# Final Paper of ADA

BI Norwegian Business School — December 7, 2022

## Assignment 1

### A

### I

| Variable Name | Num of Missing Values |
|---|---|
| PIN | 0 |
| CIN | 0 |
| continent | 0 |
| country | 0 |
| city | 166 |
| countrycode | 641 |
| partnum | 20 |
| partcode | 15 |
| sample | 7795 |
| sex | 0 |
| age | 0 |
| religious | 64 |
| religion | 497 |
| relstat | 171 |
| relstat2 | 171 |
| rellength | 7072 |
| ideal_intelligence | 36 |
| ideal_kindness | 45 |
| ideal_health | 67 |
| ideal_physatt | 351 |
| ideal_resources | 401 |
| mate_age | 5479 |
| popsize | 0 |
| country_religion | 0 |

| | |
|---|---|
| lattitude | 0 |
| gem1995 | 3302 |
| gdi1995 | 1842 |
| gii | 281 |
| gdi2015 | 251 |
| gggi | 0 |
| gdp_percap | 0 |
| infect_death | 0 |
| infect_yll | 0 |
| cmc_yll | 0 |
| gb_path | 9453 |

```
1 data <- read.csv("ReplicationProcessedfinaldata04202018.csv")
2 miss <- data.frame(columns=colnames(data),missing.num=NA)
3 for (column in miss$columns){
4   miss[miss$columns==column,2] <- sum(is.na(data[,column]))
5 }
6 rows.before <- nrow(data)
7 miss
8 # here we didn't use the commands given in the hints because we
      think directly creating a dataframe would be more convenient.
```

## II

We have removed 40.1764% of the total observations.

```
1 data <- data[-c(which((data$mate_age<12)|(is.na(data$mate_age))))
    ,]
2 rows.after <- nrow(data)
3 proportion.removed <- (rows.before-rows.after)/rows.before
4 proportion.removed #0.401764
```

## III

| Country Name | Num of Observations |
|---|---|
| Hungary | 839 |
| Pakistan | 474 |
| Poland | 380 |
| Slovenia | 476 |

```r
countries <- data.frame(country.name=unique(data$country),
    nCountries=NA)
for (name in countries$country.name) {
  countries[countries$country.name==name,2]=sum(data$country==name
    )
}
countries <- countries[-c(which(countries$nCountries<350)),]
```
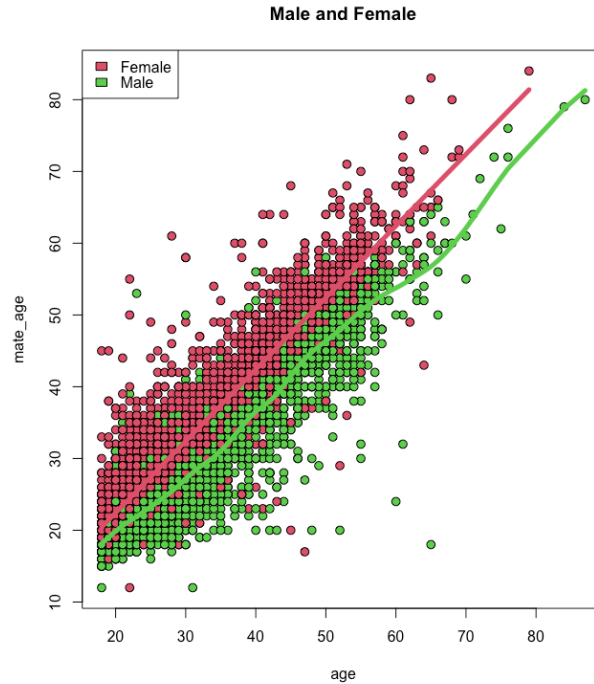
**B**

**I**



Figure 1: Scatter Plot of Male and Female

After testing the fitness of simple, quadratic and cubic linear models for both male and female groups. We identify the model as shown in equation 1 below which implements simple linear for female and quadratic linear for male. Averagely speaking, the male tend to find a mate around 3 years younger and females' mates tend to be approximate 2.5 years older. Looking at the plot for females, the regression line shoots above the line of mate_age=age. It indicates that females tend to choose older mates than their age. In contrast, for males, the opposite happens, which indicates that males tend to prefer younger partners. All the p-values are statistically significant, meaning that we should include them all in our model.
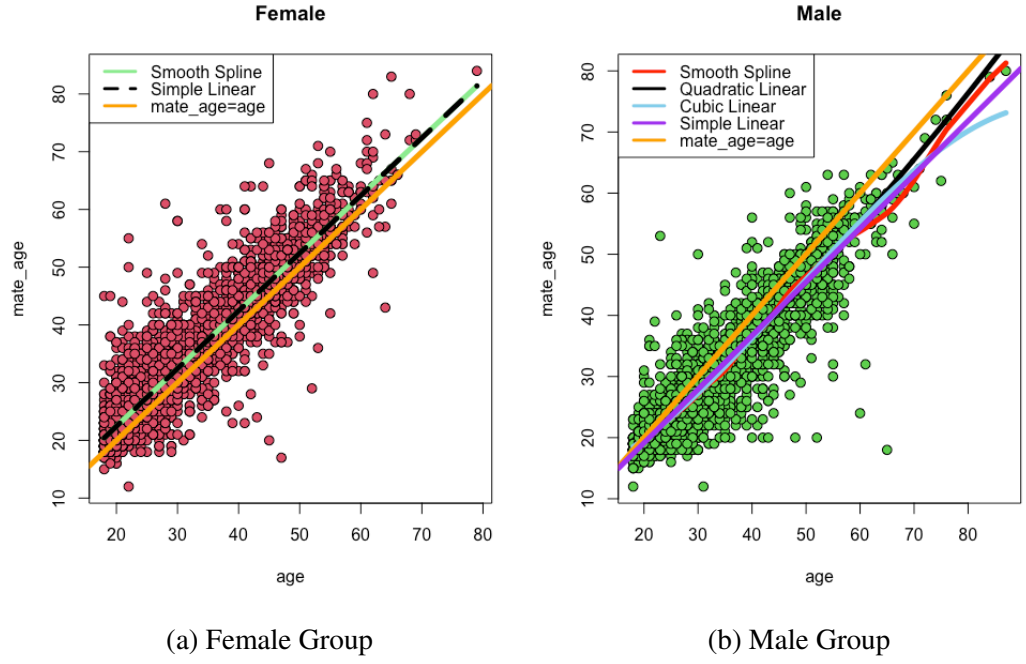
(a) Female Group

(b) Male Group

Figure 2: Two Groups

$$mate\_age = \beta_0 + \beta_1 age + \beta_2 sex + \beta_3 age \times sex + \beta_4 sex \times age^2 \tag{1}$$

