3.3. Estimates are sorted. Calculating the upper and lower quantiles would give values close to confidence interval. See Figure 3.2. While the former show symmetric distribution, this is not. 0.25 and 0.975 quantile might be inappropriate. In this case, we should pick the *shortest interval*

with 95%. Best critical region leads to the shortest length of CI given  $\alpha$ , so we are finding this one.

**Algorithm 16:** Empirical confidence interval by Monte Carlo method

```

input : distribution  $f$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   | Generate  $X_1^{(m)}, \dots, X_n^{(m)} \stackrel{iid}{\sim} f$ ;
3   | Compute  $\hat{\theta}^{(m)} = \hat{\theta}(\mathbf{X}^{(m)})$ ;
4 end
5 if Distribution of  $\{\hat{\theta}^{(m)}\}_1^M$  symmetric then
6   | Sort  $\{\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(M)}\}$  in decreasing order, i.e.  $\{\hat{\theta}_{(1)}^{(1)}, \dots, \hat{\theta}_{(M)}^{(M)}\}$ ;
7   | Compute  $LB = \frac{\alpha}{2}$  sample quantile and  $UB = 1 - \frac{\alpha}{2}$  sample quantile;
8 else
9   | foreach  $lb < 0.05$  with  $ub - lb = 1 - \alpha$  do
10  |   | Candidate interval  $(lb, ub)$ ;
11  |   | calculate length  $l_i = ub - lb$ ;
12  | end
13  |  $(LB, UB)$ : pick up the interval with the smallest length  $l_i$ ;
14 end
output:  $(LB, UB)$ 

```

### 3.3.2 Empirical confidence level

On the contrary, we can estimate confidence level given confidence interval.

**Example 3.3** (Confidence interval for variance). If  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then

$$T = \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

Thus,  $100(1 - \alpha)\%$  confidence interval is given by

$$\left(0, \frac{(n-1)S^2}{\chi_{\alpha}^2(n-1)}\right)$$

For each MC sample, compute confidence interval. Just check if *known true parameter* is in the interval. Its proportion becomes the confidence level. It is simpler than estimate confidence interval itself.

**Algorithm 17:** Empirical confidence level by Monte Carlo method

```

input : distribution  $f$  with parameter  $\theta$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   | Generate  $X_1^{(m)}, \dots, X_n^{(m)} \stackrel{iid}{\sim} f$ ;
3   | Compute the confidence interval  $C_m$ ;
4   | Compute  $Y_j = I(\theta \in C_m)$ , i.e. whether  $\theta$  is in the CI;
5 end
6 Empirical confidence level  $\bar{Y} = \sum_{m=1}^M Y_m$ ;
output:  $\bar{Y}$ 

```

Let  $\mu = 0$ ,  $\sigma = 2$ ,  $N = 20$ , and let  $M = 1000$ .

```

ci_var <- function(x, variance, alpha) {
  n <- length(x)
  s2 <- var(x)

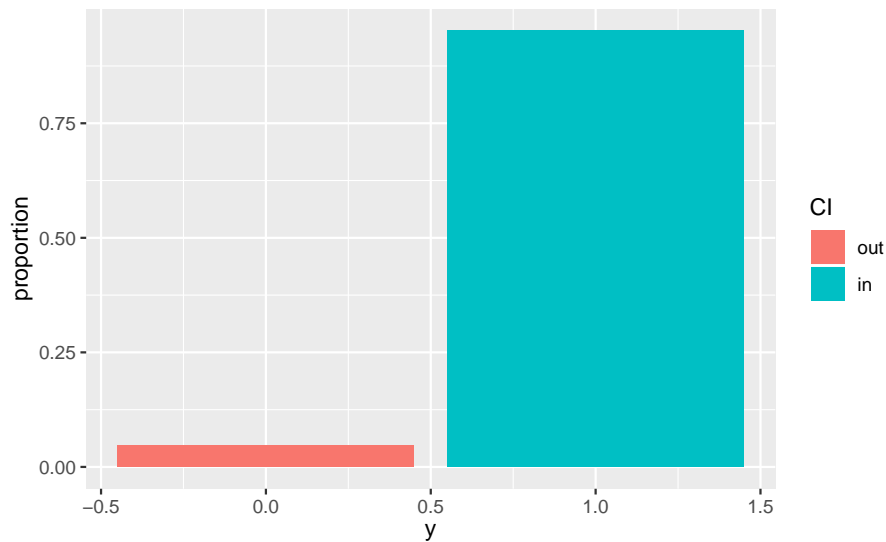
```

```

  (n - 1) * s2 / qchisq(alpha, df = n - 1) > variance
}
#-----
ci_lev <-
  mc_data(rnorm, N = 20, M = 1000, mean = 0, sd = 2)[,
    .(hat = mean(ci_var(x, variance = 4, alpha = .05)),
      by = sam]

ci_lev[,
  .N,
  by = hat][,
    proportion := N / sum(N)] %>%
  ggplot(aes(x = hat, y = proportion, fill = factor(hat))) +
  geom_bar(stat = "identity") +
  scale_fill_discrete(
    name = "CI",
    labels = c("out", "in")
  ) +
  xlab(expression(y))

```

Figure 3.4: Proportion of  $\sigma^2$  in confidence intervals

This leads to empirical confidence level, i.e. *sample proportion*. Just follow the last step 6 of Algorithm 17.

```

(ci_lev <-
  ci_lev[,
    .(level = mean(hat))])
#>   level
#> 1: 0.952

```

It is very close to 0.95. One of advantages of simulation study is we can assume various situation. For example, *violation of Gausiannity*.

**Example 3.4** (Violation of Normal distribution assumption). Refer to Example 3.3. This has assumed that  $X_i \stackrel{iid}{\sim} N(\mu = 2, \sigma^2 = 4)$ . What if not? For instance,

$$X_1, \dots, X_n \stackrel{iid}{\sim} \chi^2(df = 2)$$

Just change random numbers.

```
ci_lev2 <-
  mc_data(rchisq, N = 20, M = 1000, df = 2)[,
    .(hat = mean(ci_var(x, variance = 4, alpha = .05))),
    by = sam][,
    .(non_normal = mean(hat))]
```

Table 3.1: Empirical confidence level for each population

Normal	Chisq
0.952	0.763

From Table 3.1, we found that *non-normality lowers confidence level* from 0.952 to 0.763.

## 3.4 Hypothesis tests

Using MC method, we have done point estimation and interval estimation. Now consider *hypothesis testing*.

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_1 : \theta \in \Theta_1$$

where  $\{\Theta_0, \Theta_1\}$  is a partition of the parameter space  $\Theta$ . First of all, we have *test statistic*

$$T(\mathbf{X}) \stackrel{H_0}{\sim} f$$

and  $f$  is called *null distribution*. Given observed data, we compute this test statistic  $T_0$ . Where  $T_0$  is located in the null distribution  $f$  decides whether we reject or accept  $H_0$ . If  $T_0$  is very far from the middle, we can say that the realized data set is very rare event under  $H_0$ . In this case, we reject  $H_0$ . Otherwise, accept it. This is why we compute the tail probability, p-value.

### 3.4.1 Empirical p-value

**Example 3.5.** Suppose that  $X_1, \dots, X_{10} \stackrel{iid}{\sim} \text{Exp}(\lambda = 1)$ , which are observed as follows

‘*rxexp*’

Let  $\theta = E(X) = \frac{1}{\lambda}$ .

$$H_0 : \theta = 0.5 \quad \text{vs} \quad H_1 : \theta > 0.5$$

Test using  $T = \frac{\bar{X} - \theta_0}{S/\sqrt{n}}$  statistic.

