

STAT 2005 – PROGRAMMING LANGUAGES FOR STATISTICS

TUTORIAL 2 RANDOM NUMBERS

2020

LIU Ran

Department of Statistics, The Chinese University of Hong Kong

1 Generate Random Numbers

1.1 Random Seed

- The pseudo random number generator in R requires a starting value, which is called 'random seed' or 'seed'.
- $\left. \begin{array}{l} \text{Same seed} \\ \text{Same command} \end{array} \right\} \Rightarrow \text{Same random numbers}$

Remark 1.1. Set seed in the first line of your model simulation study in order to get reproducible result.

Remark 1.2. Usually, we set an ordinary number to be the seed.

1.2 Random Sample

`sample(x, size, replace = FALSE, prob = NULL)`

```

1 > set.seed(2005)
2 > sample(10)
3 [1] 4 1 5 7 10 2 9 3 8 6
4 > sample(10)
5 [1] 7 5 1 9 2 6 3 8 4 10
6
7 > set.seed(2005)
8 > sample(10)
9 [1] 4 1 5 7 10 2 9 3 8 6
10
11 > set.seed(2005)
12 > a = 1 # non-random process will not affect the seed
13 > sample(10)
14 [1] 4 1 5 7 10 2 9 3 8 6

```

For different parameters:

```

1 > set.seed(2005)
2 > sample(10)
3 [1] 4 1 5 7 10 2 9 3 8 6
4 > sample(5)
5 [1] 5 1 4 2 3
6
7 > set.seed(2005)
8 > sample(5)
9 [1] 4 1 5 3 2
10 > sample(10)
11 [1] 4 2 5 7 3 8 1 10 9 6

```

For different types:

```

1 > set.seed(2005)
2 > runif(3)
3 [1] 0.83248656 0.08690543 0.89776514
4 > sample(10)
5 [1] 1 5 7 8 4 2 10 3 9 6
6
7 > set.seed(2005)

```

```

8 > sample(10)
9 [1] 4 1 5 7 10 2 9 3 8 6
10 > runif(3)
11 [1] 0.006695024 0.599185211 0.826918751

```

An exception? (optional and no need to be remembered)

```

1 > set.seed(2005)
2 > sample(5)
3 [1] 4 1 5 3 2
4
5 > set.seed(2005)
6 > runif(1)
7 [1] 0.8324866
8 > sample(5)
9 [1] 4 1 5 3 2
10
11 > set.seed(2005)
12 > runif(1)
13 [1] 0.8324866
14 > runif(1)
15 [1] 0.08690543
16 > sample(5)
17 [1] 4 1 5 3 2

```

1.3 Random Sample From Commonly Used Distribution

Distribution	R name	Additional arguments = default value
Uniform	unif	min=0, max=1
Binomial	binom	size, prob
Exponential	exp	rate=1
Poisson	pois	lambda
Normal	norm	mean=0, sd=1

Remark 1.3. Take STAT2001 or google for them if you are not familiar with these distributions.

1.3.1 r+R Name: Random Numbers

```

1 > set.seed(2005)
2 > x = rnorm(1000)
3 > c(M=mean(x), V=var(x))
4           M           V
5 0.005746068 0.901851609

```

1.4 d+R Name: Density (pdf/pmf)

```

1 > dunif(0.3, 0, 2)
2 [1] 0.5
3 > dbinom(1, size = 5, prob = 0.7)
4 [1] 0.02835
5 # choose(n, k) Number of combinations
6 > choose(5, 1) * (0.3)^4 * (0.7)^1
7 [1] 0.02835

```

1.5 p+R name: cdf

```

1 > pnorm(0) # P(X<=0)
2 [1] 0.5
3 > pbinom(3, size = 5, prob = 0.7) # P(X<=3)
4 [1] 0.47178

```

```

5 > index = c(0,1,2,3)
6 > sum(choose(5, index) * (0.3)^(5-index) * (0.7)^index)
7 [1] 0.47178

```

1.6 q+R name: Quantiles

```

1 > qnorm(0.5) # P(X<=0) = 0.5
2 [1] 0
3 > qbinom(0.17, size = 5, prob = 0.7) # P(X<=3) >= 0.17
4 [1] 3
5 > qbinom(0.16, size = 5, prob = 0.7) # P(X<=2) >= 0.16
6 [1] 2
7 > pbinom(seq(0,5), size = 5, prob = 0.7) # P(X<=seq(0,5))
8 [1] 0.00243 0.03078 0.16308 0.47178 0.83193 1.00000

```

Remark 1.4. It returns $\min_{q \in I} \{q : P(X \leq q) \geq \alpha\}$, where α is the probability you input, I is the distribution support, q is the quantile you want to search.

