

# MA678 Homework 3

JingjianGao

9/27/2022

## 4.4 Designing an experiment

You want to gather data to determine which of two students is a better basketball shooter. You plan to have each student take  $N$  shots and then compare their shooting percentages. Roughly how large does  $N$  have to be for you to have a good chance of distinguishing a 30% shooter from a 40% shooter?

```
#Difference in percentage is 10%  
#Critical value of this experiment is 1.96, 95% confidence interval  
#In order to have a good chance of distinguishing, we want 0.1/sd to be bigger than 1.96.  
#Thus, 0.1/sqrt(0.1*0.9/N) > 1.96  
#(1.96*0.3/0.1)^2 < N  
#Therefore N has to be bigger than 35.  
#Hence, N has to be roughly 35 for me to have a good chance of  
#distinguishing 30% shooter from a 40% shooter
```

## 4.6 Hypothesis testing

The following are the proportions of girl births in Vienna for each month in girl births 1908 and 1909 (out of an average of 3900 births per month):

```
birthdata <- c(.4777,.4875,.4859,.4754,.4874,.4864,.4813,.4787,.4895,.4797,.4876,.4859,  
               .4857,.4907,.5010,.4903,.4860,.4911,.4871,.4725,.4822,.4870,.4823,.4973)
```

The data are in the folder **Girls**. These proportions were used by von Mises (1957) to support a claim that that the sex ratios were less variable than would be expected under the binomial distribution. We think von Mises was mistaken in that he did not account for the possibility that this discrepancy could arise just by chance.

(a)

Compute the standard deviation of these proportions and compare to the standard deviation that would be expected if the sexes of babies were independently decided with a constant probability over the 24-month period.

```
library(dplyr)  
  
##  
## Attaching package: 'dplyr'  
  
## The following objects are masked from 'package:stats':  
##  
##   filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union
```

```
library(ggplot2)
births <- 3900
girls_1908 <- c(.4777,.4875,.4859,.4754,.4874,.4864,.4813,.4787,.4895,.4797,.4876,.4859)
girls_1909 <- c(.4857,.4907,.5010,.4903,.4860,.4911,.4871,.4725,.4822,.4870,.4823,.4973)
standardDeviation <- sd(birthdata)
girls_mean <- mean(birthdata)
expectedStandardDeviation <- sqrt(girls_mean*(1-girls_mean)/births)
#population > 20 so we use sqrt(P*(1-P)/n)
standardDeviation
```

```
## [1] 0.006409724
```

```
expectedStandardDeviation
```

```
## [1] 0.008003121
```

(b)

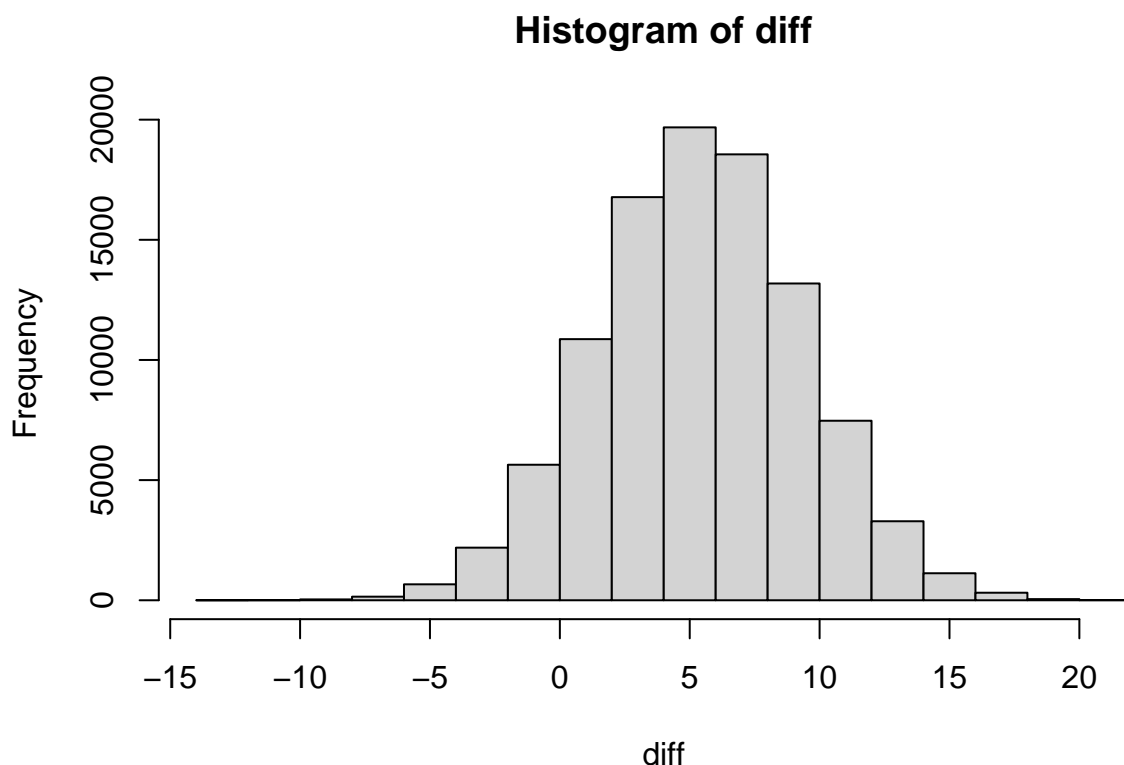
The observed standard deviation of the 24 proportions will not be identical to its theoretical expectation. In this case, is this difference small enough to be explained by random variation? Under the randomness model, the actual variance should have a distribution with expected value equal to the theoretical variance, and proportional to a  $\chi^2$  random variable with 23 degrees of freedom; see page 53.

