# Random Numbers

## STAT 133

Gaston Sanchez

github.com/ucb-stat133/stat133-fall-2016

# Random Numbers in R

Generation of random numbers is at the heart of many statistical methods

# Use of Random Numbers

## Some uses of random numbers

- ▶ Sampling procedures

- ▶ Simulation studies of stochastic processes

- ▶ Analytically intractable mathematical expressions

- ▶ Simulation of a population distribution by resampling from a given sample from that population

- ▶ More general: Simulation, Monte Carlo, Resampling

# Random Samples

- ▶ Many statistical methods rely on random samples:
  - – Sampling techniques
  - – Design of experiments
  - – Surveys
- ▶ Hence, we need a source of random numbers
- ▶ Before computers, statisticians used *tables of random numbers*
- ▶ Now we use computers to generate "random" numbers
- ▶ The random sampling required in most analyses is usually done by the computer

# Generating Random Numbers

- We cannot generate truly random numbers on a computer

- Instead, we generate **pseudo-random** numbers

- i.e. numbers that have the appearance of random numbers

- they *seem* to be randomly drawn from some known distribution

- There are many methods that have been proposed to generate pseudo-random numbers

A very important advantage of using pseudo-random numbers is that, because they are deterministic, they can be reproduced (i.e. repeated)

# Multiple Recursion

- Generate a sequence of numbers in a manner that appears to be random

- Use a deterministic generator that yields numbers recursively (in a fixed sequence)

- The previous $k$ numbers determine the next one

$$x_i = f(x_{i-1}, \ldots, x_{i-k})$$

# Simple Congruential Generator

- Congruential generators were the first reasonable class of pseudo-random number generators
- The congruential method uses modular arithmetic to generate "random" numbers

# Ingredients

- An integer $m$
- An integer $a$ such that $a < m$
- A starting integer $x_0$ (a.k.a. *seed*)

# Simple Congruential Generator

The first number is obtained as:
$x_1 = (a \times x_0) \bmod m$

The rest of the pseudorandom numbers are generated as:
$x_{n+1} = (a \times x_n) \bmod m$

# Simple Congruential Generator

For example if we take $m = 64$, and $a = 3$, then for $x_0 = 17$ we have:

$$x_1 = (3 \times 17) \bmod 64 = 51$$
$$x_2 = (3 \times 51) \bmod 64 = 25$$
$$x_3 = (3 \times 25) \bmod 64 = 11$$
$$x_4 = (3 \times 11) \bmod 64 = 33$$
$$x_5 = (3 \times 33) \bmod 64 = 35$$
$$x_6 = (3 \times 35) \bmod 64 = 41$$
$$\vdots$$

# Congruential algorithm

```r
a <- 3;  m <- 64;  seed <- 17

x <- numeric(20)

x[1] <- (a * seed) %% m

for (i in 2:20) {
  x[i] <- (a * x[i-1]) %% m
}

x[1:16]; x[17:20]

##  [1] 51 25 11 33 35 41 59 49 19 57 43  1  3  9 27 17
## [1] 51 25 11 33
```
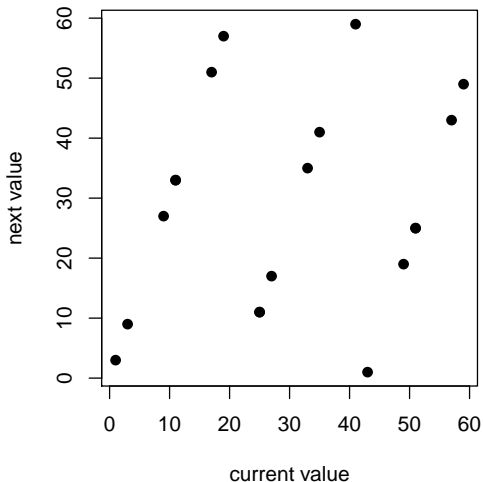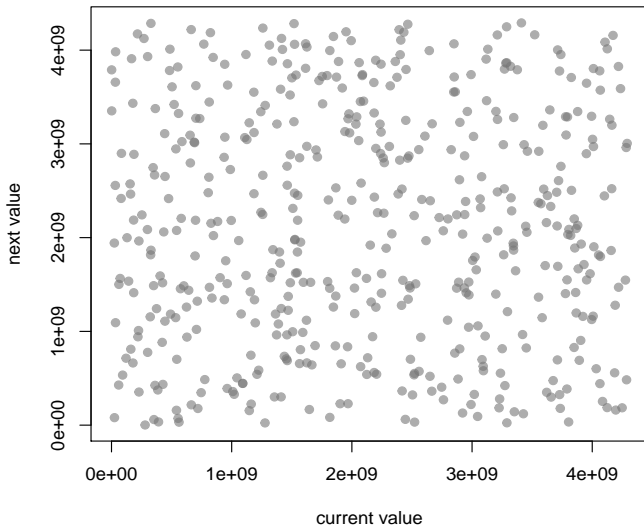
# Congruential algorithm

```
cong <- function(n, a = 69069, m = 2^32, seed = 17) {
  x <- numeric(n)
  x[1] <- (a * seed) %% m
  for (i in 2:n) {
    x[i] <- (a * x[i-1]) %% m
  }
  x
}

y <- cong(20, a = 3, m = 64, seed = 17)
```

