

# Homework 06

Simulation

*Frank Dong*

*Nov 2, 2018*

## Discrete probability simulation:

suppose that a basketball player has a 60% chance of making a shot, and he keeps taking shots until he misses two in a row. Also assume his shots are independent (so that each shot has 60% probability of success, no matter what happened before).

1. Write an R function to simulate this process.

```
i=2
shot <- NA
shot[1] <- rbinom(1,1,.6)
while(i>1) {
  shot[i] <- rbinom (1, 1, .6)
  if(shot[i]==0 & shot[i-1]==0) break
  i=1+i
}
```

2. Put the R function in a loop to simulate the process 1000 times. Use the simulation to estimate the mean, standard deviation, and distribution of the total number of shots that the player will take.

```
p <- 0.6
n.sims <- 1000
n.balls <- rep(NA,n.sims)
n.succ <- rep(NA,n.sims)
for (s in 1:n.sims) {
  i=2
  shot <- NA
  shot[1] <- rbinom(1,1,.6)
  while(i>1) {
    shot[i] <- rbinom (1, 1, .6)
    if(shot[i]==0 & shot[i-1]==0) break
    i=1+i
  }
  n.balls[s]=i
  n.succ[s] <- sum(shot==1)
}
mean(n.balls)
```

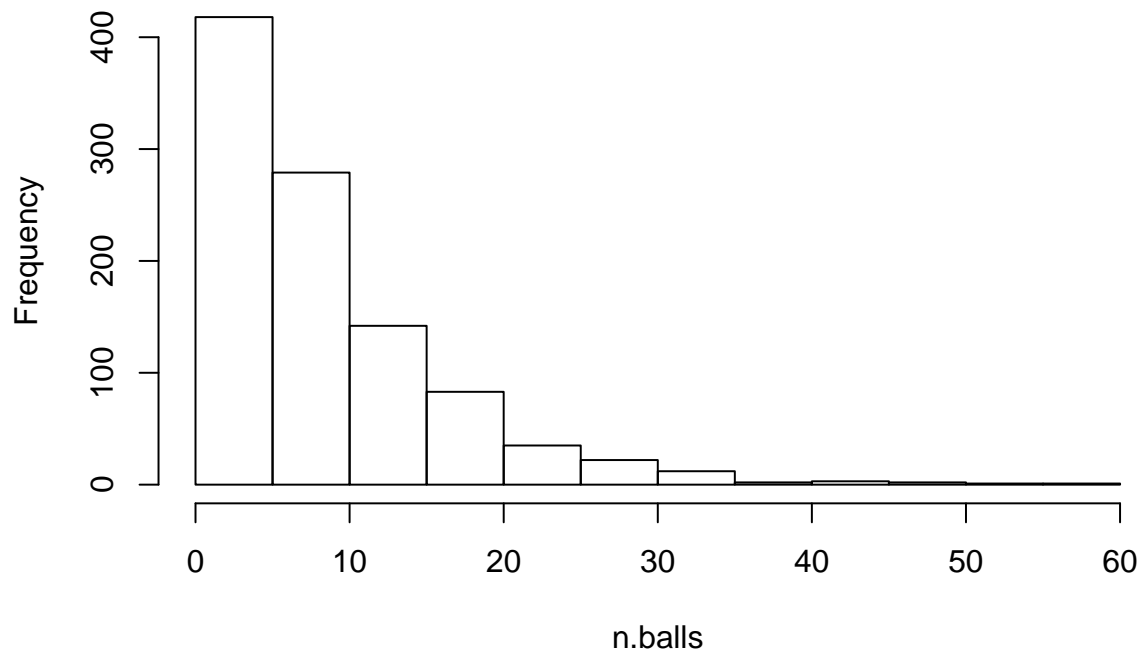
```
## [1] 9.002
```

```
sd(n.balls)
```

```
## [1] 7.622079
```

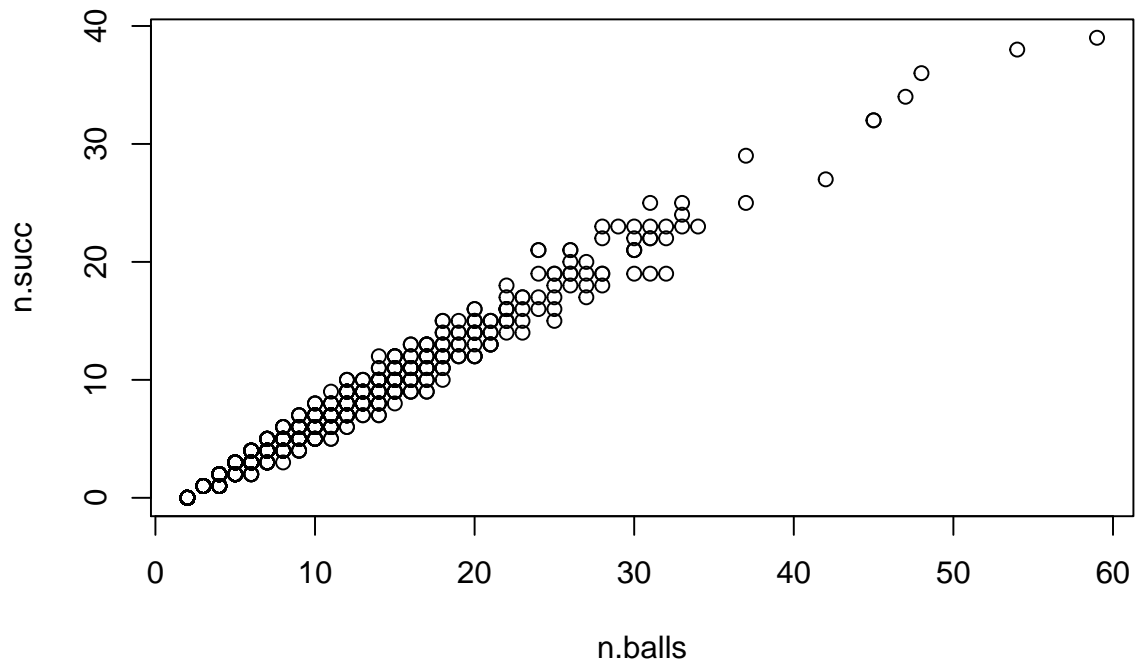
```
hist (n.balls)
```

**Histogram of n.balls**



3. Using your simulations, make a scatterplot of the number of shots the player will take and the proportion of shots that are successes.

```
plot(n.balls, n.succ)
```