Table 3: Regression Outcomes

| | *Dependent variable:* | | | | | |
|---|---|---|---|---|---|---|
| | Mate_Age | | | Mate_Age | | Mate_Age |
| | *Female* | | | *Male* | | *Both* |
| | lmFe | lmFeQu | lmMe | lmMeQu | lmMeCu | lmJoin |
| age | 0.999*** | 1.027*** | 0.878*** | 0.678*** | 0.017 | 0.999*** |
| | (0.005) | (0.036) | (0.005) | (0.032) | (0.110) | (0.005) |
| $age^2$ | | $-0.0004$ | | 0.003*** | 0.020*** | |
| | | (0.001) | | (0.0004) | (0.003) | |
| $age^3$ | | | | | $-0.0001^{***}$ | |
| | | | | | (0.00002) | |
| sex | | | | | | 2.389*** |
| | | | | | | (0.603) |
| age×sex | | | | | | $-0.321^{***}$ |

4

(0.034)

|  | | | | | | |
|---|---|---|---|---|---|---|
| sex$\times$age$^2$ | | | | | | 0.003*** |
| | | | | | | (0.0005) |
| 5 | | | | | | |
| Constant | 2.449*** | 2.018*** | 1.576*** | 4.839*** | 12.527*** | 2.449*** |
| | (0.175) | (0.584) | (0.181) | (0.553) | (1.336) | (0.169) |
| Observations | 4,752 | 4,752 | 3,862 | 3,862 | 3,862 | 8,614 |
| $R^2$ | 0.876 | 0.876 | 0.872 | 0.873 | 0.875 | 0.878 |
| Adjusted $R^2$ | 0.876 | 0.876 | 0.872 | 0.873 | 0.874 | 0.878 |
| Residual SE | 4.160 | 4.161 | 3.856 | 3.837 | 3.818 | 4.019 |

*Note:* $^*$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01
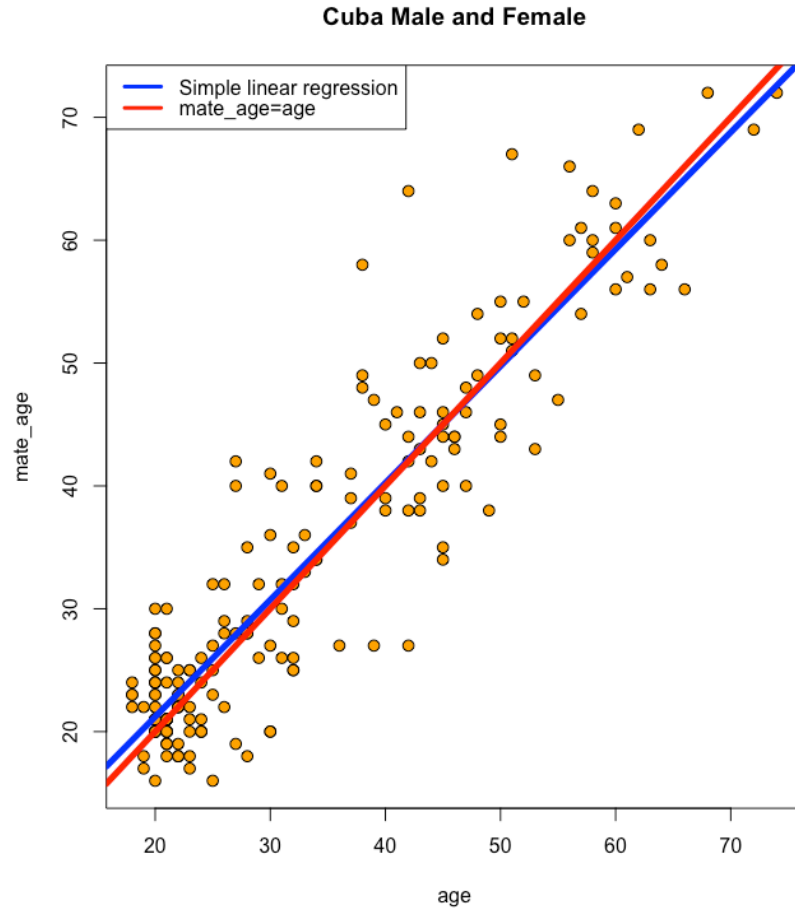
**Cuba Male and Female**



Figure 3: Scatter Plot of Cuba Male and Female

The F-test suggests for the Cuba dataset, without considering the term gender, we have no evidence to reject the null hypothesis which means there is no mate age preference for each gender.

```
1 Analysis of Variance Table
2
3 Model 1: mate_age ~ age
4 Model 2: mate_age ~ offset(1 * age) - 1
5   Res.Df  RSS Df Sum of Sq     F Pr(>F)
6 1    186 5850
7 2    188 5992 -2   -142.03 2.258 0.1074
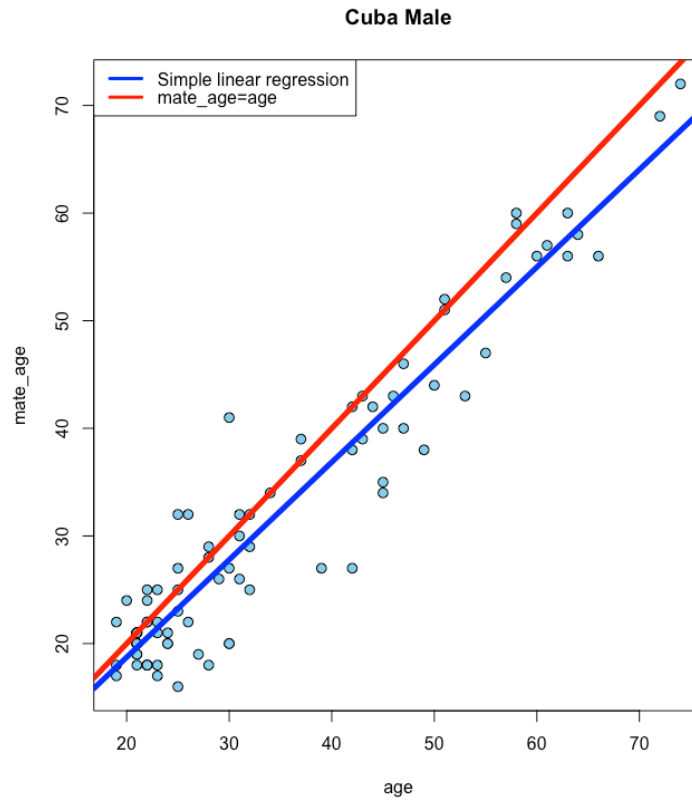```

Listing 1: Anova Analysis of Both Gender Dataset

Figure 4: Scatter Plot of Cuba Male

The F-test is significant for the male data group. Combined with Figure 4 we have the evidence to conclude that males tend to prefer younger mates on average.