# Congruential algorithm

```
plot(y[1:(n-1)], y[2:n], pch = 19,
     xlab = 'current value', ylab = 'next value')
```

```
cong(n, a = 69069, m = 2^32, seed = 17)
```

# Ingredients

- $m$ is the modulus $(0 < m)$
- $a$ is the multiplier $(0 < a < m)$
- $x_0$ is the seed $(0 \leq x_0 \leq m)$

The period length of a Random Generator Number is at most $m$, and for some choices of $a$ much less than that.

# About the seed

- You can reproduce your simulation results by controlling the seed
- `set.seed()` allows to do this
- Every time you perform simulation studies indicate the random number generator and the seed that was used

# About the seed

```
# set the seed
set.seed(69069)

runif(3)  # call the uniform RNG

## [1] 0.1648855 0.9564664 0.3345479

runif(3)  # call runif() again

## [1] 0.01109596 0.18654873 0.94657805

# set the seed back to 69069
set.seed(69069)
runif(3)

## [1] 0.1648855 0.9564664 0.3345479
```

# Simple Congruential Generator

We can use a congruential algorithm to generate uniform numbers

We'll describe one of the simplest methods for simulating independent uniform random variables on the interaval [0, 1]

# Generating Uniform Numbers

The generator proceeds as follows

$$x_1 = (ax_0) \bmod m$$
$$u_1 = x_1/m$$

$u_1$ is the first pseudo-random number, taking some value between 0 and 1.

# Ingredients

The second pseudorandom number is obtained as:

$$x_2 = (ax_1) \bmod m$$
$$u_2 = x_2/m$$

$u_2$ is another pseudorandom number.

# Generating Uniform Numbers

- If $m$ and $a$ are chosen properly, it is difficult to predict the value $u_2$ given $u_1$
- For most practical purposes $u_2$ is approximately independent of $u_1$

# Simple Congruential Generator

For example if we take $m = 7$, and $a = 3$, then for $x_0 = 2$ we have:

$$x_1 = (3 \times 2) \bmod 7 = 6, \quad u_1 = 0.857$$
$$x_2 = (3 \times 6) \bmod 7 = 4, \quad u_2 = 0.571$$
$$x_3 = (3 \times 4) \bmod 7 = 5, \quad u_3 = 0.714$$
$$x_4 = (3 \times 5) \bmod 7 = 1, \quad u_4 = 0.143$$
$$x_5 = (3 \times 1) \bmod 7 = 3, \quad u_5 = 0.429$$
$$x_6 = (3 \times 3) \bmod 7 = 2, \quad u_6 = 0.286$$
$$\vdots$$

# Random Numbers in R

## R uses a pseudo random number generator

- It starts with a **seed** and an **algorithm** (i.e. function)
- The seed is plugged into the algorithm and a number is returned
- That number is then plugged into the algorithm and the next number is created
- The algorithm is such that the produced numbers behave like random numbers

# RNG functions in R

| Function | Description |
|---|---|
| runif() | Uniform |
| rbinom() | Binomial |
| rmultinom() | Multinomial |
| rnbinom() | Negative binomial |
| rpois() | Poisson |
| rnorm() | Normal |
| rbeta() | Beta |
| rgamma() | Gamma |
| rchisq() | Chi-squared |
| rcauchy() | Cauchy |

See more info: ?Distributions

# Random Number Functions

`runif(n, min = 0, max = 1)` sample from the uniform distribution on the interval (0,1)

The chance the value drawn is:
- between 0 and $1/3$ has chance $1/3$
- between $1/3$ and $1/2$ has chance $1/6$
- between $9/10$ and 1 has chance $1/10$

# Random Number Functions

`rnorm(n, mean = 0, sd = 1)` sample from the normal distribution with center = `mean` and spread = `sd`

# Random Number Functions

`rnorm(n, mean = 0, sd = 1)` sample from the normal distribution with center = `mean` and spread = `sd`

`rbinom(n, size, prob)` sample from the binomial distribution with number of trials = `size` and chance of success = `prob`

# Bootstrapping

# Let's Generalize

- A **statistic** is just a function of a random sample

- Statistics are used as estimators of quantities of interest about the distribution, called **parameters**

- Statistics are random variables

- Parameters are NOT random variables

# Let's Generalize

- In simple cases, we can study the *sampling distribution* of the statistic analytically

- e.g. we can prove under mild conditions that the distribution of the sample proportion is close to normal for large sample sizes

- In more complicated cases we can turn to simulation

# Sampling Distributions

- In our example $X_1, X_2, \ldots, X_n$ are independent observations from the same distribution

- The distribution has center (mean value) $\mu$ and spread (standard deviation) $\sigma$

- e.g. interest in the distribution of $median(X_1, X_2, \ldots, X_n)$

- We take many samples of size $n$, and study the behavior of the sample medians

# Some Limitations

▶ Consider *t-test* procedures for inference about means
▶ Most classical methods rest on the use of Normal Distributions
▶ However, most real data are not Normal
▶ We cannot use $t$ confidence intervals for strongly skewed data (unless samples are large)
▶ What about inference for a *ratio* of means? (no simple traditional inference)

## Fundamental Reasoning

- Apply computer power to relax some of the conditions needed in traditional tests
- Have tools to do inference in new settings
- What would happen if we applied this method many times?