## Continuous probability simulation:

the logarithms of weights (in pounds) of men in the United States are approximately normally distributed with mean 5.13 and standard deviation 0.17; women with mean 4.96 and standard deviation 0.20. Suppose 10 adults selected at random step on an elevator with a capacity of 1750 pounds. What is the probability that the elevator cable breaks?

```
library(base)
## construct a function
simulation <- function(){
  weight <- rep(NA,10)
  a <- sample(1:10000,10)

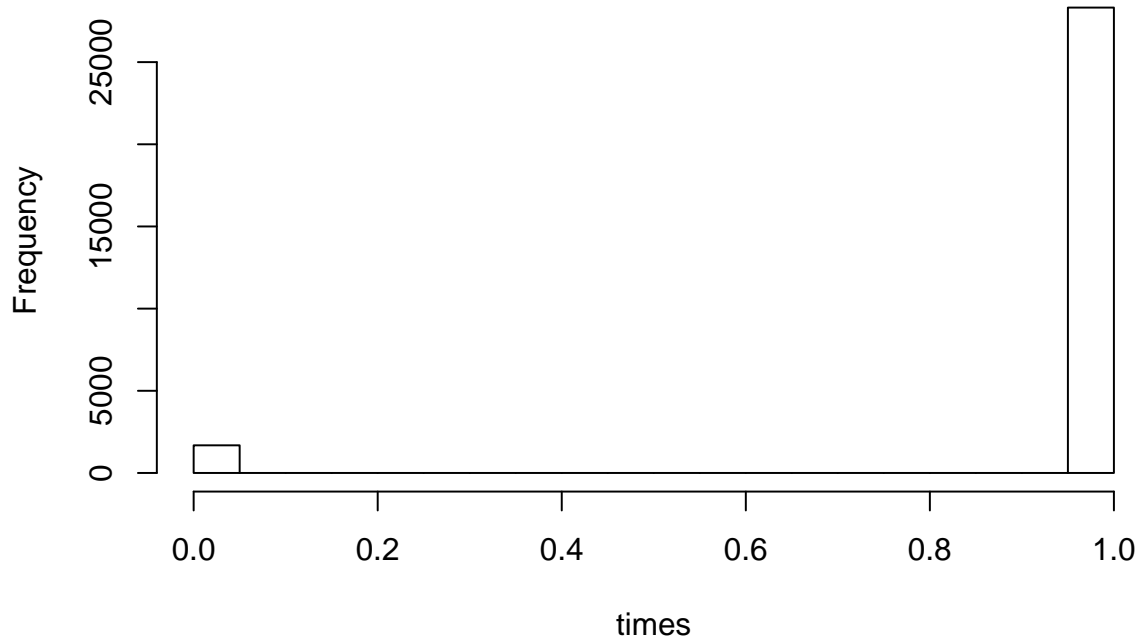
  for (i in 1:10) {
    logweight <- ifelse(a[i]%%2==0,rnorm(1,5.13,0.17),rnorm(1,4.96,0.2))
    weight[i] <- exp(logweight)
  }
  s=sum(weight)
  if (s<=1750){
    t=1}
  else if(s>1750){
    t=0
  }
  return(t)
}

#use the function for simulation
count <- 30000
times <- rep(NA,count)
for (i in 1:count) {
  test1<- simulation()
  times[i] <- test1
}

#probability
p <- sum(times)/count

hist(times)
```

## Histogram of times



p

```
## [1] 0.9439667
```

## Predictive simulation for linear regression:

take one of the models from previous excessive that predicts course evaluations from beauty and other input variables. You will do some simulations.

```
prof <- read.csv("http://www.stat.columbia.edu/~gelman/arm/examples/beauty/ProfEvaltnsBeautyPublic.csv")

# convert into factors
prof$profnumber <- as.factor(prof$profnumber)
prof$female <- as.factor(prof$female)

# convert dummy `class*` variables into a factor
dummies <- prof[, 18:47]
prof$class <- factor(apply(dummies, FUN=function(r) r %>% 1:30, MARGIN=1))

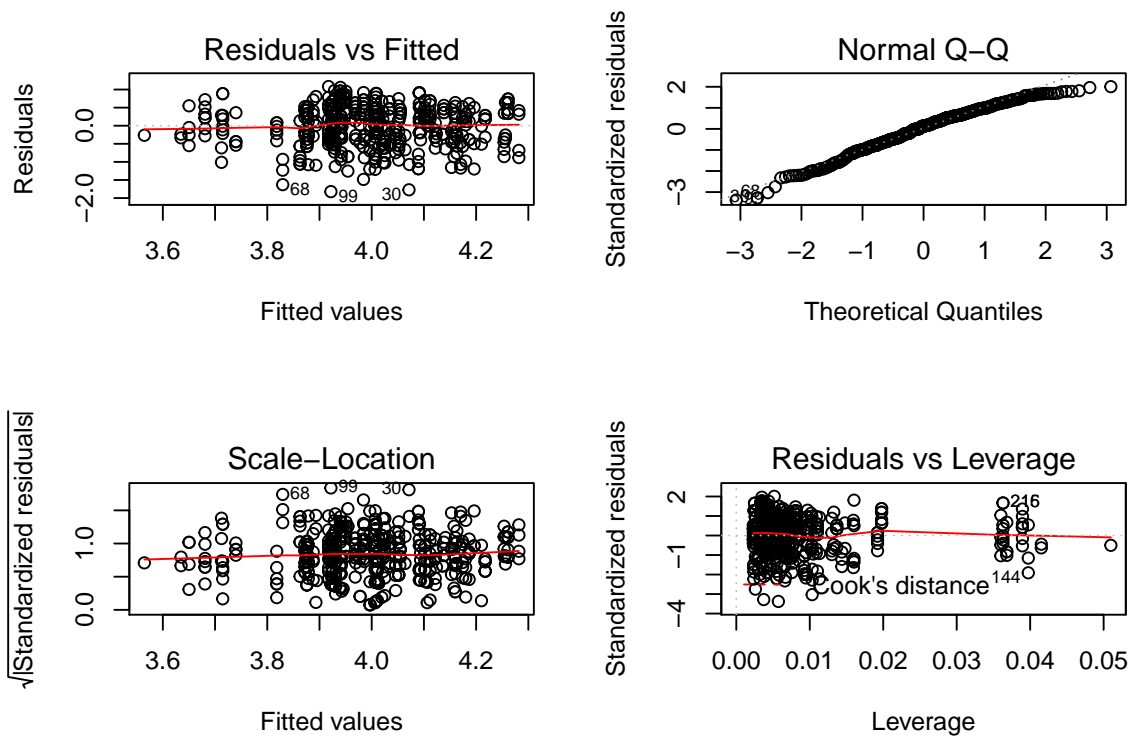
# remove dummy variables
prof <- prof[-c(18:47)]

# normalise and centre professor evaluation (all other predictors are binary)
prof$c.profevaluation <- prof$profevaluation - mean(prof$profevaluation) / (2 * sd(prof$profevaluation))
```