Before looking at p-value, briefly look at *empirical null distribution* of test statistic.

```
mc_data(rexp, rate = 2)[,
  .(tstat = t.test(x, mu = .5)$statistic),
  by = sam] %>%
  ggplot(aes(x = tstat)) +
  geom_histogram(bins = 30, col = gg_hcl(1), alpha = .7) +
  geom_vline(xintercept = t.test(xexp, mu = .5)$statistic, col = I("red")) + # xexp: observed data
  geom_vline(xintercept = -t.test(xexp, mu = .5)$statistic, col = I("red")) +
  xlab("T")
```

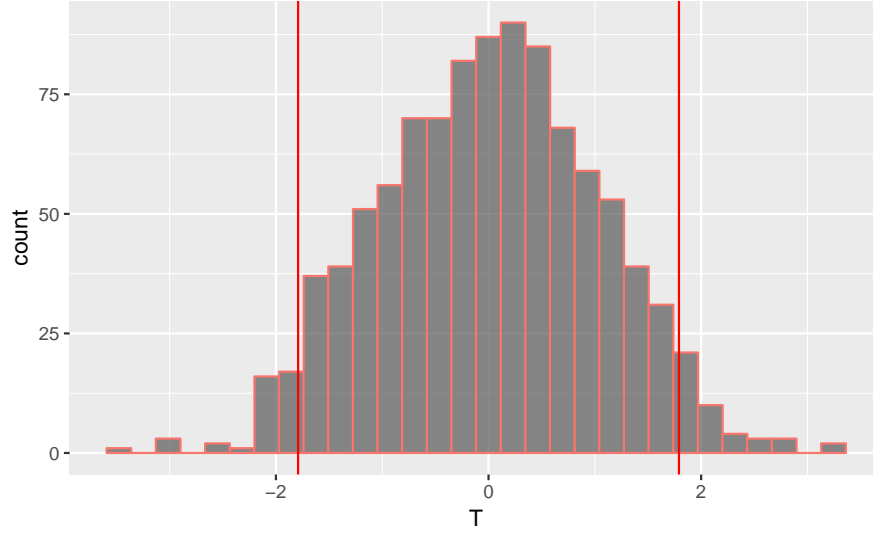


Figure 3.5: Empirical Null Distribution

By proceeding the similar way, we can get empirical distribution of test statistics. Some are out of observed  $T_0$ , some are not. Motivation is that we just count these. Proportion of these would estimate p-value. Recap what p-value is.

**Definition 3.2** (p-value). Let  $T$  be test statistic and let  $T_0$  be observed test statistic given data. Then p-value is

$$p - value := \begin{cases} P(|T| \geq T_0 \mid H_0) & \text{both sided} \\ P(T \geq T_0 \mid H_0) & \text{one sided} \\ P(T \leq T_0 \mid H_0) & \text{one sided} \end{cases}$$

Denote that p-value is probability. So in MC setting, we can estimate this by computing *sample mean of identity function*.

**Lemma 3.1** (Empirical p-value). Let  $T_0$  be observed test statistic and let  $\{T_1, \dots, T_M\}$  be test statistic computed in each MC sample.

$$\text{Empirical } p\text{-value} = \begin{cases} \frac{|\{T_j : (T_j > |T_0|) \text{ or } (T_j < -|T_0|)\}|}{M} & \text{both-sided} \\ \frac{|\{T_j : (T_j > T_0)\}|}{M} \text{ or } \frac{|\{T_j : (T_j < T_0)\}|}{M} & \text{one-sided} \end{cases}$$

**Algorithm 18:** Empirical p-value by Monte Carlo method

<b>input</b> : Given observed data, compute $T_0$	
1	<b>for</b> $m \leftarrow 1$ <b>to</b> $M$ <b>do</b>
2	Generate $X_1^{(m)}, \dots, X_n^{(m)} \stackrel{H_0}{\sim} f$ ;
3	Compute $T_m(\mathbf{X}^{(m)})$ ;
4	<b>end</b>
5	Empirical p-value $\hat{p} = \begin{cases} \frac{ \{T_j : (T_j >  T_0 ) \text{ or } (T_j < - T_0 )\} }{M} & \text{both-sided;} \\ \frac{ \{T_j : (T_j > T_0)\} }{M} \text{ or } \frac{ \{T_j : (T_j < T_0)\} }{M} & \text{one-sided} \end{cases}$
<b>output:</b> $\hat{p}$	



Go back to Example 3.5. Only left is computing 5 of Algorithm 18. (Denote that `xexp` in the code is vector object of observed data).

```
(tt_exp <-
  mc_data(rexp, rate = 2)[,
    .(tstat = t.test(x, mu = .5)$statistic),
    by = sam)[,
    .(pval = mean(tstat > abs(t.test(xexp, mu = .5)$statistic)))])
#>      pval
#> 1: 0.046
```

It is smaller than 0.05, so we reject  $H_0$ .

### 3.4.2 Comparing several tests

MC method would be used in comparing tests rather than conducting test itself. By generating random number, we can evaluate tests.

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_1 : \theta \in \Theta_1$$

As mentioned earlier,  $\{\Theta_0, \Theta_1\}$  is a partition of the parameter space  $\Theta$ . For this test, we can perform several tests. Test method 1, test method 2, et cetera. All these methods produce error, but these errors might be different. So we try to compare this.

what is true	accept $H_0$	reject $H_0$
$H_0$	correct decision	Type I error
$H_1$	Type II Error	correct decision

In most tests, we aim to reject  $H_0$ . By rejecting it, we can evidently say that  $H_0$  is not true. In this sense, we treat type I error more importantly than type II error in general. Test strategy becomes to control type I error probability first and then lower type II error probability.

**Definition 3.3** (Power function). Let  $\theta \in \Theta$  be a parameter of a test.

$$\beta(\theta) := P(\text{reject } H_0 \mid \theta)$$

With this power function, each type I error and type II error probability is given.

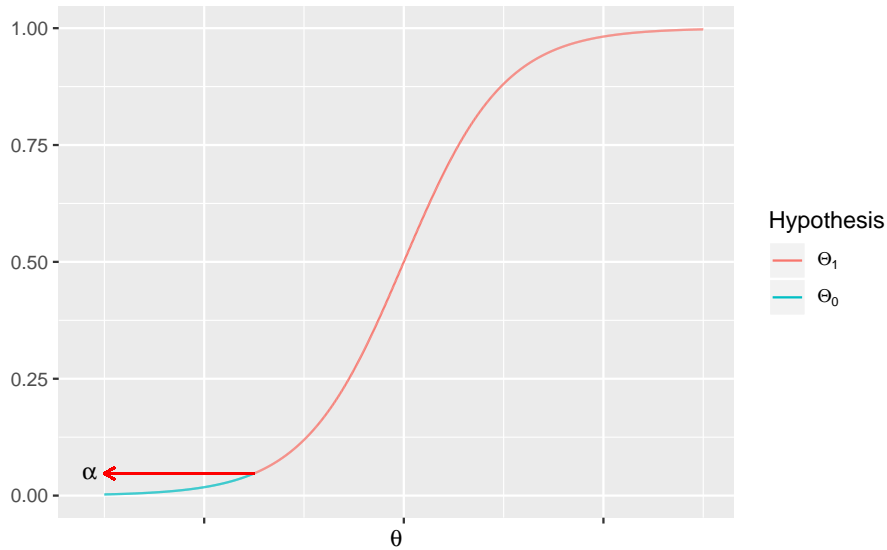
**Lemma 3.2** (typeerr). 1.  $P(\text{Type I error}) = \beta(\theta_0)$ ,  $\theta_0 \in \Theta_0$

2. Power  $\beta(\theta_1) = 1 - P(\text{Type II error})$ ,  $\theta_1 \in \Theta_1$

Following our test strategy, fixing  $P(\text{Type I error})$  and maximizing  $\beta(\theta_1)$ , we construct following test.

**Definition 3.4** (Size  $\alpha$  Test). A test with  $\beta(\theta)$  is called size  $\alpha$  test if and only if

$$\alpha := \sup_{\theta \in \Theta_0} \beta(\theta), \quad 0 \leq \alpha \leq 1$$

Figure 3.6: Size  $\alpha$  Test

Then how to compare tests? Look at the following example. Three columns of the middle part are type I error rate.

test methods	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.01$	Power
Test 1	0.09	0.04	0.01	0.7
Test 2	0.11	0.06	0.01	0.65
Test 3	0.15	0.07	0.02	0.9

Here, we will choose **Test 1**.

1. Type I error rate  $\approx \alpha$ 
  - before looking at power, this should be satisfied.
  - So Test 3 is excluded
2. Larger power
  - Thus, we select Test 1.

### 3.4.3 Empirical type-I error rate

Recall Lemma 3.2. As in p-value, we just compute sample proportion for each type I error rate and power under null and alternative distribution.