# Bootstrap Idea

- ▶ Statistical inference is based on the sampling distributions of sample statistics
- ▶ A sampling distribution is based on many random samples from the population
- ▶ The bootstrap is a way of finding the sampling distribution (approximately)

# Bootstrap Samples

```r
x <- c(3.15, 0, 1.58, 19.65, 0.23, 2.21)
mean(x)

## [1] 4.47
```

# Bootstrap Samples

```
x <- c(3.15, 0, 1.58, 19.65, 0.23, 2.21)
mean(x)

## [1] 4.47
```

```
(x1 <- sample(x, size = 6, replace = TRUE))

## [1]  3.15  0.00  2.21  3.15 19.65  0.00

mean(x1)

## [1] 4.693333
```

# Bootstrap Samples

```
(x2 <- sample(x, size = 6, replace = TRUE))

## [1] 3.15 2.21 3.15 1.58 1.58 1.58

mean(x2)

## [1] 2.208333

(x3 <- sample(x, size = 6, replace = TRUE))

## [1]  1.58  2.21  0.00  0.23 19.65  1.58

mean(x3)

## [1] 4.208333
```

# Procedure for Bootstrapping

- Repeatedly sampling **with replacement** from a random sample

- Each bootstrap sample is the same size as the original sample

- Calculate the statisc of interest (e.g. mean, median, sd)

- Draw hundreds or thousands of samples

- Obtain a bootstrap distribution

# Bootstrap Samples

```
bootstrap <- numeric(1000)

for (b in 1:1000) {
  boot_sample <- sample(x, size = 6, replace = TRUE)
  bootstrap[b] <- mean(boot_sample)
}
```

# Bootstrap Distribution

```
hist(bootstrap, col = 'gray70', las = 1)
```

**Histogram of bootstrap**

# How does Bootstrapping work?

- We are not using the resamples as if they were real data

- Bootstrap samples is not a substitute to gather more data to improve accuracy

- The bootstrap idea is to use the resample statistics to estimate how the sample statistic varies from the studied random sample

- The bootstrap distribution approximates the sampling distribution of the statistic

# Bootstrap samples

## Computing the bootstrap distribution implies

- Calulate the statistic for each sample

- The distribution of these resample statistics is the bootstrap distribution

- A bootstrap sample is the same size as the original random sample

# Another Example

# Bootstrap resampling

```r
# Iris Virginica subset (the "population")
virginica <- subset(iris, Species == 'virginica')

# random sample of Petal.Length (size = 5)
set.seed(7359)
rand_sample <- sample(virginica$Petal.Length, size = 5)
rand_sample

## [1] 5.1 5.6 5.8 5.7 5.8
```

# Bootstrap resampling

```r
# create 500 bootstrap samples of size 5 with replacement
resamples <- 500
n <- length(rand_sample)

boot_stats <- numeric(resamples)

for (i in 1:resamples) {
  boot_sample <- sample(rand_sample, size = n, replace = TRUE)
  boot_stats[i] <- mean(boot_sample)
}
```

# Bootstrap resampling

```r
# "population" mean
mean(virginica$Petal.Length)

## [1] 5.552

# bootstrap mean
mean(boot_stats)

## [1] 5.60028
```
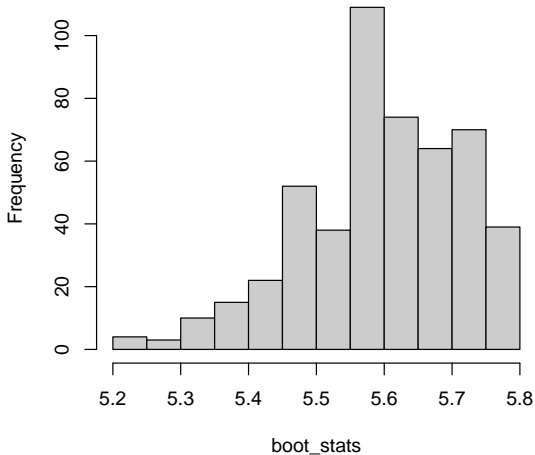
# Bootstrap Distribution



**Histogram of boot_stats**

# Bootstrap standard error

The bootstrap standard error is just the standard deviation of
the bootstrap samples

```
# descriptive statistics
summary(boot_stats)

##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    5.20    5.52    5.60    5.60    5.70    5.80

# Bootstrap Standard Error
sd(boot_stats)

## [1] 0.1167183
```