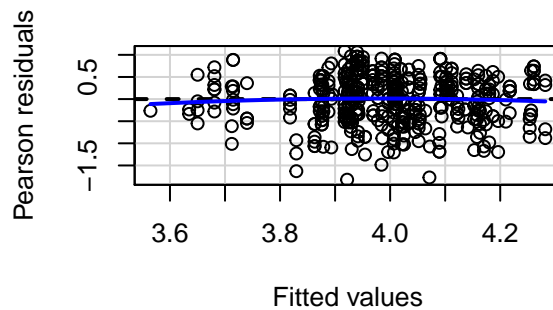
1. Instructor A is a 50-year-old woman who is a native English speaker and has a beauty score of 1. Instructor B is a 60-year-old man who is a native English speaker and has a beauty score of - .5. Simulate 1000 random draws of the course evaluation rating of these two instructors. In your simulation, account for the uncertainty in the regression parameters (that is, use the `sim()` function) as well as the predictive uncertainty.

```
m2 <- lm(data=prof, courseevaluation~prof$btystdave+nonenglish+age)
```

```
par(mfrow=c(2,2))
plot(m2)
```



```
residualPlot(m2)
```



We found that this model fits not bad. So we can simulate using this model.

```
sim2 <- sim(m2,1000)
```

```
coef_prof <- coef(sim2)
```

```
InstructorA <- coef_prof[,1]+coef_prof[,2]*1+0+50*coef_prof[,4]
```

```
InstructorB <- coef_prof[,1]+coef_prof[,2]*(-0.5)+0+60*coef_prof[,4]
```

```
mean(InstructorA)
```

```
## [1] 4.165806
```

```
sd(InstructorA)
```

```
## [1] 0.04531986
```

```
mean(InstructorB)
```

```
## [1] 3.967746
```

```
sd(InstructorB)
```

```
## [1] 0.03875464
```

2. Make a histogram of the difference between the course evaluations for A and B. What is the probability that A will have a higher evaluation?

```
#we simulate 10000 times
```

```
sim2 <- sim(m2,10000)
```

```
coef_prof <- coef(sim2)
```

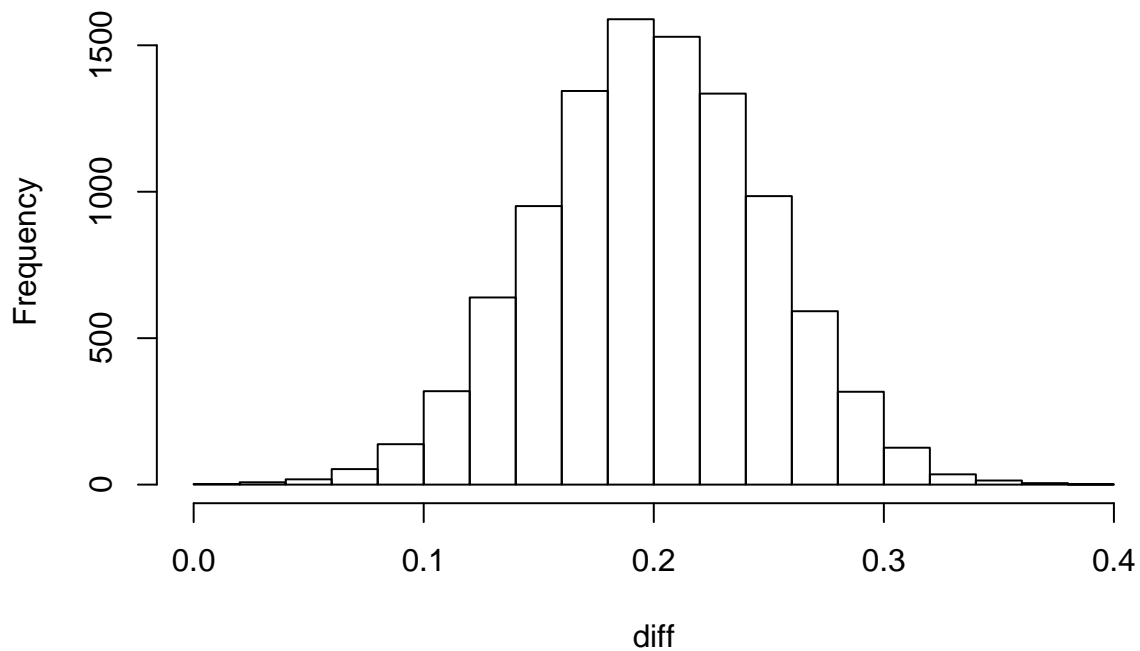
```
InstructorA <- coef_prof[,1]+coef_prof[,2]*1+0+50*coef_prof[,4]
```

```
InstructorB <- coef_prof[,1]+coef_prof[,2]*(-0.5)+0+60*coef_prof[,4]
```

```
diff <- InstructorA-InstructorB
```

```
hist(diff)
```

**Histogram of diff**



```
#probability that A will have a higher evaluation
```

```
t <- diff[which(diff>=0)]
```

```
p <- length(t)/length(diff)
```

```
p
```