**Lemma 3.3.** Consider  $H_0 : \theta \in \Theta_0$  vs  $H_1 : \theta \in \Theta_1$ .

Define  $I(\mathbf{X})$  by

$$I(\mathbf{X}) = \begin{cases} 1 & H_0 \text{ is rejected} \mid H_0 \\ 0 & \text{otherwise} \end{cases}$$

For each MC sample, compute this statistic  $I_m = I(\mathbf{X}^m)$ . Then empirical type I error rate can be computed as

$$\frac{1}{M} \sum_{m=1}^M I_m$$

Table 3.4: Empirical Type I error for Z and T

Z-test	T-test
0.048	0.051

**Algorithm 19:** Empirical type I error rate by Monte Carlo method

```

input :  $H_0 : \theta \in \Theta_0$  vs  $H_1 : \theta \in \Theta_1$ 
1 for  $m \leftarrow 1$  to  $M$  do
2   Generate  $X_1^{(m)}, \dots, X_n^{(m)} \stackrel{H_0}{\sim} f$ ;
3   Compute  $T_m(\mathbf{X}^{(m)})$ ;
4   Compute  $I_m = \begin{cases} 1 & H_0 \text{ is rejected} \mid H_0 \\ 0 & \text{otherwise} \end{cases}$ ;
5 end
6 Empirical Type I error rate  $\hat{\alpha} = \frac{1}{M} \sum_{m=1}^M I_m$ ;
output: compare  $\hat{\alpha}$  with  $\alpha$ 

```

**Example 3.6** (Testing normal mean). Suppose that  $X_1, \dots, X_{20} \stackrel{iid}{\sim} N(\mu, \sigma^2 = 100)$ . Test

$$H_0 : \mu = 500 \quad \text{vs} \quad H_1 : \mu > 500$$

$$1. \text{ Z-test: } Z = \frac{\bar{X} - 500}{\sigma^2 / \sqrt{20}} \stackrel{H_0}{\sim} N(0, 1)$$

$$2. \text{ t-test: } T = \frac{\bar{X} - 500}{S / \sqrt{20}} \stackrel{H_0}{\sim} t(20 - 1)$$

```

test_list <- function(x, mu, sig, a = .05) {
  n <- length(x)
  xbar <- mean(x) - mu
  list(
    z = xbar / (sig / sqrt(n)) > qnorm(a, lower.tail = FALSE),
    t = xbar / (sd(x) / sqrt(n)) > qt(a, df = n - 1, lower.tail = FALSE)
  )
}
#-----
err_mc <-
  mc_data(rnorm, N = 20, mean = 500, sd = 10)[,
    lapply(.SD, test_list, mu = 500, sig = 10) %>%
      unlist() %>%
      as.list(),
    by = sam][,
      lapply(.SD, mean),
      .SDcols = -"sam"]

```

Both test have Type I error close to  $\alpha$ , but Z-test seems bit better.

### 3.4.4 Empirical power

Next step is power. See Figure 3.6. Power is different in that this is computed in *alternative distribution*, *not null distribution*.

$$\beta(\theta_1) = P(\text{reject } H_0 \mid \theta_1 \in \Theta_1)$$

**Lemma 3.4.** Consider  $H_0 : \theta \in \Theta_0$  vs  $H_1 : \theta \in \Theta_1$ .

Define  $I(\mathbf{X})$  by

$$I(\mathbf{X}) = \begin{cases} 1 & H_0 \text{ is rejected} \mid H_1 \\ 0 & \text{otherwise} \end{cases}$$

For each MC sample, compute this statistic  $I_m = I(\mathbf{X}^m)$ . Then empirical power can be computed as

$$\frac{1}{M} \sum_{m=1}^M I_m$$

Process will be same but we test under  $H_1$ . However, this makes a lot difference due to structure of each hypothesis. In many cases,  $H_0$  is simple, i.e.  $\mu = 500$ . In 2 of Algorithm 19, we can consider only  $N(\mu = 500, 100)$ . Since  $\Theta_0$  and  $\Theta_1$  form partition, alternative hypothesis usually is not simple. In this example,  $\mu > 500$ . We cannot specify one distribution for alternative. How to deal with this?

Trying many points for  $\mu_1 \in \Theta_1 = \{\mu : \mu > 500\}$  might be possible. Our goal is finding larger power. So find test with larger power for all points in  $\Theta_1$ .

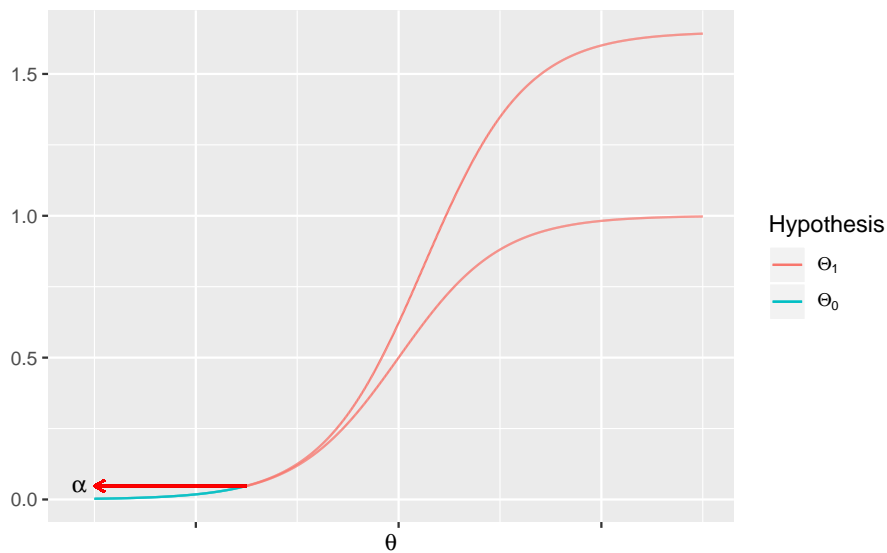


Figure 3.7: Comparing power between two test methods

See Figure 3.7. One test method has higher  $\beta(\theta)$  function curve in  $\Theta_1$ . This test is *powerful than the other*. We would choose this test in this step. So what we have to do is choose some points  $\theta_1 \in \Theta_1$ , and draw the

power curve.

**Algorithm 20:** Empirical power by Monte Carlo method

```

input :  $H_0 : \theta \in \Theta_0$  vs  $H_1 : \theta \in \Theta_1$ 
1 foreach  $\theta_1 \in \Theta_1$  do
2   for  $m \leftarrow 1$  to  $M$  do
3     Generate  $X_1^{(m)}, \dots, X_n^{(m)} \stackrel{H_0}{\sim} f$ ;
4     Compute  $T_m(\mathbf{X}^{(m)})$ ;
5     Compute  $I_m = \begin{cases} 1 & H_0 \text{ is rejected} \mid H_1; \\ 0 & \text{otherwise} \end{cases}$ ;
6   end
7   Empirical power  $\hat{\beta} = \frac{1}{M} \sum_{m=1}^M I_m$ ;
8 end
9 Draw a power curve  $\hat{\beta}$  against  $\theta_1$  output: curve and  $\{\hat{\beta}\}$ 

```

In fact, we can try every  $\theta \in \Theta$  and *draw entire power curve*. Refer to Example 3.6.

```

pw_mc <-
  lapply(seq(450, 650, by = 10), function(mu) {
    mc_data(rnorm, N = 20, mean = mu, sd = 10)[,
                                                h1 := mu]
  })
pw_mc <- rbindlist(pw_mc)

```

One column is added from previous process. This is group for  $H_1$ . So we should specify `by = .(h1, sam)`.

```

pw_mc <-
  pw_mc[,
    lapply(.SD, test_list, mu = 500, sig = 10) %>%
      unlist() %>%
      as.list(),
    by = .(h1, sam)][,
      lapply(.SD, mean),
      by = h1, .SDcols = -"sam"]

```

```

pw_mc %>%
  melt(id.vars = "h1", variable.name = "test") %>%
  ggplot(aes(x = h1, y = value, colour = test)) +
  geom_path() +
  geom_point() +
  scale_colour_discrete(
    name = "Test",
    labels = c("Z", "T")
  ) +
  labs(
    x = expression(mu),
    y = expression(beta)
  )