*#In this case, the difference is small enough to be explained by random variation*

## 5.5 Distribution of averages and differences

The heights of men in the United States are approximately normally distributed with mean 69.1 inches and standard deviation 2.9 inches. The heights of women are approximately normally distributed with mean 63.7 inches and standard deviation 2.7 inches. Let  $x$  be the average height of 100 randomly sampled men, and  $y$  be the average height of 100 randomly sampled women. In R, create 1000 simulations of  $x - y$  and plot their histogram. Using the simulations, compute the mean and standard deviation of the distribution of  $x - y$  and compare to their exact values.

```
set.seed(1000)
N <- 100
men_mean <- 69.1
women_mean <- 63.7
men_sd <- 2.9
women_sd <- 2.7
simulations <- 1000
men_observation <- replicate(simulations,rnorm(N,mean=men_mean,sd=men_sd))
women_observation <- replicate(simulations,rnorm(N,mean=women_mean,sd=women_sd))
men_matrix <- matrix(men_observation,nrow=N,ncol=simulations)
women_matrix <- matrix(women_observation,nrow=N,ncol=simulations)
diff <- men_matrix - women_matrix
hist(diff)
```



```
mean(diff)
```

```
## [1] 5.394219
```

```
sd(diff)
```

```
## [1] 3.953435
```

## 5.8 Coverage of confidence intervals:

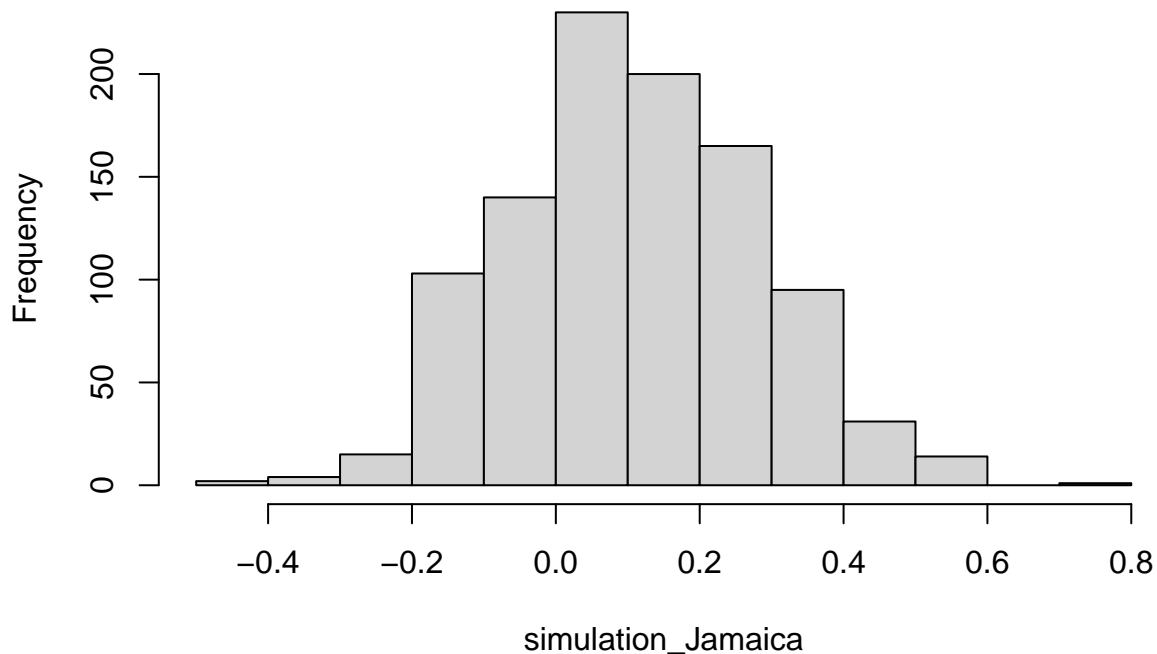
On page 15 there is a discussion of an experimental study of an education-related intervention in Jamaica, in which the point estimate of the treatment effect, on the log scale, was 0.35 with a standard error of 0.17. Suppose the true effect is 0.10—this seems more realistic than the point estimate of 0.35—so that the treatment on average would increase earnings by 0.10 on the log scale. Use simulation to study the statistical properties of this experiment, assuming the standard error is 0.17.