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

## How many simulation draws are needed:

take the model from previous exercise that predicts course evaluations from beauty and other input variables. Use `display()` to summarize the model fit. Focus on the estimate and standard error for the coefficient of beauty.

```
beauty <- read.csv("http://www.stat.columbia.edu/~gelman/arm/examples/beauty/ProfEvaltnsBeautyPublic.csv")
```

1. Use `sim()` with `n.sims = 10000`. Compute the mean and standard deviations of the 1000 simulations of the coefficient of beauty, and check that these are close to the output from `display`.

```
display(m2)
```

```
## lm(formula = courseevaluation ~ prof$btystdave + nonenglish +
##      age, data = prof)
##              coef.est coef.se
## (Intercept)    4.02    0.13
## prof$btystdave  0.14    0.03
## nonenglish     -0.33    0.11
## age            0.00    0.00
## ---
## n = 463, k = 4
## residual sd = 0.54, R-Squared = 0.06
```

```
sim2 <- sim(m2,10000)
coef_prof <- coef(sim2)
beautm <- coef_prof[,2]
#mean
mean(beautm)
```

```
## [1] 0.1353915
```

```
sd(beautm)
```

```
## [1] 0.03332276
```

Notice that the output from the display is 0.14, while the coefficient from the simulation is 0.1352.  $0.1352 > 0.14 - 0.03 \times 2$ , so they are close to the output from display.

2. Repeat with `n.sims = 1000`, `n.sims = 100`, and `n.sims = 10`. Do each of these a few times in order to get a sense of the simulation variability.

```
sim3 <- sim(m2,1000)
coef_prof <- coef(sim3)
beautm <- coef_prof[,2]
#mean
mean(beautm)
```

```
## [1] 0.1338241
```

```
sd(beautm)
```

```
## [1] 0.03498679
```

```
sim4 <- sim(m2,100)
coef_prof <- coef(sim4)
beautm <- coef_prof[,2]
```

```
#mean
mean(beautm)
```

```
## [1] 0.135522
```

```
sd(beautm)
```

```
## [1] 0.03209861
```

```
sim5 <- sim(m2,10)
coef_prof <- coef(sim5)
beautm <- coef_prof[,2]
```

```
#mean
mean(beautm)
```

```
## [1] 0.1262112
```

```
sd(beautm)
```

```
## [1] 0.02947679
```

3. How many simulations were needed to give a good approximation to the mean and standard error for the coefficient of beauty?

Depending on what we simulate, we need 10000 times to give a good approximation to the mean and std error.

## Predictive simulation for linear regression:

using data of interest to you, fit a linear regression model. Use the output from this model to simulate a predictive distribution for observations with a particular combination of levels of all the predictors in the regression.

I'm interested in the beauty data, so I still use this linear model to do prediction.

```
m4 <- lm(data=prof, courseevaluation~prof$btystdave+nonenglish+age)
```

```
xtilde <- cbind(prof$age,prof$btystdave,prof$nonenglish)
pred <- function(xpred,lmfit){
  npred <- dim(xpred)[1]
  simlm <- sim(lmfit,1)
  coef_prof <- coef(simlm)
```

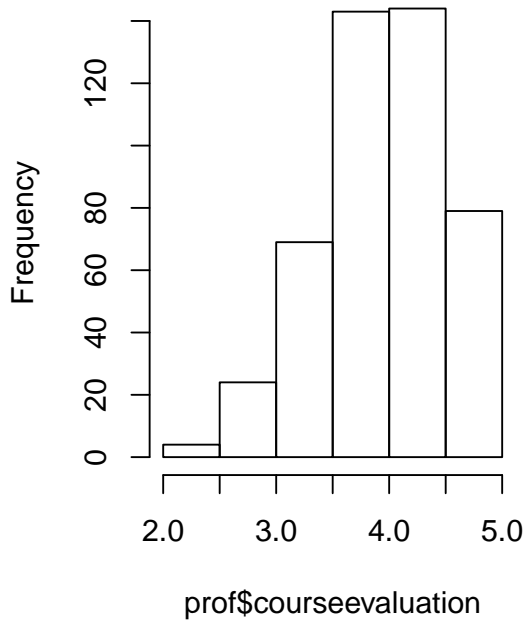
```
  ypred <- rnorm(npred,coef_prof[,1]+coef_prof[,2]*xpred[,2]+xpred[,1]*coef_prof[,4]+coef_prof[,3]*xpred[,3])
}
```

```
yprediction <- replicate(1000,pred(xtilde,m4))
yprediction <- rowSums(yprediction)/1000
par(mfrow=c(1,2))
```

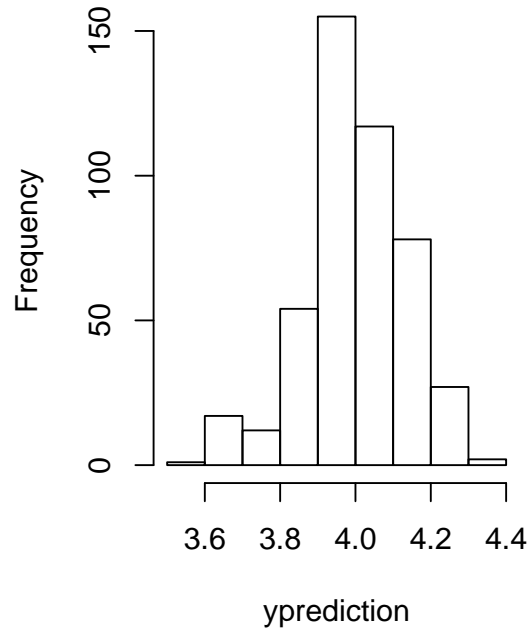


```
hist(prof$courseevaluation)
hist(yprediction)
```

**Histogram of prof\$courseevaluation**



**Histogram of yprediction**



We find that these two histogram look alike. So the affectness is not bad.