```
1 Analysis of Variance Table
2
3 Model 1: mate_age ~ age
4 Model 2: mate_age ~ offset(1 * age) - 1
5   Res.Df    RSS Df Sum of Sq      F     Pr(>F)
6 1     86 1524.8
7 2     88 2292.0 -2   -767.15 21.633 2.452e-08 ***
```
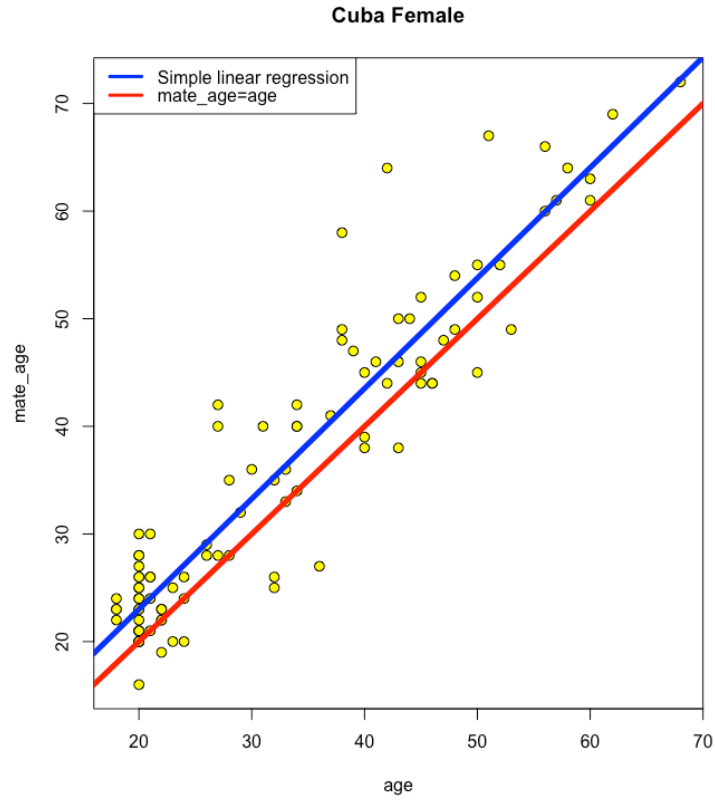
Listing 2: Anova Analysis of Male Dataset

Figure 5: Scatter Plot of Cuba Female

The F-test for female group is also significant, meaning on average, females tend to prefer older mates.

```
Analysis of Variance Table

Model 1: mate_age ~ age
Model 2: mate_age ~ offset(1 * age) - 1
  Res.Df    RSS Df Sum of Sq     F    Pr(>F)
1     98 2573.1
2    100 3700.0 -2   -1126.9 21.46 1.863e-08 ***
```

Listing 3: Anova Analysis of Female Dataset

The fact that male and female have opposite age preference with regard to choosing mate could justify the insignificance in Listing 1 since the two opposite age preference could be offset when male and female merge as one large group
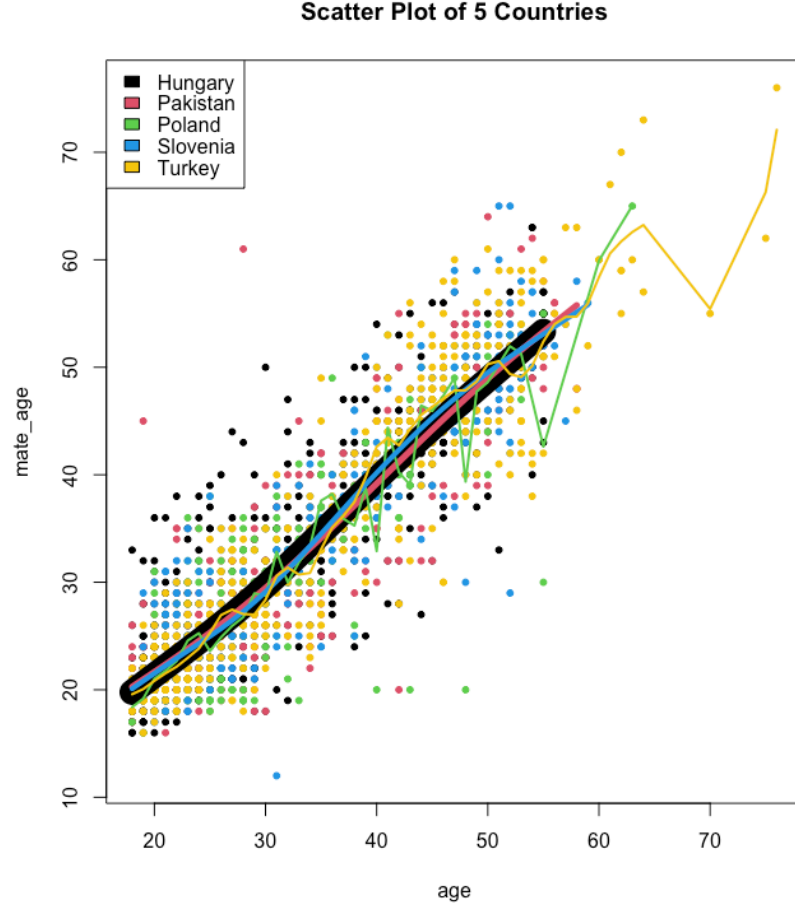
**Scatter Plot of 5 Countries**



Figure 6: Scatter Plot and Smooth Splines

The scatterplot suggests a highly overlapped common trend for Hungary, Pakistan, and Slovenia. As a result, we treat these countries as base case. For Turkey and Poland, we introduce two indicator variables. From our analysis, one can draw that the best fit models for base case and Poland are both simple linear models, whereas a cubic linear model for Turkey. The term *Poland* would be no longer significant after introducing its interaction term with *age*. The result of F-test between models before and after introducing interaction term is insignificant indicating that the interaction term is unnecessary. Hence we remove the interaction term of *Poland*, and we end up with the final model as shown in equation 2 and regression result shown in table 5 (lmJoint2) which is statistically more plausible.
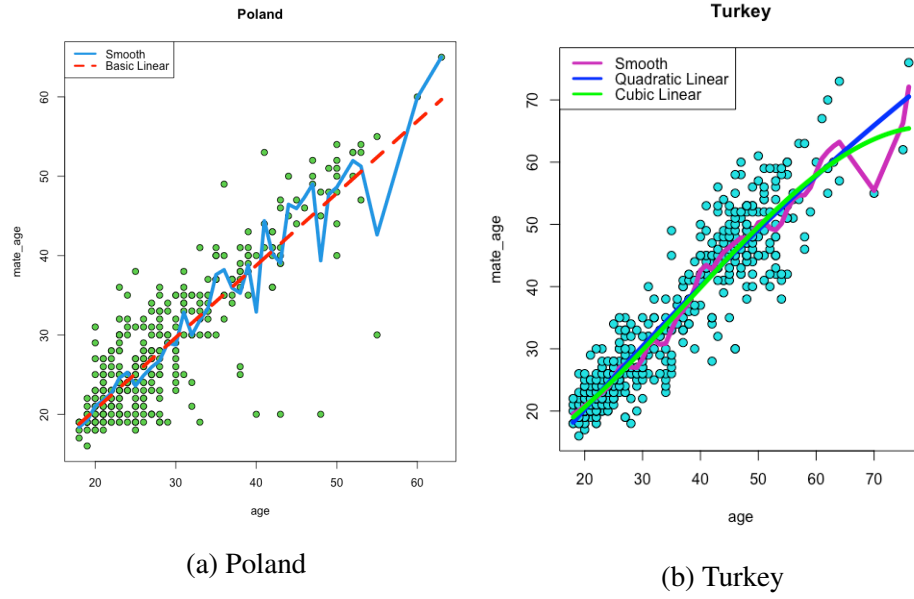
(a) Poland

(b) Turkey

Figure 7: Scatter plot seperately

Table 4

| | Dependent variable: | | | | | |
|---|---|---|---|---|---|---|
| | Mate_Age | Mate_Age | | Mate_Age | | |
| | *ThreeCountry* | *Poland* | | *Turkey* | | |
| | lmThree | lmPo | lmPoQu | lmTu | lmTuQu | lmTuCu |
| age | 0.933*** | 0.909*** | 0.822*** | 0.952*** | 1.138*** | 0.200 |
| | (0.008) | (0.026) | (0.170) | (0.015) | (0.093) | (0.367) |
| age$^2$ | | | 0.001 | | −0.003** | 0.022** |
| | | | (0.002) | | (0.001) | (0.009) |
| age$^3$ | | | | | | −0.0002*** |
| | | | | | | (0.0001) |
| Constant | 2.170*** | 2.405*** | 3.744 | 1.539*** | −1.485 | 9.519** |
| | (0.254) | (0.793) | (2.726) | (0.517) | (1.568) | (4.452) |
| Observations | 2,733 | 380 | 380 | 564 | 564 | 564 |
| R$^2$ | 0.839 | 0.758 | 0.758 | 0.883 | 0.884 | 0.886 |
| Adjusted R$^2$ | 0.839 | 0.757 | 0.756 | 0.883 | 0.884 | 0.885 |
| Residual SE | 4.52 | 4.59 | 4.60 | 4.40 | 4.39 | 4.36 |

$$
\begin{aligned}
mate\_age =& \beta_0 + \beta_1 age + \beta_2 Poland + \beta_3 Turkey \\
& + \beta_4 Turkey \times age^2 + \beta_5 Turkey \times age^3
\end{aligned}
\tag{2}
$$