2 apply

- **apply**: Use one function to all elements in one(some) dimension(s) of a matrix(array).

```

1 > (y = matrix(rep(1:3,rep(5,3)), nrow=5))
2      [,1] [,2] [,3]
3 [1,]  1   2   3
4 [2,]  1   2   3
5 [3,]  1   2   3
6 [4,]  1   2   3
7 [5,]  1   2   3
8 > apply(y, 1, mean) # operate on each row
9 [1] 2 2 2 2 2
10 > apply(y, 2, mean) # operate on each column
11 [1] 1 2 3

```

```

1 > (z = array(rep(seq(20, by=-2, length=12), 2), dim=c(3, 4, 2)))
2 , , 1
3
4      [,1] [,2] [,3] [,4]
5 [1,]  20  14   8   2
6 [2,]  18  12   6   0
7 [3,]  16  10   4  -2
8
9 , , 2
10
11     [,1] [,2] [,3] [,4]
12 [1,]  20  14   8   2
13 [2,]  18  12   6   0
14 [3,]  16  10   4  -2
15
16 > apply(z, c(1, 3), mean) # one dimension left
17      [,1] [,2]
18 [1,]  11  11
19 [2,]   9   9
20 [3,]   7   7
21
22 > z[1,,1] # each element is the result that 'mean' operates on a vector
23 [1] 20 14 8 2
24
25 > c(mean(z[1,,1]), mean(z[2,,1]), mean(z[3,,1]))
26 [1] 11 9 7
27
28 > c(mean(z[1,,2]), mean(z[2,,2]), mean(z[3,,2]))
29 [1] 11 9 7

```

```

1 > apply(z, 1, mean) # two dimensions left
2 [1] 11 9 7
3
4 > z[1,,] # each element is the result that 'mean' operates on a matrix
5      [,1] [,2]
6 [1,]  20  20
7 [2,]  14  14
8 [3,]   8   8
9 [4,]   2   2
10
11 > c(mean(z[1,,]), mean(z[2,,]), mean(z[3,,]))
12 [1] 11 9 7

```

Remark 2.1. MARGIN argument: Dimensions you want to keep.

3 Monte Carlo Method

3.1 Main theory

The main theoretical basis for Monte Carlo Method is the Law Of Large Numbers(i.i.d):

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} E(X) = \int xp(x)dx.$$

And the central limit theorem could give the approximate confidence intervals for $\mu = E(X)$:

$$\sqrt{n} \frac{(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

Usually, we use the Standard Deviation(SD) of samples to replace σ . ($\sigma = sd(X)$)

Remark 3.1. Standard Deviation is different from Standard Error:

$$\text{SD: } \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2} \quad \text{SE: } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Remark 3.2. More rigorous version for LLN and CLT, please go to the wiki.

3.2 Monte Carlo for integral

Suppose we can draw samples from the distribution $p(x)$ defined on $[a,b]$, and want to calculate the following integral:

$$\int_a^b f(x)dx = \int_a^b \frac{f(x)}{p(x)}p(x)dx = E\left(\frac{f(X)}{p(X)}\right).$$

So, the Monte Carlo estimation for the integral is

$$E\left(\frac{f(X)}{p(X)}\right) \approx \frac{1}{n} \sum_{i=1}^n \left(\frac{f(x_i)}{p(x_i)}\right),$$

where x_i is i.i.d sample from $p(x)$.

Exercise 3.1. Calculate the integral $\int_0^2 \cos(x)dx$:

```

1 > set.seed(2005)
2 > x = runif(1000000, min = 0, max = 2)
3 > mean(cos(x)*2) # p(x) = 1/2
4 [1] 0.9093535
5 > sin(2) - sin(0)
6 [1] 0.9092974

```

4 Exercises

1. Toss a biased coin for 10 times, with $P(\text{Head}) = 0.7$. Simulate the result.
2. Generate a 3×4 matrix of uniform random values between 5 and 9 and find the average of each column.
3. Suppose $Z \sim N(0, 1)$. Find constant c s.t.
 - (a) $P(Z \leq c) = 0.1151$;
 - (b) $P(1 < Z \leq c) = 0.1525$;
 - (c) $P(-c < Z \leq c) = 0.8164$.