```

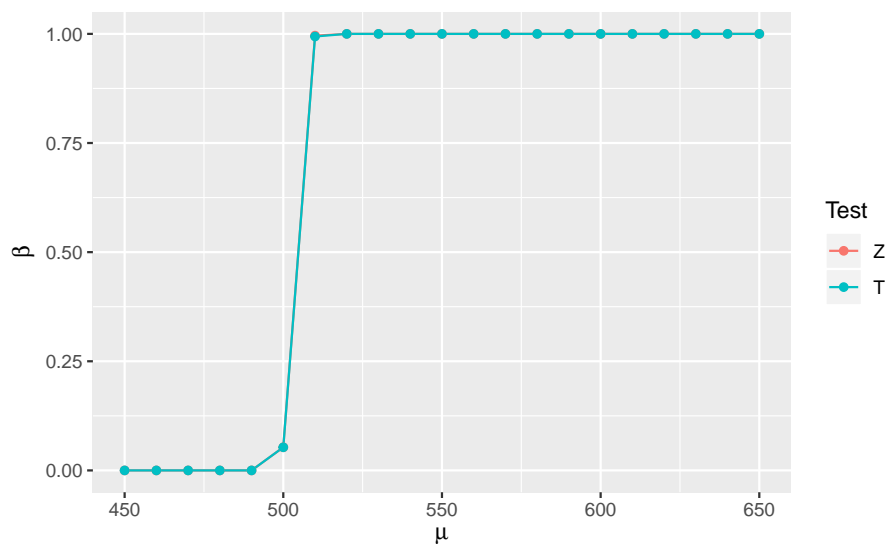


Figure 3.8: Empirical power curve of each z-test and t-test

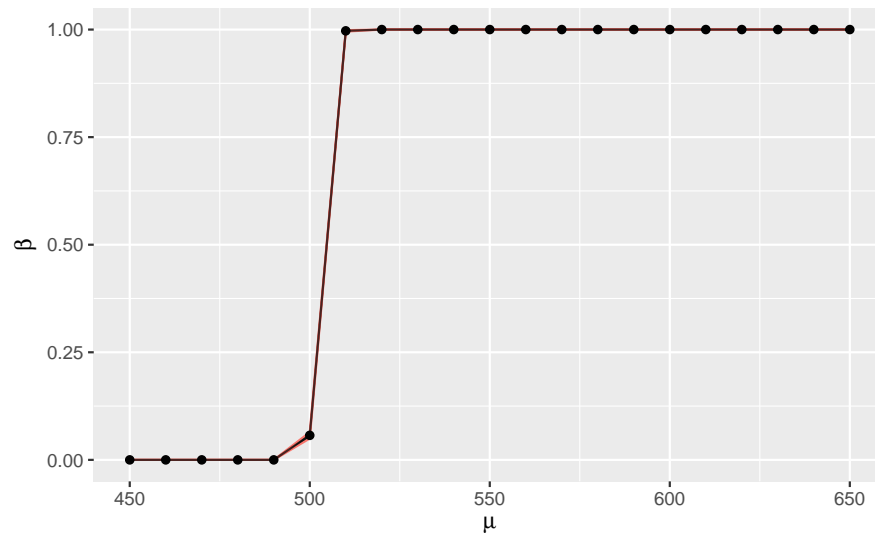
Recall that we are estimating power. Instead of `mean()`, we can use `sd()`. This would give us *standard error* of our estimator for power. Since it is sample proportion,

$$\widehat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{M}}$$

Consider *T*-test.

```
pw_mc2 <- lapply(seq(450, 650, by = 10), function(mu) {
  mc_data(rnorm, N = 20, mean = mu, sd = 10)[,
    h1 := mu]
})
pw_mc2 <- rbindlist(pw_mc2)
#-----
pw_mc2 <-
  pw_mc2[,
    .(te = t.test(x, alternative = "greater", mu = 500)$p.value <= .05),
    by = .(h1, sam)][,
    .(te = mean(te)),
    by = h1][,
    se := sqrt(te * (1 - te) / 1000)]

pw_mc2 %>%
  ggplot(aes(x = h1, y = te)) +
  geom_ribbon(aes(ymin = te - se, ymax = te + se), col = gg_hcl(1)) +
  geom_path(alpha = .7) +
  geom_point() +
  labs(
    x = expression(mu),
    y = expression(beta)
  )
```

Figure 3.9: Empirical power curve  $\hat{p} \pm \widehat{SE}(\hat{p})$  for t-test

### 3.4.5 Count Five test for equal variance

Commonly, F-test is used for equality of two population variances. McGrath and Yeh (2005) suggests nonparametric testing without Normal assumption, so called *Count Five*. Instead, this method requires some conditions.

1. same mean
2. same sample size

#### Algorithm 21: Count Five test

```

input :  $X_1, \dots, X_{n_x} \perp\!\!\!\perp Y_1, \dots, Y_{n_y}$ 
          $H_0 : \sigma_X^2 = \sigma_Y^2$ 
1 Compute  $C_X = \left| \{i : |X_i - \bar{X}| > \max_j |Y_j - \bar{Y}|\} \right|$ ;
2 if  $C_X \geq 5$  then
3   | return reject  $H_0$ ;
4 else
5   | return accept  $H_0$ ;
6 end

```

```

gauss <-
  tibble(
    x1 = rnorm(20, mean = 0, sd = 1),
    x2 = rnorm(20, mean = 0, sd = 1.5)
  )

```

```

gauss %>%
  gather(key = "variable", value = "value") %>%
  ggplot(aes(x = variable, y = value, fill = variable)) +
  geom_boxplot() +
  geom_point(alpha = .5)

```

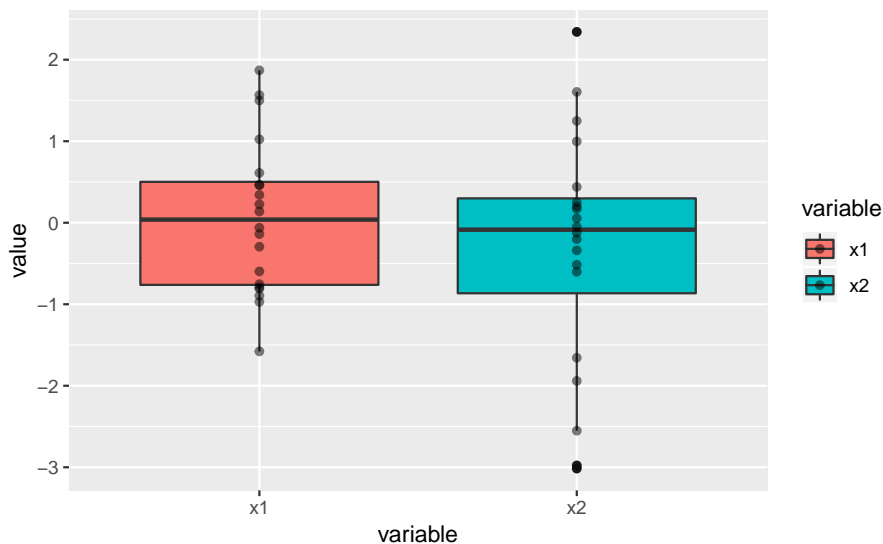


Figure 3.10: Side-by-side boxplot

We would perform *Count Five* test for multiple simulated data sets such as in Figure 3.10.

$$X_1^{(m)}, \dots, X_{20}^{(m)} \sim N(0, 1) \perp\!\!\!\perp Y_1^{(m)}, \dots, Y_{20}^{(m)} \sim N(0, 1.5)$$

```
count5test <- function(x, y) {
  X <- x - mean(x)
  Y <- y - mean(y)
  outx <- sum(X > max(Y)) + sum(X < min(Y))
  outy <- sum(Y > max(X)) + sum(Y < min(X))
  max(c(outx, outy)) > 5
}
```

Apply MC method to get *empirical type I error*.

```
mc_data(rnorm, N = 20, M = 1000)[,
  x2 := rnorm(20 * 1000)][,
  .(chat = count5test(x = x, y = x2)),
  by = sam][,
  .(chat = mean(chat))]

#>      chat
#> 1: 0.026
```

## 3.5 Statistical Methods

## 3.6 Bootstrap

### 3.6.1 Resampling

Bootstrap is a class of nonparametric Monte Carlo methods that estimate the distribution of a population by *resampling*. Different with previous MC method, we do not know population distribution. Instead, resampling methods treat an observed sample as a finite population. This is called *pseudo-population* in that this is regarded as having the same characteristics as the true population.