```
1 Analysis of Variance Table
2
3 Model 1: mate_age ~ age + Poland + Poland * age + Turkey + Turkey
     * I(age^2) +
4    Turkey * I(age^3) - I(age^3) - I(age^2)
5 Model 2: mate_age ~ age + Poland + Turkey + Turkey * I(age^2) +
     Turkey *
6    I(age^3) - I(age^3) - I(age^2)
7  Res.Df   RSS Df Sum of Sq      F Pr(>F)
8 1   2726 55520
9 2   2727 55528 -1   -8.5642 0.4205 0.5167
```

Listing 4: Anova Analysis of Models with and without Interaction Term

Table 5

|  | *Dependent variable:* | |
|---|---|---|
|  | mate_age | |
|  | lmJoint1 | lmJoint2 |
| age | 0.927*** | 0.925*** |
|  | (0.010) | (0.009) |
| Poland | −0.058 | −0.577** |
|  | (0.841) | (0.255) |
| Turkey | −1.692*** | −1.734*** |
|  | (0.588) | (0.584) |
| age×Poland | −0.018 | |
|  | (0.028) | |
| Turkey×age$^2$ | 0.004*** | 0.004*** |
|  | (0.001) | (0.001) |
| Turkey×age$^3$ | −0.0001*** | −0.0001*** |
|  | (0.00002) | (0.00002) |

11

| | | |
|---|---|---|
| Constant | 2.462*** | 2.531*** |
| | (0.318) | (0.299) |
| | | |
| Observations | 2,733 | 2,733 |
| $R^2$ | 0.840 | 0.840 |
| Adjusted $R^2$ | 0.840 | 0.840 |
| Residual Std. Error | 4.513 | 4.512 |

*Note:*  $^*$p<0.1; $^{**}$p<0.05; $^{***}$p<0.01

# Assignment 2

## A

We already know that:

$$x = a + by \tag{3}$$
$$y = c + dx \tag{4}$$

Now we use equation 4 to substitute the y in equation 3, so we can get the following equation:

$$x = a + b(c + dx)$$
$$x = a + bc + bdx$$
$$(1 - bd)x = a + bc$$

Since $bd \neq 1$, we can perform the following transformation:

$$x = \frac{a + bc}{1 - bd} \tag{5}$$

In the same way, we can compute the y as follows:

$$y = c + d(a + by)$$
$$y = c + ad + bdy$$
$$(1 - bd)y = c + ad$$
$$y = \frac{c + ad}{1 - bd}$$
$$y = \frac{c + ad + bdc - bdc}{1 - bd}$$
$$y = \frac{c(1 - bd) + d(a + bc)}{1 - bd}$$
$$y = c + d\frac{a + bc}{1 - bd} \tag{6}$$

## B

From the text we know that if $i = 2j$ for an integer $j$ which means $i$ is an even integer, then:

$$\mathcal{I} = \{i - 1\}$$

And if $i = 2j - 1$ which means i is an odd integer, then:

$$\mathcal{I} = \{i + 1\}$$

Thus, if we want to prove $\mathcal{I} = \{i - (-1)^i\}$, we just need to prove $\mathcal{I} = \{i - (-1)^{2j}\}$ for $i$ being even and $\mathcal{I} = \{i - (-1)^{2j-1}\}$ for $i$ being odd. Since we know that $(-1)^2 = 1$ and $(-1)^{2j} = [(-1)^2]^j$ and $(-1)^{2j-1} = [(-1)^2]^j/(-1)$, we can conclude that:

$$\mathcal{I} = \{i - (-1)^{2j}\}$$
$$= \{i - 1\} \tag{7}$$
$$\mathcal{I} = \{i - (-1)^{2j-1}\}$$
$$= \{i + 1\} \tag{8}$$

Equation 7 and 8 confirms the correctness of $\mathcal{I} = \{i - (-1)^i\}$.

## C

According to the network model and the equation we proved in problem B, we can draw the conclusion that $\mid \mathcal{I}(i) = 1 \mid$ always stands. And by using $i - (-1)^i$ to substitute the $k$, we can draw the following equation:

$$y_i = \beta_0 + \beta_1\left(\frac{1}{\mid \mathcal{I}(i) \mid} \sum_{k \in \mathcal{I}(i)} y_k\right) + \mu_i$$
$$= \beta_0 + \beta_1 y_{i-(-1)^i} + \mu_i \tag{9}$$

To be able to use equation 5 and 6, we should use $i - (-1)^i$ to substitute $i$ to create a new equation as follows:

$$y_{i-(-1)^i} = \beta_0 + \beta_1 y_{i-(-1)^i-(-1)^{i-(-1)^i}} + \mu_{i-(-1)^i} \tag{10}$$

Since the following two scenarios:

1. If i an odd integer, then $(-1)^i = -1$ and $i - (-1)^i$ is even which leads to the result that $(-1)^{i-(-1)^i}$ being 1 and $i - (-1)^i - (-1)^{i-(-1)^i}$ being $i$ because the second term and third term are offset with each other.

2. If i an even integer, then $(-1)^i = 1$ and $i - (-1)^i$ is odd which leads to the result that $(-1)^{i-(-1)^i}$ being -1 and $i - (-1)^i - (-1)^{i-(-1)^i}$ being $i$ because the second term and third term are offset with each other.

equation 10 can also be written as:

$$y_{i-(-1)^i} = \beta_0 + \beta_1 y_i + \mu_{i-(-1)^i} \tag{11}$$

Use equation 11 to substitute $y_{i-(-1)^i}$ in equation 9 we can get the following Equation:

$$
\begin{aligned}
y_i =& \beta_0 + \beta_1 y_{i-(-1)^i} + \mu_i \\
=& \beta_0 + \beta_1(\beta_0 + \beta_1 y_i + \mu_{i-(-1)^i}) + \mu_i \\
=& \beta_0 + \beta_0\beta_1 + \beta_1^2 y_i + \beta_1\mu_{i-(-1)^i} + \mu_i \\
(1-\beta_1^2)y_i =& \beta_0 + \beta_0\beta_1 + \beta_1\mu_{i-(-1)^i} + \mu_i \\
(1+\beta_1)(1-\beta_1)y_i =& \beta_0(1+\beta_1) + \beta_1\mu_{i-(-1)^i} + \mu_i \\
y_i =& \frac{\beta_0}{1-\beta_1} + \frac{1}{1-\beta_1^2}(\mu_i + \beta_1\mu_{i-(-1)^i}) \quad\quad (12)
\end{aligned}
$$

# D

## I

Using equation $y_i = \beta_0 + \beta_1 y_{i-(-1)^i} + \mu_i$ to simulate $y_1, y_2...y_n$ is nearly impossible since if we want to calculate $y_1$, we need to know $y_2$. If we would like to know $y_2$ then we need to know $y_1$, which is a dead loop.

## II

- Line 1-6: Set the basic parameters for the afterward simulation.

- Line 7: Build the index vector so we can access both $\mu_i$ and $\mu_{i-(-1)^i}$ in each calculation.

- Line 8: Simulate 200 random normal numbers as the error terms.

- Line 9: "beta0/(1-beta1)" corresponds to the calculation of $\frac{\beta_0}{1-\beta_1}$. Using the iSeq as index we can assess both $\mu_i$ and $\mu_{i-(-1)^i}$ and "u[iSeq]+beta1*u[iSeq-(-1)^iSeq]" corresponds to $\mu_i + \beta_1\mu_{i-(-1)^i}$. "(1-beta1^2)" corresponds to $\frac{1}{1-\beta_1^2}$.

## III