## Repeat the previous exercise using a logistic regression example.

For this question, we use the switching wells data.

```
wells <- read.table("http://www.stat.columbia.edu/~gelman/arm/examples/arsenic/wells.dat", header=TRUE)
wells_dt <- data.table(wells)

m5 <- glm(switch ~ (dist)+wells_dt$arsenic+wells_dt$educ, data=wells_dt, family=binomial(link="logit"))

xtilde <- cbind(1,wells_dt$dist, wells_dt$arsenic, wells_dt$educ)

n.sims <- 1000
sim.logit <- sim(m5,n.sims)
coef.logit <- coef(sim.logit)

n.tilde <- nrow(xtilde)
ytilde <- array(NA,c(n.sims,n.tilde))
for (s in 1:n.sims) {
  p.tilde <- invlogit(xtilde%%coef.logit[s,])
  ytilde[s,] <- rbinom(n.tilde,1,p.tilde)
}
```

```

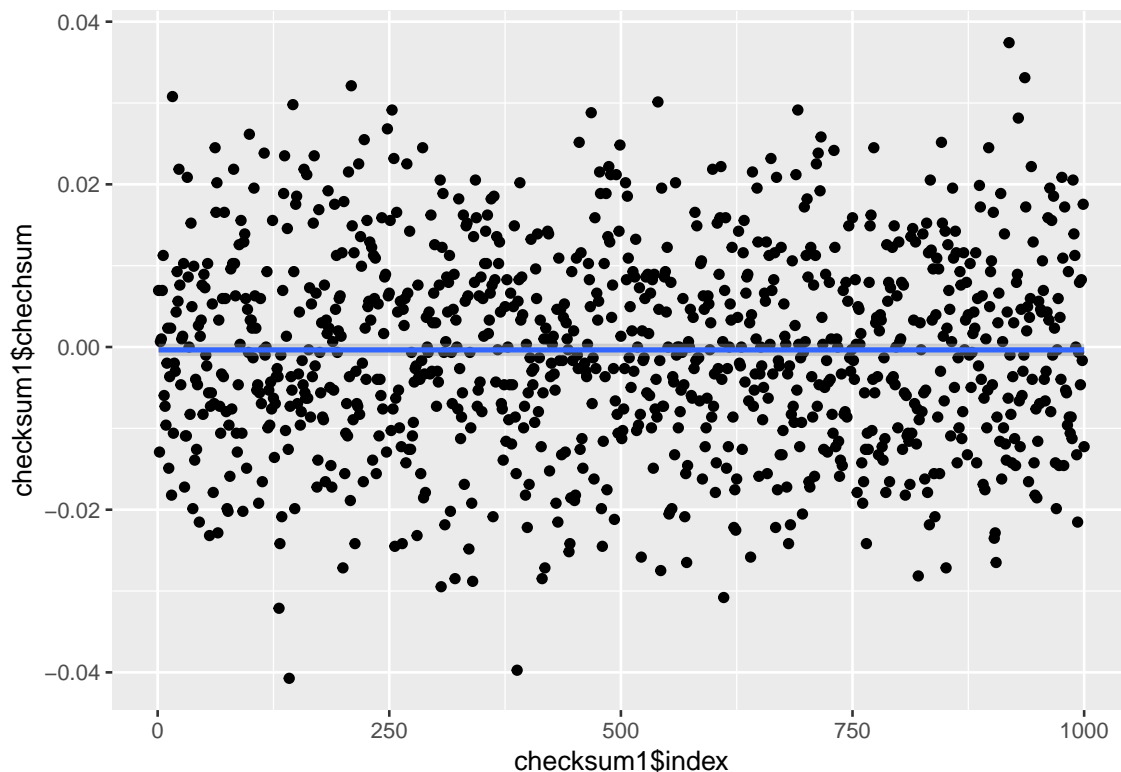
#plot, now let's check the goodness of fit of simulation
check <- array(NA,c(n.sims,n.tilde))
##For every loop, calculate the difference between simulation result and the real result.
for (s in 1:n.sims){
  check[s,] <- (t(ytilde[s,])-wells_dt$switch)
}

##calculate average difference of each y.
checksum <- rowSums(check)/3020

#draw the plot of difference. Actually it equals the average residual of each y.
checksum1 <- data.frame("index"=seq(1,1000),"checksum"=checksum)
ggplot(data=checksum1,mapping = aes(x=checksum1$index,y=checksum1$checksum))+
  geom_point()+
  geom_smooth()

## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'

```



So from the smooth line above, we see that the simulation fits very well.

## Repeat the previous exercise using a Poisson regression example.

We use the risk behavior data as our input or poisson model:

```

risky_behaviors<-read.dta("http://www.stat.columbia.edu/~gelman/arm/examples/risky.behavior/risky_behav
risky_behaviors$fupacts <- round(risky_behaviors$fupacts)