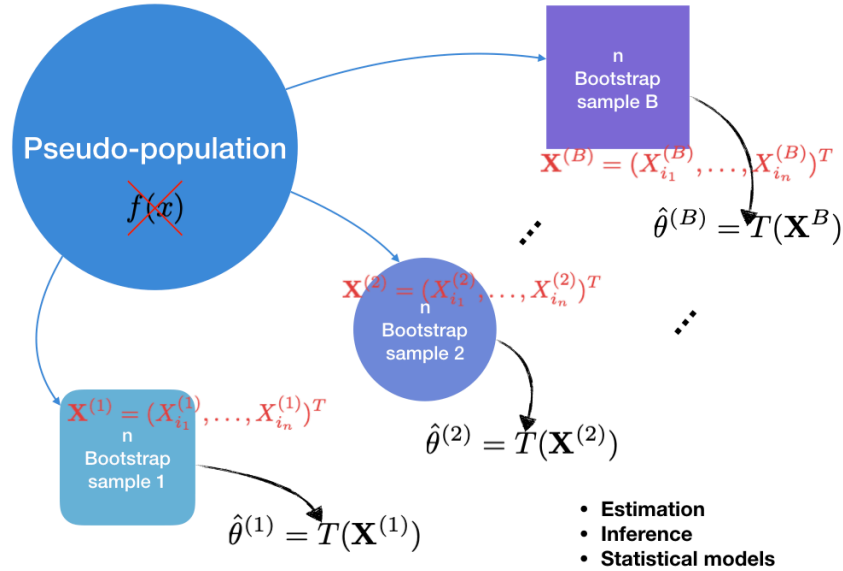


Figure 3.11: Resampling

See Figure 3.11. From the observed sample, which is pseudo-population, resampling generates multiple bootstrap samples by *sampling with replacement*. Surprisingly, this simple sampling procedure approximate the true population distribution quite successful.

### 3.6.2 Approximations in bootstrap

How does bootstrap work? Efron and Gong (1983) provides simple example  $T = \bar{x}$ , i.e. sample average.

**Example 3.7** (Estimation of sample mean). Having observed  $X_1 = x_1, \dots, X_n = x_n$ , compute

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Using bootstrap, we try to see the empirical distribution of this.

Note that

$$X_1, \dots, X_n \stackrel{iid}{\sim} F$$

where  $F$  is true unknown distribution. Having observed  $X_1 = x_1, \dots, X_n = x_n$ , we get *empirical distribution function* by computing the sample average.

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x) \quad (3.1)$$

This works for estimator of  $F$ .  $\hat{F}$  endows mass  $\frac{1}{n}$  on each observed point  $x_i$ ,  $i = 1, \dots, n$ . In other words, pseudo-distribution becomes *discrete uniform*.

$$\hat{F} \stackrel{d}{=} \text{unif}(x_1, \dots, x_n) \quad (3.2)$$

We have set the population distribution which is bogus. Now we can apply previous MC sampling with  $\frac{1}{n}$ . Write  $\{X_1^*, \dots, X_n^*\}$  as *bootstrap sample* by resampling. Then

$$P(X^* = x_i) = \frac{1}{n}$$

i.e.

$$X^* \stackrel{iid}{\sim} \text{unif}(x_1, \dots, x_n) \quad (3.3)$$

This gives *bootstrap cdf*, cdf of bootstrap samples, by

$$F^*(x) = \text{cdf of } \text{unif}(x_1, \dots, x_n) = \hat{F} \quad (3.4)$$

One proceeds in a similar way for estimating cdf that

$$\hat{F}^* = \frac{1}{n} \sum_{i=1}^n I(X^* \leq x_i) \quad (3.5)$$

This is called *ecdf of bootstrap replicates*.

*Remark.* For any data points  $X_1, \dots, X_n$ ,

1. **Empirical Distribution Function**  $\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x_i)$
2. **Bootstrap cdf**  $F^*(x) = \hat{F}$
3. **ECDF of bootstrap replicates**  $\hat{F}^* = \frac{1}{n} \sum_{i=1}^n I(X^* \leq x_i)$

This remark can explain Figure 3.11 in a distribution way.

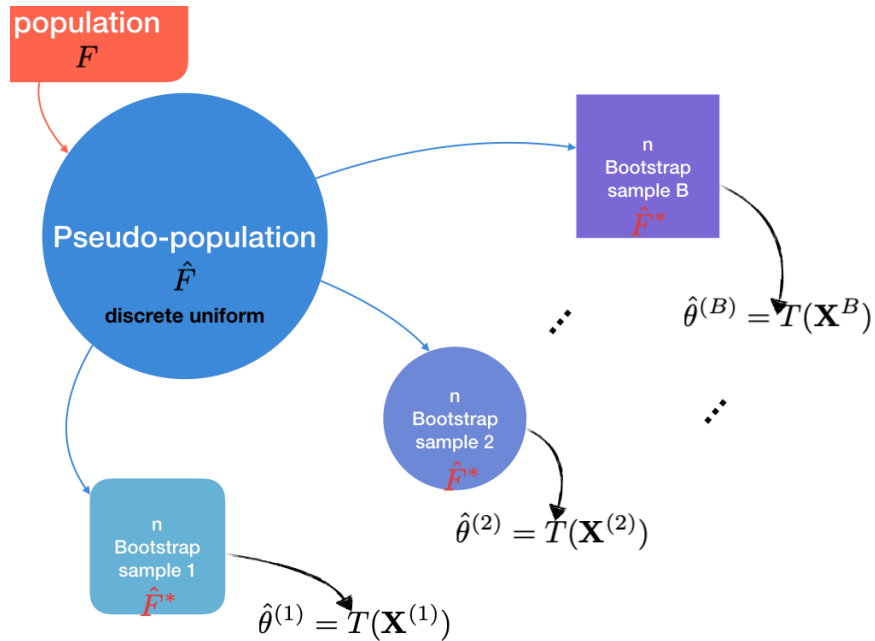


Figure 3.12: Empirical distribution of bootstrap

Figure 3.12 shows how each sample is distributed, approximately. We get the data set from true  $F$ . From this finite population, we first estimate  $F$  by  $\hat{F}$ . Resampling multiple bootstrap samples, each sample estimates  $\hat{F}$  by  $\hat{F}^*$ .

**Theorem 3.1** (Two approximations in bootstrap). *There are two approximations in bootstrap. For large  $B$ , bootstrap samples approximate bootstrap replicates. For large  $n$ , bootstrap replicates approximate true population.*

$$\begin{array}{c} \mathbf{X}^{(b)} \rightarrow \hat{f} \rightarrow f \\ \uparrow \quad \uparrow \\ \text{large } B \quad \text{large } n \end{array}$$

*Proof.* Since

$$E\hat{F} = \frac{1}{n} \sum_{i=1}^n P(X_i \leq x_i)$$

$$\hat{F} \xrightarrow{a.s.} F$$

as  $n \rightarrow \infty$  by the strong law of large numbers. Let

$$\overline{\hat{F}^*} := \frac{1}{B} \sum_b \hat{F}_b^*$$

where  $\hat{F}_b^*$  is ecdf of  $b$ -th bootstrap replicate. Since

$$E\hat{F}_b^* = \frac{1}{n} \sum_{i=1}^n P(X^{(b)} \leq x_i)$$

$$\overline{\hat{F}^*} \xrightarrow{a.s.} F^* = \hat{F}$$

as  $B \rightarrow \infty$  by S.L.L.N. □

Denote that Theorem 3.1 can be also expressed as

$$\begin{array}{c} \overline{\hat{F}^*} \rightarrow \hat{F} \rightarrow F \\ \uparrow \quad \uparrow \\ \text{large } B \quad \text{large } n \end{array} \tag{3.6}$$

This approximation not only justifies the bootstrap procedure but also shows the problem of it. We can always increase  $B$  if we want. Then we earn  $\overline{\hat{F}^*} \approx \hat{F}$ , i.e. bootstrap samples approximate pseudo-population. However,  $n$  is fixed. For  $\hat{F}$  to approximate  $F$ , large  $n$  is required. It is not under control. If small  $n$  is given,  $\hat{F}$  will not be close to  $F$ . Then the bootstrap samples will not be close to  $F$  finally.

**Corollary 3.1.** *Resampling the large number of replicates, i.e. large  $B$  produces a good estimates of  $\hat{F}$  but it does not guarantee a good estimate of  $F$ .*

**Example 3.8** (Poisson population). Suppose that  $\{2, 2, 1, 1, 5, 4, 4, 3, 2, 1\} \sim \text{Poisson}(\lambda = 2)$ . Resampling from this pseudo-population, can we appropriately explain the population?