```
check <- NULL
y.calculated <- NULL
for (i in (1:n)) {
  check[i] <- (beta0+beta1*y[i-(-1)^i]+u[i]-y[i])<1e-10
  y_calculated[i] <- beta0+beta1*y[i-(-1)^i]+u[i]
}
sum(check)==n #True
plot(y.calculated,y,cex=1.5,pch=21,main = "Simulated y and
    calculated y",xlab="Y Calculated",ylab="Y Simulated")
abline(a=0,b=1,col="red",lwd=4)
# which verifies that every y satisfies the equation
```
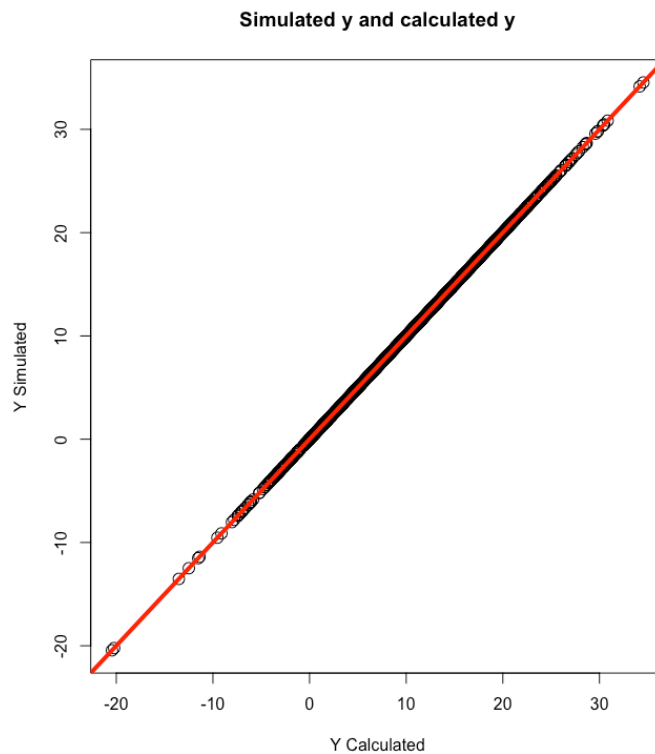
15

Figure 8: Scatter Plot of Cuba Male

## E

## I

Simulating 1000 times the sample mean of y and comparing it with $\mu$ shows that the sample mean of y is exactly the same as the population mean mu.

```
numRep=1000
m <- 100
n <- 2*m
beta0 = 1
beta1 = 0.9
sigma.u = 1
mu = beta0/(1-beta1)
mean.y.collect <- NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  mean.y.collect[i] <- mean(y)}
hist(mean.y.collect,main = "Histogram of y means",col = "white")
abline(v=mu,col="red",lwd=4)
```
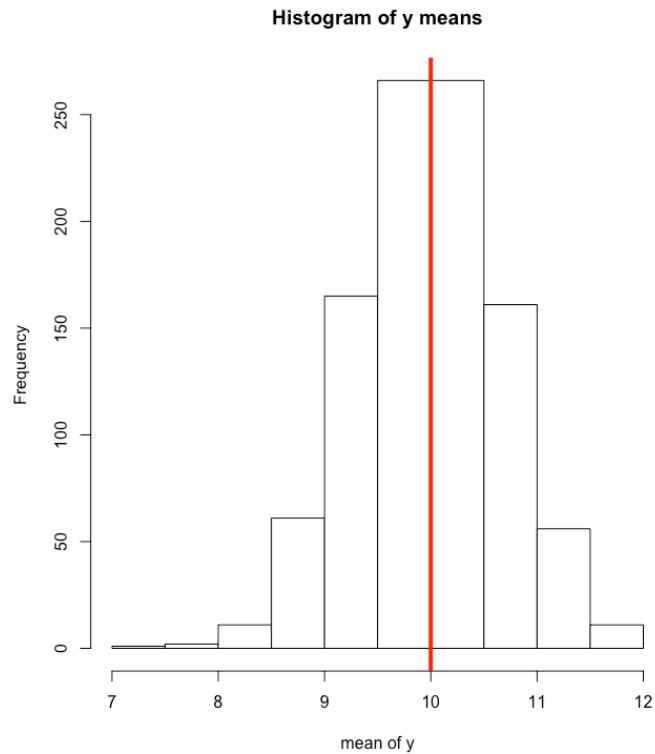
Figure 9: Histogram of y means

## II

Simulating with 1000 repetitions, to approximate the coverage rate of the confidence interval, indicates a 83.7% coverage rate, which is indeed not close to 95%.