```

```

m6 <- glm(fupacts~women_alone+couples+bs_hiv,data=risky_behaviors,family = poisson)

xtilde <- cbind(1,risky_behaviors$women_alone,risky_behaviors$couples,risky_behaviors$bs_hiv)

n.sims <- 1000
sim.logit <- sim(m5,n.sims)
coef.logit <- coef(sim.logit)

n.tilde <- nrow(xtilde)
ytilde <- array(NA,c(n.sims,n.tilde))
for (s in 1:n.sims) {
  p.tilde <- invlogit(xtilde%%coef.logit[s,])
  ytilde[s,] <- rbinom(n.tilde,1,p.tilde)
}

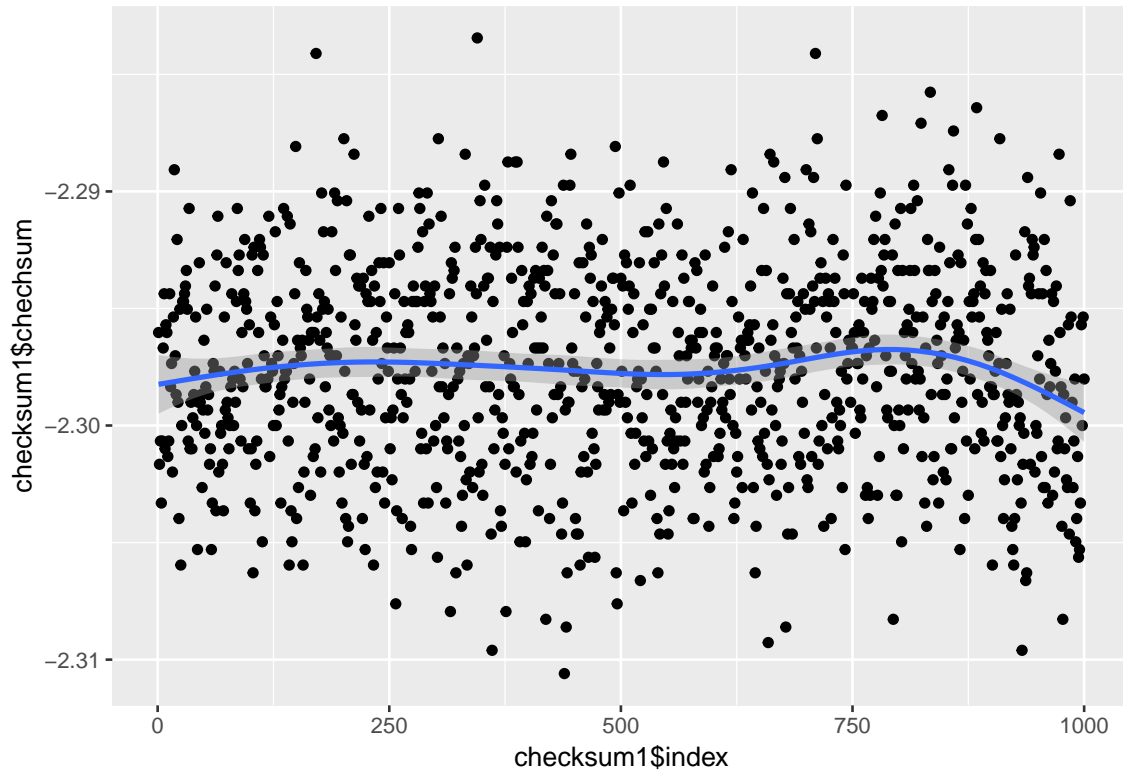
#plot, now let's check the goodness of fit of simulation
check <- array(NA,c(n.sims,n.tilde))
##For every loop, calculate the difference between simulation result and the real result.
for (s in 1:n.sims){
  check[s,] <- (t(ytilde[s,])-risky_behaviors$fupacts)
}

##calculate average difference of each y.
checksum <- rowSums(check)/3020

#draw the plot of difference. Actually it equals the average residual of each y.
checksum1 <- data.frame("index"=seq(1,1000),"chechsum"=checksum)
ggplot(data=checksum1,mapping = aes(x=checksum1$index,y=checksum1$chechsum))+
  geom_point()+
  geom_smooth()

## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'

```



## Inference for the ratio of parameters:

a (hypothetical) study compares the costs and effectiveness of two different medical treatments. - In the first part of the study, the difference in costs between treatments A and B is estimated at \$600 per patient, with a standard error of \$400, based on a regression with 50 degrees of freedom. - In the second part of the study, the difference in effectiveness is estimated at 3.0 (on some relevant measure), with a standard error of 1.0, based on a regression with 100 degrees of freedom. - For simplicity, assume that the data from the two parts of the study were collected independently.

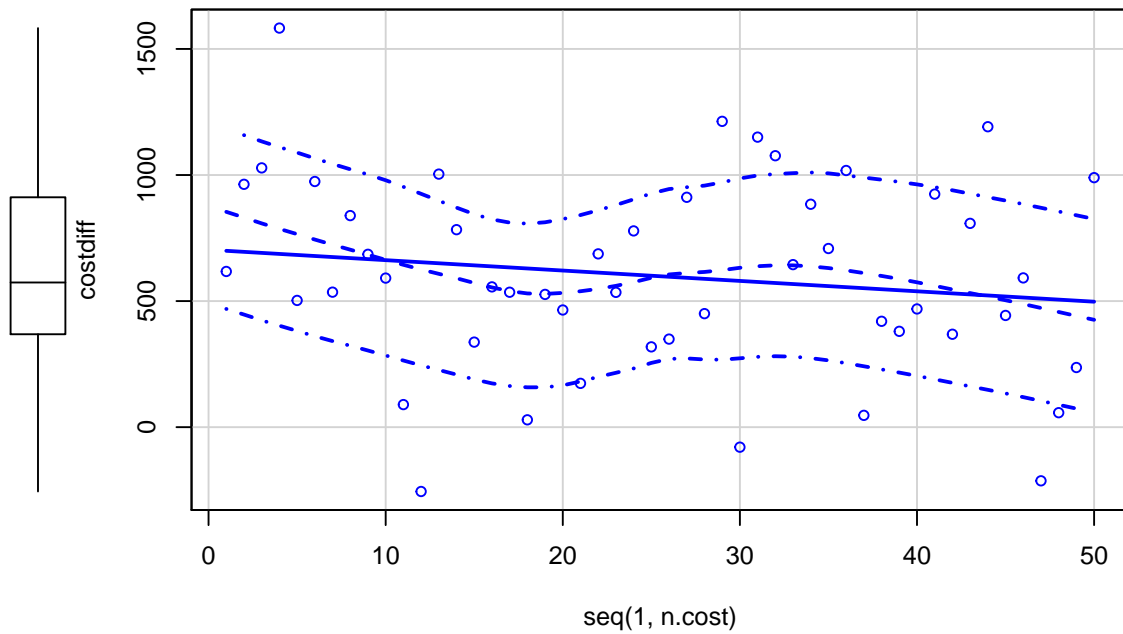
Inference is desired for the incremental cost-effectiveness ratio: the difference between the average costs of the two treatments, divided by the difference between their average effectiveness. (This problem is discussed further by Heitjan, Moskowitz, and Whang, 1999.)

1. Create 1000 simulation draws of the cost difference and the effectiveness difference, and make a scatterplot of these draws.