Table 3.5: Empirical distribution

1	2	3	4	5
0.3	0.3	0.1	0.2	0.1

Table 3.5 is just a result of averaging, empirical distribution  $\hat{F}$ . General bootstrap will resample by distribution. We can see the problem at once.

$$P(X = 0) = e^{-2} = 0.135$$

However, we did not observe 0, so the bootstrap sample will never include zero, i.e. not a full domain. In sum,

$$\overline{\hat{F}^*} \rightarrow \hat{F} \not\rightarrow F$$

### 3.6.3 Bootstrap standard error

Look at the Figure 3.12 again. For each sample, we calculate estimator of our interest. For instance, Example 3.7 - sample average  $\bar{x}$ . After that, we get empirical distribution of  $\bar{x}$  such as *standard error*.

Before looking at the empirical distribution, let's review sample estimation.

**Theorem 3.2** (Standard error of sample mean). *Standard error of  $\bar{X} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ , i.e. the root mean squared error is estimated by*

$$\hat{\sigma} = \left[ \frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{\frac{1}{2}}$$

Here,  $n-1$  was divided for *unbiasedness*. Recall that

$$Var(\bar{X}) = \frac{\sigma^2}{n}$$

It follows that

$$\hat{\sigma}^2 = \frac{s^2}{n}$$

and we know that  $s^2$  should be divided by  $n-1$  to be unbiased. Bootstrap generalizes point estimation process 3.2 a bit differently (Efron and Gong, 1983). From Equation (3.1) to (3.4), we resample bootstrap sample with empirical probability distribution  $\hat{F}$  so that

$$X_1^*, \dots, X_n^* \stackrel{iid}{\sim} \hat{F} \tag{3.7}$$

where  $\hat{F}$  is discrete uniform in each observed data point. In this sample, compute the estimate of interest, e.g. average

$$\bar{X}^* = \frac{1}{n} \sum_{i=1}^n X_i^*$$

From Theorem 3.2, this  $\bar{X}^*$  has *estimated variance* of

$$Var \bar{X}^* = \frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2 \quad (3.8)$$

In fact, this is a variance *under sampling scheme* (3.7), i.e. indicates one-time-sampling from  $\hat{F}$ . Using this, we can construct **bootstrap estimate of standard error for sample mean**, which come by sampling  $B$  times.

$$\hat{\sigma}_B = \left[ \frac{1}{B-1} \sum_{b=1}^B (\bar{X}_b^* - \bar{X}^*) \right]^{\frac{1}{2}} \quad (3.9)$$

where  $\bar{X}_b^*$  is a sample mean of  $b$ -th bootstrap replicate and  $\bar{X}^*$  is average of every  $\bar{X}_b^*$ .

**Theorem 3.3** (Bootstrap standard error). *Bootstrap estimate of standard error for any estimator  $\hat{\theta}(X_1, \dots, X_n)$  is*

$$\hat{\sigma}_B = \left[ \frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_b^* - \bar{\theta}^*) \right]^{\frac{1}{2}}$$

where  $\hat{\theta}_b^*$  is independent bootstrap replications and

$$\bar{\theta}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^*$$

Now we format Figure 3.12 to practical algorithm.

**Algorithm 22:** Bootstrap algorithm

```

input :  $n$  observations  $x_1, \dots, x_n$ 
1 for  $b \leftarrow 1$  to  $B$  do
2   Sampling with replacement  $X_1^{(b)}, \dots, X_n^{(b)}$  from the observed sample;
3   Compute estimate
      
$$\hat{\theta}(X_1^{(b)}, \dots, X_n^{(b)}) \equiv \hat{\theta}_b^*$$

   ;
4 end
5  $\bar{\theta}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^*$ ;
6 Bootstrap standard error
      
$$\hat{\sigma}_B = \left[ \frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_b^* - \bar{\theta}^*) \right]^{\frac{1}{2}}$$

   ;
output:  $\hat{\sigma}_B$ 
```

Efron and Tibshirani (1994) suggests that  $B = 40$  is usually enough to estimate standard error well. It rarely require  $B > 200$ . On the other hand, much larger  $B$  is needed in interval estimation.

First we try perform bootstrap without package doing bootstrap such as **boot** or **bootstrap**. The following is observed sample  $n = 50$ , i.e. finite population in bootstrap literature.

```
X
#> # A tibble: 50 x 1
#>       x
#>   <dbl>
#> 1 35.2
#> 2  3.24
#> 3 11.5
#> 4  6.00
#> 5 48.9
#> 6  2.71
#> 7 38.7
#> 8 34.0
#> 9 19.8
#> 10 40.8
#> # ... with 40 more rows
```

Sample mean of this sample is

```
X %>%
  summarise(x = mean(x))
#> # A tibble: 1 x 1
#>       x
#>   <dbl>
#> 1 22.2
```

```
MC_CORES <- parallel::detectCores() - 1 # parallelization
```

```
resample <- function(data, statistic = mean) {
  # sampling with replacement
  xb <-
    data %>%
    data.table() %>%
    .[sample(1:.N, size = .N, replace = TRUE)]
  # estimator
  xb[,
    lapply(.SD, statistic) %>%
    as.numeric()
  ]
}
```

To fasten the process, we implement parallelization `parallel::mclapply`. This function is based on fork mechanism of Unix OS. So this is not available in Windows OS.

```
B <- 40
#-----
Xse <-
  parallel::mclapply(
    1:B,
    function(x) {resample(data = X, statistic = mean)},
    mc.cores = MC_CORES
  ) %>%
  unlist()
```

```
tibble(se = Xse) %>%
  ggplot(aes(x = se)) +
  geom_histogram(bins = 30) +
  xlab(expression(hat(sigma)[B]))
```

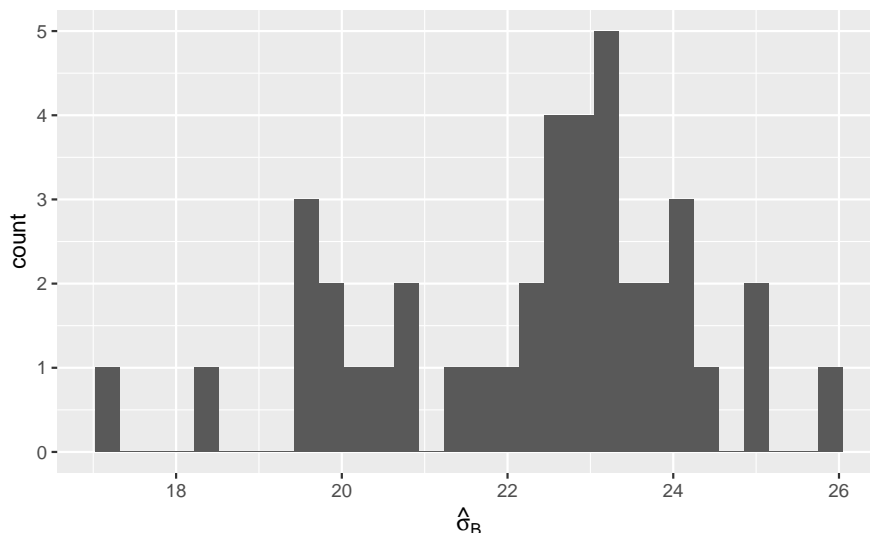


Figure 3.13: Bootstrap replicates of sample mean done by mclapply

Bootstrap standard error is given by

```
sd(Xse)
#> [1] 1.91
```

Another way and possible also in Windows OS is `foreach::foreach`. Using `%dopar%` with pre-specified workers instead of `%do%`, we can parallelize the jobs. `.inorder` argument enable the task done in order different with the other functions. Default is `FALSE` and it is more stable.