```
numRep=1000
m <- 100
n <- 2*m
beta0 = 1
beta1 = 0.9
sigma.u = 1
counter <- NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  lo.lim <- mean(y) - 1.96*sqrt(var(y))/sqrt(n)
  up.lim <- mean(y) + 1.96*sqrt(var(y))/sqrt(n)
  counter[i] <- (mu>lo.lim) & (mu<up.lim)
}
sum(counter)/length(counter) #0.837
```

# Assignment 3

## A

### I

Let's first compute $\bar{x}_n$ first:

$$
\begin{aligned}
\bar{x}_n =& \frac{1}{n} \sum_{i=1}^{n} x_i \\
=& \frac{1}{n} \sum_{i=1}^{n} (\mu + \frac{1}{1 - \beta_1} \mu_i) \\
=& \frac{1}{n} \sum_{i=1}^{n} \mu + \frac{1}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^{n} \mu_i \\
=& \mu + \frac{1}{1 - \beta_1} \bar{\mu}_i
\end{aligned}
\tag{13}
$$

Now, let's look at the $\bar{y}_n$. As the equation of $y_i$ shown in equation 12, so the $\bar{y}_n$ could be formulated as follows:

$$
\begin{aligned}
\bar{y}_n =& \frac{1}{n} \sum_{i=1}^{n} (\frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1^2} (\mu_i + \beta_1 \mu_{i-(-1)^i})) \\
=& \frac{1}{n} (\sum_{i=1}^{m} (\frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1^2} (\mu_{2i} + \beta_1 \mu_{2i-(-1)^{2i}})) \\
& + \sum_{i=1}^{m} (\frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1^2} (\mu_{2i-1} + \beta_1 \mu_{2i-1-(-1)^{2i-1}}))) \\
=& \frac{1}{n} (\sum_{i=1}^{m} (\frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1^2} (\mu_{2i} + \beta_1 \mu_{2i-1})) \\
& + \sum_{i=1}^{m} (\frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1^2} (\mu_{2i-1} + \beta_1 \mu_{2i}))) \\
=& \frac{1}{n} ((2m \frac{\beta_0}{1 - \beta_1}) + \frac{1}{1 - \beta_1^2} \sum_{i=1}^{m} (\mu_{2i} + \beta_1 \mu_{2i-1} + \mu_{2i-1} + \beta_1 \mu_{2i})) \\
=& \frac{1}{n} ((n \frac{\beta_0}{1 - \beta_1}) + \frac{1}{1 - \beta_1^2} \sum_{i=1}^{m} (\mu_{2i}(1 + \beta_1) + \mu_{2i-1}(1 + \beta_1))) \\
=& \frac{\beta_0}{1 - \beta_1} + \frac{1}{n} \frac{1 + \beta_1}{(1 - \beta_1)(1 + \beta_1)} \sum_{i=1}^{m} (\mu_{2i} + \mu_{2i-1}) \\
=& \frac{\beta_0}{1 - \beta_1} + \frac{1}{n} \frac{1}{1 - \beta_1} \sum_{i=1}^{n} \mu_i \\
=& \frac{\beta_0}{1 - \beta_1} + \frac{1}{n} \frac{n}{1 - \beta_1} \frac{\sum_{i=1}^{n} \mu_i}{n} \\
=& \frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1} \bar{\mu}_i
\end{aligned}
$$

Since we already know that $\mu = \frac{\beta_0}{1-\beta_1}$, thus the final equation of $\bar{y}_n$ can be written as:

$$\bar{y}_n = \mu + \frac{1}{1-\beta_1}\bar{\mu}_i \tag{14}$$

By comparing equation 13 and 14 we can draw the conclusion that $\bar{x}_n = \bar{y}_n$.

## II

In order to determine whether the means of y and the mean of x are equal. We simulate 1000 times whether the difference between the mean of x and the mean of y is less than 1e-10 (an extremely small number). Comparing the True and False values with our number of simulations gives us positive feedback. The difference between the mean of x and the mean of y is less than 1e-10 at 100% of our simulation. Thus, we can conclude that the mean of x and the mean of y is approximately equal.

```
1  m <- 1000
2  n <- 2*m
3  numRep <- 1000
4  beta0 <- 1
5  beta1 <- 0.9
6  sigma.u <- 1
7  mu <- beta0/(1-beta1)
8  xy.collect <- NULL
9  for (i in (1:numRep)) {
10    iSeq <- (1:n)
11    u <- rnorm(n, mean=0, sd=sigma.u)
12    y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
        beta1^2)
13    x <- mu + 1/(1-beta1)*u
14    xy.collect[i] <- (mean(x)-mean(y))<1e-10
15  }
16  sum(xy.collect) #1000 which equals to the length of xy.collect, so
        for all the 1000 simulations, mean(x)=mean(y)
```

## B

## I

Looking at the plot (Figure 10), one can see that the density curve precisely follows the histogram, which justifies that the equation holds.

```
1  m <- 1000
2  n <- 2*m
3  numRep <- 5000
4  beta0 <- 1
5  beta1 <- 0.9
6  sigma.u <- 1
7  mu <- beta0/(1-beta1)
```