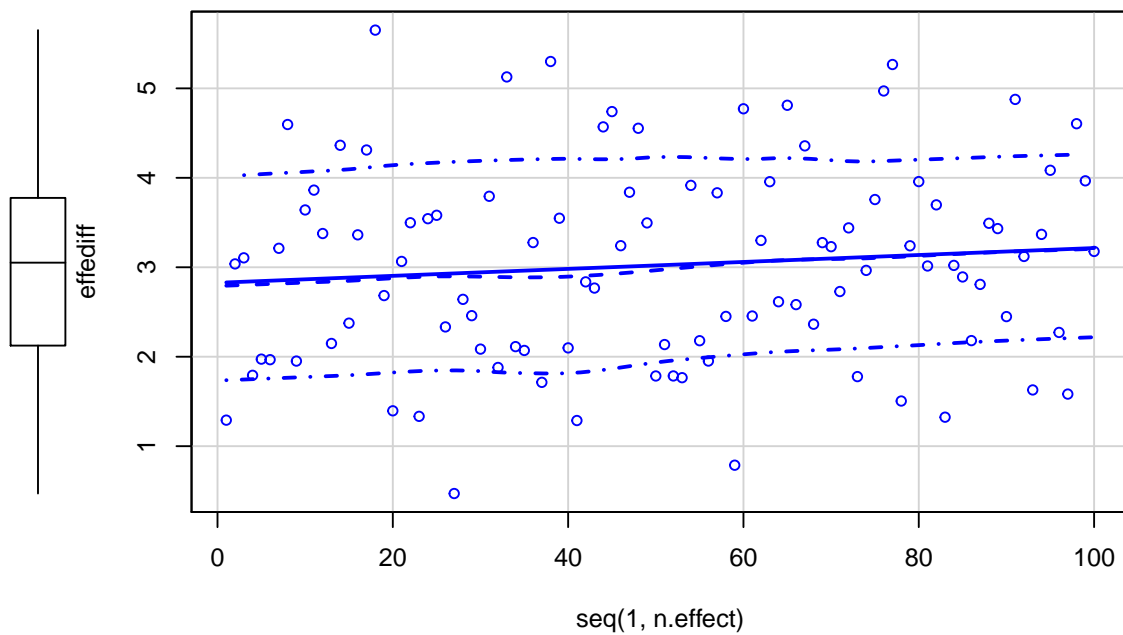
```
# cost difference
n.cost <- 50
costdiff <- rnorm(n.cost,600,sd=400)

#effectiveness difference
n.effect <- 100
effediff <- rnorm(n.effect,3,1)

scatterplot(x=seq(1,n.cost),y=costdiff)
```



```
scatterplot(x=seq(1,n.effect),y=effediff)
```



2. Use simulation to come up with an estimate, 50% interval, and 95% interval for the incremental cost-effectiveness ratio.

```
n.sims <- 1000
diffratio <- rep(NA,n.sims)

for (i in 1:n.sims ) {
  # cost difference
```

```

n.cost <- 50
costdiff <- rnorm(n.cost,600,sd=400)
costdiff <- mean(costdiff)
#effectiveness difference
n.effect <- 100
effediff <- rnorm(n.effect,3,1)
effediff <- mean(effediff)
diffratio[i] <- costdiff/effediff
}

```

```
quantile(diffratio,c(0.25,0.75))
```

```
##      25%      75%
## 187.2812 214.0647
```

```
quantile(diffratio,c(0.025,0.975))
```

```
##      2.5%     97.5%
## 162.3373 238.7889
```

3. Repeat this problem, changing the standard error on the difference in effectiveness to 2.0.

```

n.sims <- 1000
diffratio <- rep(NA,n.sims)

for (i in 1:n.sims ) {
  # cost difference
  n.cost <- 50
  costdiff <- rnorm(n.cost,600,sd=400)
  costdiff <- mean(costdiff)
  #effectiveness difference
  n.effect <- 100
  effediff <- rnorm(n.effect,3,2)
  effediff <- mean(effediff)
  diffratio[i] <- costdiff/effediff
}

```

```
quantile(diffratio,c(0.25,0.75))
```

```
##      25%      75%
## 185.3818 217.4919
```

```
quantile(diffratio,c(0.025,0.975))
```

```
##      2.5%     97.5%
## 160.6382 249.0866
```

## Predictive checks:

using data of interest to you, fit a model of interest. 1. Simulate replicated datasets and visually compare to the actual data.