When we do `foreach` parallization, we should distribute the jobs to workers manually.

```
cl <- parallel::makeCluster(MC_CORES)
doParallel::registerDoParallel(cl, cores = MC_CORES)
parallel::clusterExport(cl, c("X", "resample"))
parallel::clusterEvalQ(cl, c(library(dplyr), library(data.table)))
#> [[1]]
#> [1] "dplyr"      "stats"      "graphics"   "grDevices"  "utils"
#> [6] "datasets"   "methods"    "base"       "data.table" "dplyr"
#> [11] "stats"      "graphics"   "grDevices"  "utils"      "datasets"
#> [16] "methods"    "base"
#>
#> [[2]]
#> [1] "dplyr"      "stats"      "graphics"   "grDevices"  "utils"
#> [6] "datasets"   "methods"    "base"       "data.table" "dplyr"
#> [11] "stats"      "graphics"   "grDevices"  "utils"      "datasets"
#> [16] "methods"    "base"
#>
#> [[3]]
#> [1] "dplyr"      "stats"      "graphics"   "grDevices"  "utils"
#> [6] "datasets"   "methods"    "base"       "data.table" "dplyr"
#> [11] "stats"      "graphics"   "grDevices"  "utils"      "datasets"
#> [16] "methods"    "base"

Xse_foreach <-
  foreach(b = 1:B, .combine = c, .inorder = FALSE) %dopar% {
    resample(data = X, statistic = mean)
```

```

}

parallel::stopCluster(cl)

tibble(se = Xse_foreach) %>%
  ggplot(aes(x = se)) +
  geom_histogram(bins = 30) +
  xlab(expression(hat(sigma)[B]))

```

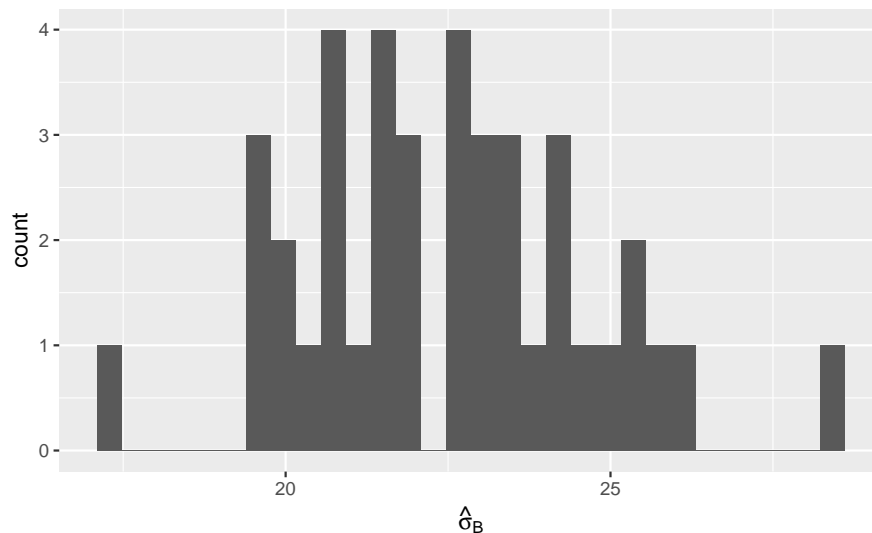


Figure 3.14: Bootstrap replicates for sample mean done by foreach

It gives bootstrap standard error as

```

sd(Xse_foreach)
#> [1] 2.2

```

Are these kinds of parallization useful?

```

microbenchmark::microbenchmark(
  "MCLAPPLY2" = {
    parallel::mclapply(
      1:B,
      function(x) {resample(data = X, statistic = mean)},
      mc.cores = 2
    ) %>%
    unlist()
  },
  "MCLAPPLY3" = {
    parallel::mclapply(
      1:B,
      function(x) {resample(data = X, statistic = mean)},
      mc.cores = 3
    ) %>%
    unlist()
  },
  "MCLAPPLY4" = {
    parallel::mclapply(

```



```

1:B,
function(x) {resample(data = X, statistic = mean)},
mc.cores = 4
) %>%
  unlist()
},
"FORLOOP" = {
  for (b in 1:B) {
    resample(data = X, statistic = mean)
  }
},
times = 5,
unit = "s"
) %>%
  autoplot()

```

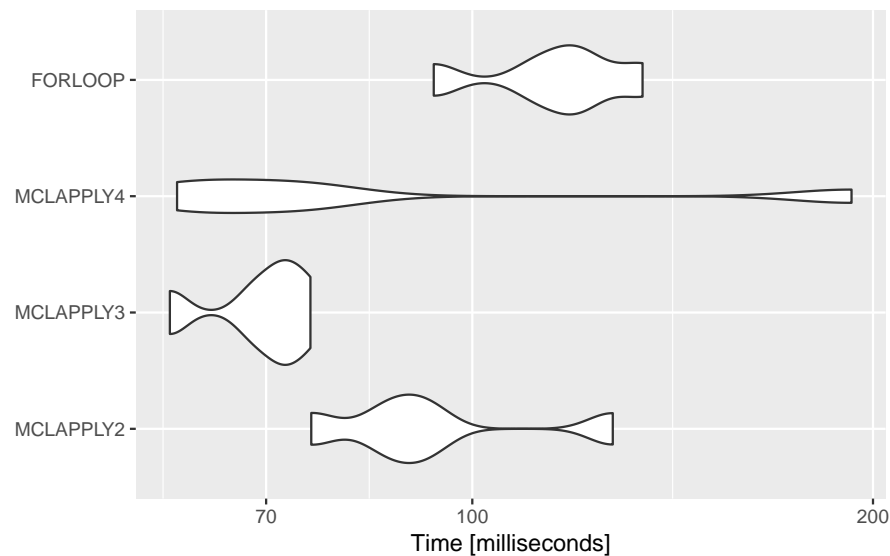


Figure 3.15: Benchmark between mclapply and for loop

Figure 3.15 is comparing for loop with each `mc.cores`. It is clear that parallelization is faster than ordinary loop. In fact, all these procedures can be done by `boot` package.

```
library(boot)
```

Before performing bootstrap, we should define a `statistic` function.

```

mean_boot <- function(x, i) {
  mean(x[i])
}
#-----
boot(data = X %>% pull(), statistic = mean_boot, R = B)
#>
#> ORDINARY NONPARAMETRIC BOOTSTRAP
#>
#>
#> Call:
#> boot(data = X %>% pull(), statistic = mean_boot, R = B)

```

```
#>
#>
#> Bootstrap Statistics :
#>   original   bias   std. error
#> t1*      22.2   0.112     2.07
```

It gives sample mean `original`, bootstrap bias `bias`, and bootstrap `std. error`. We will cover bias later. Due to the programming fact, this is much more faster than the previous one. Also, we can parallelize with this function. `parallel = c("no", "multicore", "snow")`. If we choose "multicore" option, we should specify `ncpus` as in `mclapply`. If "snow", cluster should be supplied in `cl` argument like in `foreach`.

```
boot(
  data = X %>% pull(),
  statistic = mean_boot,
  R = B,
  parallel = "multicore",
  ncpus = MC_CORES
)
#>
#> ORDINARY NONPARAMETRIC BOOTSTRAP
#>
#>
#> Call:
#> boot(data = X %>% pull(), statistic = mean_boot, R = B, parallel = "multicore",
#>       ncpus = MC_CORES)
#>
#>
#> Bootstrap Statistics :
#>   original   bias   std. error
#> t1*      22.2  -0.102     3.04
```

```
microbenchmark::microbenchmark(
  "MCLAPPLY" = {
    parallel::mclapply(
      1:B,
      function(x) {resample(data = X, statistic = mean)}},
      mc.cores = MC_CORES
    ) %>%
    unlist()
  },
  "BOOT" = {
    boot(
      data = X %>% pull(),
      statistic = mean_boot,
      R = B,
      parallel = "multicore",
      ncpus = MC_CORES
    )
  },
  times = 5,
  unit = "s"
) %>%
autoplot()
```

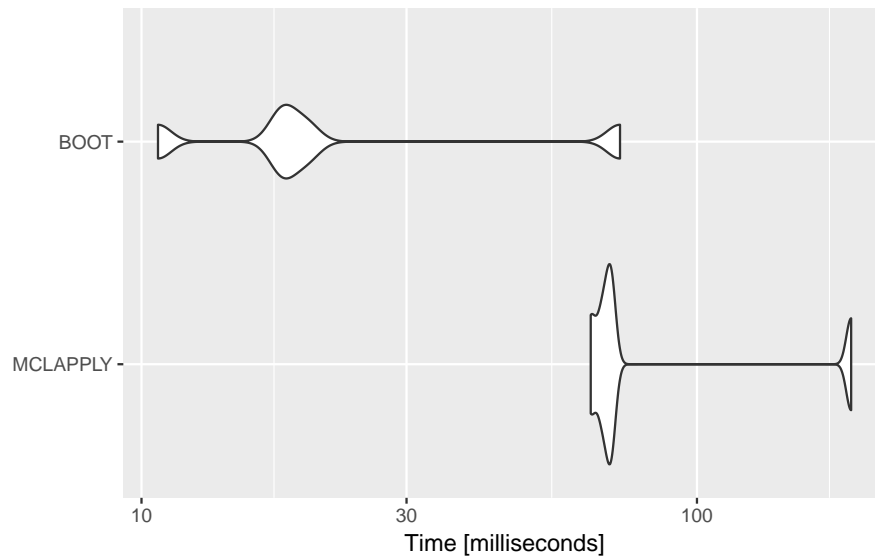


Figure 3.16: Benchmark between mclapply and boot

In Figure 3.16, we can see the difference of the speed.