```
8  mean_collect=NULL
9  for (i in (1:numRep)) {
10   iSeq <- (1:n)
11   u <- rnorm(n, mean=0, sd=sigma.u)
12   y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
       beta1^2)
13   mean_collect[i] <- mean(y)}
14 hist(mean_collect,freq = F,main = "Histogram of the means of y",
      col = "white",xlab = "mean of y")
15 x.plot<-seq(mu -3*sqrt((sigma.u/(1-beta1))^2/n),mu + 3*sqrt((sigma
      .u/(1-beta1))^2/n),length.out = 400)
16 y.plot<-(1/(sqrt(2*pi)*(sqrt((sigma.u/(1-beta1))^2/n))))*exp(-0.5*
      (x_plot - mu)^2/((sigma.u/(1-beta1))^2/n))
17 points(x.plot,y.plot, type="l", col="red",lwd=4)
```
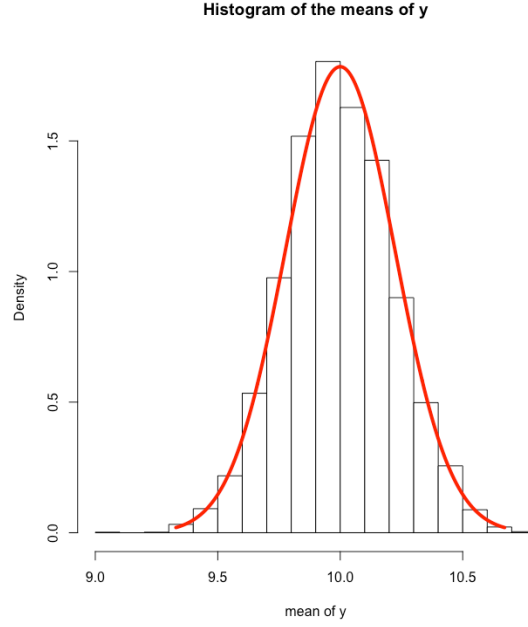


Figure 10: Histogram of y means

**II**

We have concluded that:

$$\bar{y_n} \sim N(\mu, [\sigma_u/(1-\beta_1)]^2 \frac{1}{n})$$

Since we know that for a normal distribution $X \sim N(\mu, \sigma^2)$, we have that:

$$Pr(\mu - 1.96\sigma \leq X \leq \mu + 1.96\sigma) = 95\%$$

Through some transformation we can get that:

$$Pr(X - 1.96\sigma \leq \mu \leq X + 1.96\sigma) = 95\% \tag{15}$$

In our case, we have that:

$$X = \bar{y}_n \qquad (16)$$

$$\sigma = \frac{1}{\sqrt{n}} \frac{\sigma_u}{1 - \beta_1} \qquad (17)$$

Substitute the $X$ and $\mu$ in equation 15 with equation 16 and 17 so we can get that:

$$Pr(\bar{y}_n - 1.96 \frac{1}{\sqrt{n}} \frac{\sigma_u}{1 - \beta_1} \leq \mu \leq \bar{y}_n + 1.96 \frac{1}{\sqrt{n}} \frac{\sigma_u}{1 - \beta_1}) = 95\% \qquad (18)$$

**III**

```
1  m <- 1000
2  n <- 2*m
3  numRep <- 1000
4  beta0 <- 1
5  beta1 <- 0.9
6  sigma.u <- 1
7  mu <- beta0/(1-beta1)
8  interval_collect=NULL
9  for (i in (1:numRep)) {
10    iSeq <- (1:n)
11    u <- rnorm(n, mean=0, sd=sigma.u)
12    y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
        beta1^2)
13    up.bound <- mean(y)+1.96/sqrt(n)*sigma.u/(1-beta1)
14    low.bound <- mean(y)-1.96/sqrt(n)*sigma.u/(1-beta1)
15    interval_collect[i] <- (mu>low.bound)&(mu<up.bound) }
16  sum(interval_collect)/length(interval_collect) # 0.942, the result
        is around 0.95
```

# C

## I

Verifying via summation whether both sides of the equality are equal. Calculating each term, we notice that the sum of the terms (50.19136) is really close to the variance of $s_y^2$ (50.16381).

```
1  set.seed(4110)
2  m <- 5000000
3  n <- 2*m
4  iSeq <- (1:n)
5  beta0 <- 1
6  beta1 <- 0.9
7  sigma.u <- 1
8  u <- rnorm(n, mean=0, sd=sigma.u)
9  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-beta1
      ^2)
```

21

```
10 term1 <- ((1+beta1^2)/(1-beta1^2)^2)*(1/(n-1))*sum(u^2)
11 term2 <- 0
12 for (i in (1:m)) {
13     term2 = term2+(1/(n-1))*u[2*i]*u[(2*i-1)]
14 }
15 term2<-4*beta1/((1-beta1)^2)*term2
16 term3 <- 1/((1-beta1)^2)*n/(n-1)*mean(u)^2
17 var(y) #50.16381
18 (term1+term2-term3) #50.19136
19 # From the two values above we can see that they are not exactly
       the same, but pretty close.
```

## II

Simulatting and comparing it with $\frac{1+\beta_1^2}{(1-\beta_1^2)^2}\sigma_\mu^2$ one can see that it is pretty close to $s_y^2$ 50.1385 , although it is far from $\sigma_\mu^2/(1-\beta_1)^2$ which equals to 100.

```
1 set.seed(4110)
2 m <- 5000000
3 n <- 2*m
4 beta0 <- 1
5 beta1 <- 0.9
6 sigma.u <- 1
7 mu <- beta0/(1-beta1)
8 iSeq <- (1:n)
9 u <- rnorm(n, mean=0, sd=sigma.u)
10 y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-beta1
       ^2)
11 var(y) #50.16381
12 (1+beta1^2)/(1-beta1^2)^2*sigma.u^2 #50.1385
13 sigma.u^2/(1-beta1)^2 #100
14 # From the three values above we can see that var(y) is close to
       (1+beta1^2)/(1-beta1^2)^2*sigma.u^2 which is 50.1385 but not
       even close to sigma.u^2/(1-beta1)^2 which is 100.
```

## D

### I

Looking at the Plot of $\beta_0$, $\beta_1$ and also the $\sigma_\mu$ estimates (Figure 11). The estimates of $\beta_0$ are pretty far from the actual value (represented by a vertical red line). The same applies to $\beta_1$ and $\sigma_\mu$ as well. Hence, it is indeed naive and far from reality estimation.