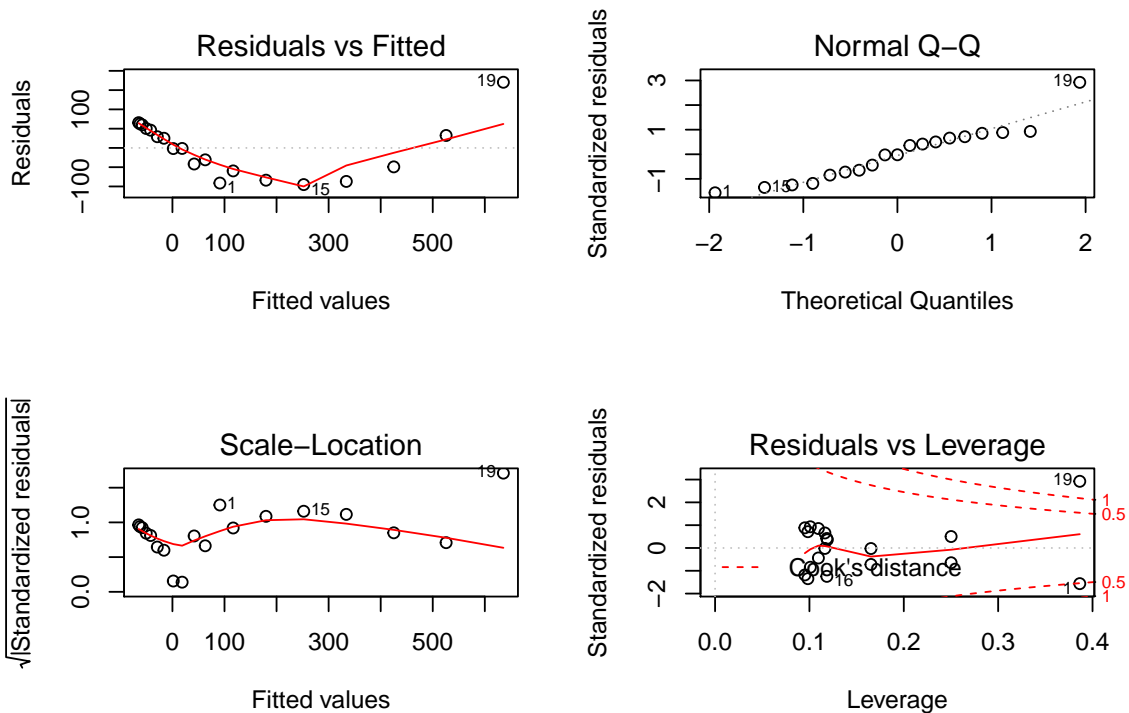
For this question, we use pressure data to do the simulation:

```
#Import data
data(pressure)

m7<- lm(formula = pressure ~ temperature+ I(temperature*temperature), data = pressure)
summary(m7)

##
## Call:
## lm(formula = pressure ~ temperature + I(temperature * temperature),
##     data = pressure)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -95.142 -54.391  -1.353   48.238 170.374
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    91.154379   46.262513     1.970  0.066354 .
## temperature    -2.706167    0.595775    -4.542  0.000333 ***
## I(temperature * temperature)  0.011718    0.001597     7.336 1.67e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 74.42 on 16 degrees of freedom
## Multiple R-squared:  0.9024, Adjusted R-squared:  0.8902
## F-statistic:    74 on 2 and 16 DF,  p-value: 8.209e-09

par(mfrow=c(2,2))
plot(m7)
```



```

sim.m7 <- sim(m7, length(pressure$pressure))

c(mean(sim.m7@coef[, 1]), sd(sim.m7@coef[, 1]))

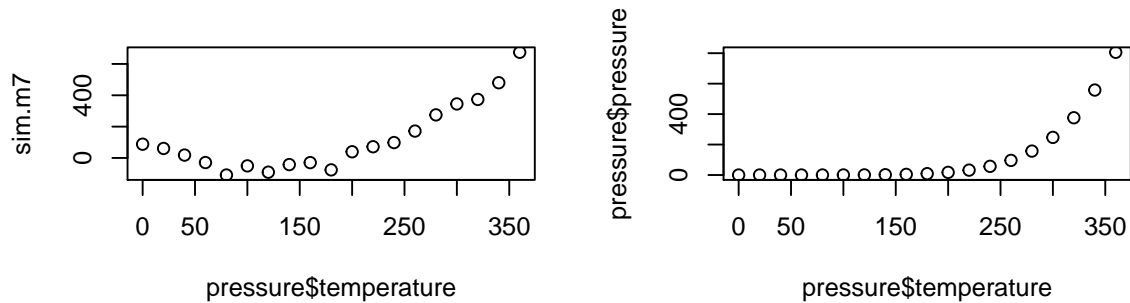
## [1] 99.26634 48.19870

c(mean(sim.m7@coef[, 2]), sd(sim.m7@coef[, 2]))

## [1] -2.8600853 0.7051724

sim.m7 <- (sim.m7@coef[,3]) * pressure$temperature*pressure$temperature + sim.m7@coef[,1]+(sim.m7@coef[
plot(pressure$temperature, sim.m7)
plot(pressure$temperature, pressure$pressure)

```



2. Summarize the data by a numerical test statistic, and compare to the values of the test statistic in the replicated datasets.

```

summary(pressure$pressure)

##      Min.   1st Qu.   Median     Mean  3rd Qu.     Max.
##  0.0002   0.1800    8.8000 124.3367 126.5000 806.0000

summary(sim.m7)

##      Min.   1st Qu.   Median     Mean  3rd Qu.     Max.
## -108.82  -35.99    60.64  119.37  223.26  674.51

```

## (optional) Propagation of uncertainty:

we use a highly idealized setting to illustrate the use of simulations in combining uncertainties. Suppose a company changes its technology for widget production, and a study estimates the cost savings at \$5 per unit, but with a standard error of \$4. Furthermore, a forecast estimates the size of the market (that is, the number of widgets that will be sold) at 40,000, with a standard error of 10,000. Assuming these two sources of uncertainty are independent, use simulation to estimate the total amount of money saved by the new product (that is, savings per unit, multiplied by size of the market).

## (optional) Fitting the wrong model:

suppose you have 100 data points that arose from the following model:  $y = 3 + 0.1x_1 + 0.5x_2 + \text{error}$ , with errors having a t distribution with mean 0, scale 5, and 4 degrees of freedom. We shall explore the implications of fitting a standard linear regression to these data.



1. Simulate data from this model. For simplicity, suppose the values of `x_1` are simply the integers from 1 to 100, and that the values of `x_2` are random and equally likely to be 0 or 1. In R, you can define `x_1 <- 1:100`, simulate `x_2` using `rbinom()`, then create the linear predictor, and finally simulate the random errors in `y` using the `rt()` function. Fit a linear regression (with normal errors) to these data and see if the 68% confidence intervals for the regression coefficients (for each, the estimates  $\pm 1$  standard error) cover the true values.
2. Put the above step in a loop and repeat 1000 times. Calculate the confidence coverage for the 68% intervals for each of the three coefficients in the model.
3. Repeat this simulation, but instead fit the model using `t` errors (use `hett::tlm`).

## (optional) Using simulation to check the fit of a time-series model:

find time-series data and fit a first-order autoregression model to it. Then use predictive simulation to check the fit of this model as in GH Section 8.4.

## (optional) Model checking for count data:

the folder `risky.behavior` contains data from a study of behavior of couples at risk for HIV;

“sex” is a factor variable with labels “woman” and “man”. This is the member of the couple that reporting sex acts to the researcher

The variables “couple” and “women\_alone” code the intervention:

couple women\_alone 0 0 control - no counselling 1 0 the couple was counselled together 0 1 only the woman was counselled

“bs\_hiv” indicates whether the member reporting sex acts was HIV-positive at “baseline”, that is, at the beginning of the study.

“bupacts” - number of unprotected sex acts reported at “baseline”, that is, at the beginning of the study

“fupacts” - number of unprotected sex acts reported at the end of the study (final report).

1. Fit a Poisson regression model predicting number of unprotected sex acts from baseline HIV status. Perform predictive simulation to generate 1000 datasets and record both the percent of observations that are equal to 0 and the percent that are greater than 10 (the third quartile in the observed data) for each. Compare these values to the observed value in the original data.
2. Repeat (1) using an overdispersed Poisson regression model.

```
# afunction to generate from quasi poisson
rqpois = function(n, lambda, phi) {
  mu = lambda
  k = mu/phi/(1-1/phi)
  return(rnbinom(n, mu = mu, size = k))
}
# https://www.r-bloggers.com/generating-a-quasi-poisson-distribution-version-2/
```

3. Repeat (2), also including gender and baseline number of unprotected sex acts as input variables.