(a)

Simulate 1000 independent replications of the experiment assuming that the point estimate is normally distributed with mean 0.10 and standard deviation 0.17.

```
set.seed(608)
Sims <- 1000
simulation_Jamaica <- rnorm(Sims, mean=0.10, sd=0.17)
hist(simulation_Jamaica)
```

## Histogram of simulation\_Jamaica



(b)

For each replication, compute the 95% confidence interval. Check how many of these intervals include the true parameter value.

```
upper_jamica <- 0.1 + qnorm(0.975)*0.17/sqrt(1000)
lower_jamica <- 0.1 - qnorm(0.975)*0.17/sqrt(1000)
```

(c)

Compute the average and standard deviation of the 1000 point estimates; these represent the mean and standard deviation of the sampling distribution of the estimated treatment effect.

```
mean(simulation_Jamica)
```

```
## [1] 0.1088936
```

```
sd(simulation_Jamica)
```

```
## [1] 0.1727075
```

## 10.3 Checking statistical significance

In this exercise and the next, you will simulate two variables that are statistically independent of each other to see what happens when we run a regression to predict one from the other. Generate 1000 data points from a normal distribution with mean 0 and standard deviation 1 by typing `var1 <- rnorm(1000,0,1)` in R. Generate another variable in the same way (call it `var2`). Run a regression of one variable on the other. Is the slope coefficient “statistically significant”? We do not recommend summarizing regressions in this way, but it can be useful to understand how this works, given that others will do so.

```
var1 <- rnorm(1000,0,1)
var2 <- rnorm(1000,0,1)
```

```
Regression <- lm(var1~var2)
summary(Regression)
```

```
##
## Call:
## lm(formula = var1 ~ var2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.8074 -0.6546  0.0172  0.6553  3.3885
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.006685   0.031100   0.215   0.830
## var2        0.034965   0.031122   1.123   0.262
##
## Residual standard error: 0.9828 on 998 degrees of freedom
## Multiple R-squared:  0.001263,    Adjusted R-squared:  0.0002624
## F-statistic: 1.262 on 1 and 998 DF,  p-value: 0.2615
```

*#The slope coefficient is statistically significant*

### 11.3 Coverage of confidence intervals

Consider the following procedure:

- Set  $n = 100$  and draw  $n$  continuous values  $x_i$  uniformly distributed between 0 and 10. Then simulate data from the model  $y_i = a + bx_i + \text{error}_i$ , for  $i = 1, \dots, n$ , with  $a = 2$ ,  $b = 3$ , and independent errors from a normal distribution.
- Regress  $y$  on  $x$ . Look at the median and mad sd of  $b$ . Check to see if the interval formed by the median  $\pm 2$  mad sd includes the true value,  $b = 3$ .
- Repeat the above two steps 1000 times.

```
n <- 100
sim_11.3 <- function(n,a=2,b=3){
  for (i in 1:100) {
    x_i <- runif(n,min=0,max=10)
    error_i <- rnorm(n,0,1)
    y_i=a+b*x_i+error_i
    reg11.3 <- lm(y_i~x_i)
  }
  return(y_i)
}
replication_11.3 <- replicate(1000,sim_11.3(100,a=2,b=3))
```

(a)

True or false: the interval should contain the true value approximately 950 times. Explain your answer.

#True, the interval should contain the true value approximately 950 times. #Because out of the 1000 simulations, 95% of the time the mean is within the upper and lower interval

(b)

Same as above, except the error distribution is bimodal, not normal. True or false: the interval should contain the true value approximately 950 times. Explain your answer.

#The interval will not contain the true value approximately 950 times. This is because bimodal distribution has two peaks and the interval will shift.

###PS #I tried my best but I still don't quite understand about this simulation thing. Please help me.