### 3.6.4 Estimation of correlation coefficient

Consider traditional bootstrap example (Efron and Gong, 1983). The dataset is GPA scores of various entering classes at American law schools in 1973.

```
law %>%
  ggplot(aes(x = LSAT, y = GPA)) +
  geom_point()
```

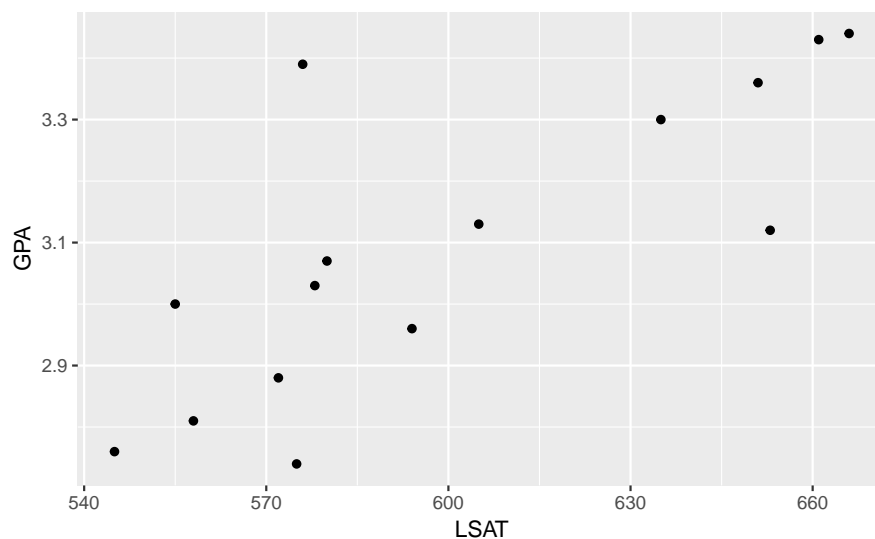


Figure 3.17: The law school data (Efron and Gong (1983))

Each point represents one law school.

- LSAT: average LSAT (national test) score of entering students

- GPA: average GPA score of entering students

**Example 3.9** (Estimation of correlation coefficient). In this  $(Y_i, Z_i) = (\text{LSAT}, \text{GPA})$  data set, we are interested in the correlation  $\rho$  of the two variables. Especially, we want to explore the distribution of  $\hat{\rho}$ .

```
law %>%
  summarise(rho = cor(LSAT, GPA))
#> # A tibble: 1 x 1
#>   rho
#>   <dbl>
#> 1 0.776
```

Let  $\mathbf{X}_i^T = (Y_i, Z_i)$  be each observation. Dependency should be kept, so we should sample  $(Y_i, Z_i)$  pairs, not individuals.

**Algorithm 23:** Estimation of correlation coefficient

```

input :  $n$  observations  $\mathbf{X}_i^T = (Y_i, Z_i)$ ,  $i = 1, \dots, n$ 
1 for  $b \leftarrow 1$  to  $B$  do
2   Sampling with replacement  $\mathbf{X}_1^{(b)}, \dots, \mathbf{X}_n^{(b)}$  from the observed sample;
3   Sample correlation coefficient

```

$$\hat{\rho}_b^* = \frac{\sum (Y_i^{(b)} - \bar{Y}^{(b)})(Z_i^{(b)} - \bar{Z}^{(b)})}{\sqrt{\sum (Y_i^{(b)} - \bar{Y}^{(b)})^2} \sqrt{\sum (Z_i^{(b)} - \bar{Z}^{(b)})^2}}$$

```

   ;
4 end
5  $\bar{\rho}^* = \frac{1}{B} \sum_{b=1}^B \hat{\rho}_b^*$ ;
6 Bootstrap standard error

```

$$\hat{\sigma}_B(\hat{\rho}) = \left[ \frac{1}{B-1} \sum_{b=1}^B (\hat{\rho}_b^* - \bar{\rho}^*)^2 \right]^{\frac{1}{2}}$$

```

   ;
output:  $\hat{\sigma}_B(\hat{\rho})$ 

```

Following Efron and Gong (1983), try  $B = 1000$ .

```
boot_cor <- function(x, i) {
  cor(x[i, 1], x[i, 2])
}
#-----
B <- 1000
(rho <- boot(law, statistic = boot_cor, R = B,
  parallel = "multicore", ncpus = MC_CORES))
#>
#> ORDINARY NONPARAMETRIC BOOTSTRAP
#>
#>
#> Call:
#> boot(data = law, statistic = boot_cor, R = B, parallel = "multicore",
#>       ncpus = MC_CORES)
#>
#>
#> Bootstrap Statistics :
#>   original   bias    std. error
```

```
#> t1*      0.776 -0.0131      0.135
```

**Lemma 3.5** (Gaussian standard error of  $\hat{\rho}$ ). *When the data follow Normal distribution, the standard error of  $\hat{\rho}$  can be estimated by*

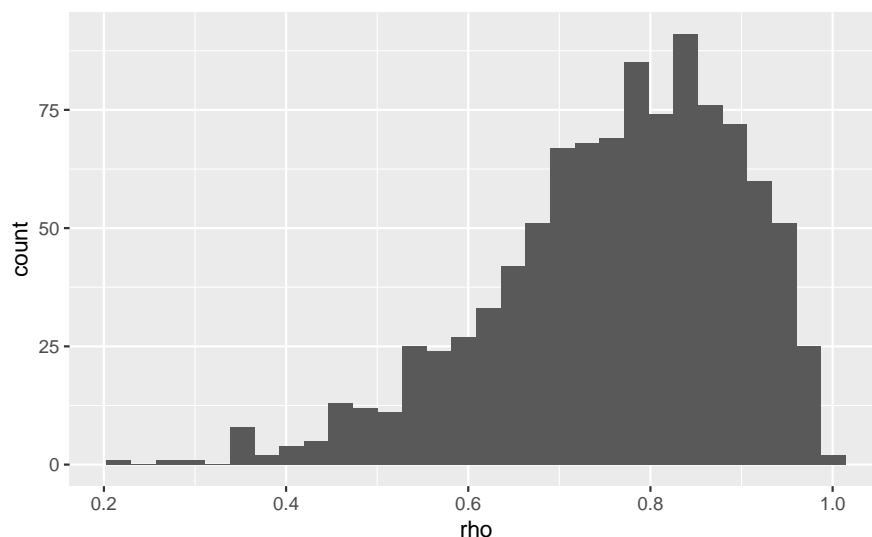
$$\hat{\sigma}_{NORM} = \frac{1 - \hat{\rho}^2}{\sqrt{n - 3}}$$

- `t0` is statistic computed from the sample, i.e. correlation coefficient of the data
- `t` is each bootstrap replicate, `matrix` object

Using `t`, we might draw empirical distribution and get standard error. In the real world, many data for scores follow normal. To check bootstrap works well, we compare the empirical distribution of `rho$t` and one with Lemma 3.5.

$$\hat{\sigma}_{NORM} = \frac{1 - \hat{\rho}^2}{\sqrt{n - 3}} = 0.115$$

```
rho_boot <- tibble(rho = rho$t[,1])
#-----
ggplot(data = rho_boot, aes(x = rho)) +
  geom_histogram(bins = 30)
```



### 3.6.5 Bootstrap bias



# Bibliography

- Efron, B. and Gong, G. (1983). A Leisurely Look at the Bootstrap, the Jackknife, and Cross-Validation. *The American Statistician*, 37(1):36–48.
- Efron, B. and Tibshirani, R. (1994). *An Introduction to the Bootstrap*. CRC Press.
- McGrath, R. N. and Yeh, A. B. (2005). A Quick, Compact, Two-Sample Dispersion Test. *The American Statistician*, 59(1):47–53.