```
1 numRep=1000
2 m <- 100
3 n <- 2*m
4 beta0 <- 1
```

```r
5  beta1 <- 0.9
6  sigma.u <- 1
7  mu <- beta0/(1-beta1)
8  beta0hat_collect=NULL
9  beta1hat_collect=NULL
10 sd.u.hat_collect=NULL
11 for (i in (1:numRep)) {
12    iSeq <- (1:n)
13    u <- rnorm(n, mean=0, sd=sigma.u)
14    y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
        beta1^2)
15    lmfit <- lm(y[iSeq] ~ y[iSeq-(-1)^iSeq])
16    beta0hat <- lmfit$coefficients[1]
17    beta1hat <- lmfit$coefficients[2]
18    sd.u.hat <- sigma(lmfit)
19    beta1hat_collect[i] <- beta1hat
20    beta0hat_collect[i] <- beta0hat
21    sd.u.hat_collect[i] <- sd.u.hat
22 }
23 hist(beta0hat_collect,xlim = c(min(beta0hat_collect),1),col="white
      ",main = "Histogram of Estimated beta0hat",xlab = "beta0hat")
24 abline(v=beta0,col="red",lwd=4)
25 hist(beta1hat_collect,xlim = c(0.9,max(beta1hat_collect)),col="
      white",main = "Histogram of Estimated beta1hat",xlab = "
      beta1hat")
26 abline(v=0.9,col="red",lwd=4)
27 hist(sd.u.hat_collect,xlim = c(min(sd.u.hat_collect),1),col="white
      ",main = "Histogram of Estimated sigmaHat",xlab = "sigmaHat")
28 abline(v=sigma.u,col="red",lwd=4)
```



(a) beta0hat  (b) beta1hat  (c) sigmaHat
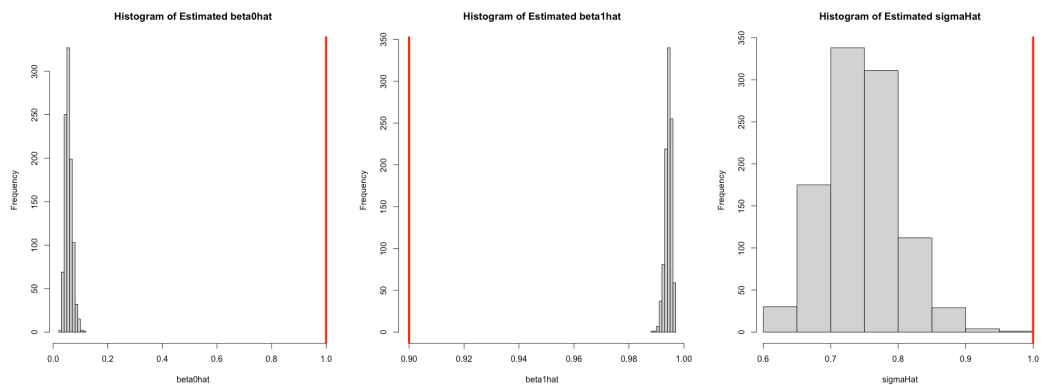
Figure 11

## II

In this case, the histograms (Figure 12) confirm that estimates for $\beta_0$, $\beta_1$, and $\sigma_\mu$ are all approximately reasonable estimates. The red vertical line indicates the actual $\beta_0$,

23

$\beta_1$, $\sigma_\mu$ respectively. One can further confirm it by looking at the mean of estimates. All of them are really close to the actual values. Hence we can verify that they are much better estimates than what we have seen previously.

```r
numRep=1000
m <- 100
n <- 2*m
beta0 <- 1
beta1 <- 0.9
sigma.u <- 1
mu <- beta0/(1-beta1)
beta1hat_collect=NULL
sd.u.hat_collect=NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  ratio <- var(y)/var(y[2*(1:m)] - y[2*(1:m)-1])
  beta1hat <- (2*ratio-sqrt(4*ratio-1))/(2*ratio-1)
  beta0hat <- mean(y)*(1-beta1hat)
  #extract residuals:
  u.hat <- y[iSeq] - (beta0hat + beta1hat*y[iSeq-(-1)^iSeq])
  sd.u.hat <- sd(u.hat)
  beta1hat_collect[i] <- beta1hat
  sd.u.hat_collect[i] <- sd.u.hat
}
hist(beta1hat_collect,freq = F,main = "Histogram of Beta1Hat")
abline(v=beta1,col="red",lwd=4)
hist(sd.u.hat_collect,main = "Histogram of Sd.u.Hat",xlab = "Sd.u.
    Hat")
abline(v=sigma.u,col="red",lwd=4)
```
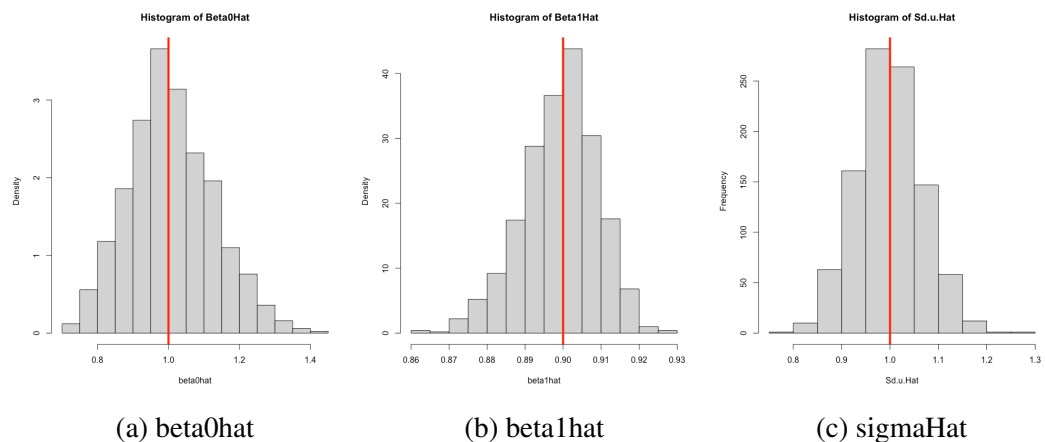


(a) beta0hat  (b) beta1hat  (c) sigmaHat

Figure 12

### III

Simulating $10^4$ times whether the upper limit of the confidence interval is larger than $\mu$, at the same time, the lower limit is smaller than $\mu$. In 95% of the cases, it is true.

```r
numRep=10000
m <- 100
beta0 = 1
n <- 2*m
beta1 = 0.9
sigma.u = 1
mu <- beta0/(1-beta1)
confinterval_collect=NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  ratio <- var(y)/var(y[2*(1:m)] - y[2*(1:m)-1])
  beta1hat <- (2*ratio-sqrt(4*ratio-1))/(2*ratio-1)
  beta0hat <- mean(y)*(1-beta1hat)
  #extract residuals:
  u.hat <- y[iSeq] - (beta0hat + beta1hat*y[iSeq-(-1)^iSeq])
  sd.u.hat <- sd(u.hat)
  up.bound <- mean(y)+1.96/sqrt(n)*sd.u.hat/(1-beta1hat)
  low.bound <- mean(y)-1.96/sqrt(n)*sd.u.hat/(1-beta1hat)
  confinterval_collect[i] <- (mu>low.bound)&(mu<up.bound)
}
sum(confinterval_collect)/length(confinterval_collect) #0.9438
    which is around 95%
```

# Appendix

```r
rm(list=ls())

# Assignment 1
## (A)
### (I)
data <- read.csv("ReplicationProcessedfinaldata04202018.csv")
miss <- data.frame(columns=colnames(data),missing.num=NA)
for (column in miss$columns){
  miss[miss$columns==column,2] <- sum(is.na(data[,column]))
}
rows.before <- nrow(data)
miss
### (II)
data <- data[-c(which((data$mate_age<12)|(is.na(data$mate_age))))
    ,]
rows.after <- nrow(data)
proportion.removed <- (rows.before-rows.after)/rows.before
proportion.removed #0.401764
### (III)
countries <- data.frame(country.name=unique(data$country),
    nCountries=NA)
for (name in countries$country.name) {
  countries[countries$country.name==name,2]=sum(data$country==name
    )
}
countries <- countries[-c(which(countries$nCountries<350)),]

## (B)
### (I)
### NB: sex = 0 are women, sex = 1 are men
with(data, plot(age, mate_age, cex = 1.2, pch = 21, bg=(sex+2),
    main="Male and Female"))
legend("topleft",c("Female","Male"),fill=c(2,3))
with(data[data$sex == 0,], lines(smooth.spline(age, mate_age), col
    = 2,lwd=5))
with(data[data$sex == 1,], lines(smooth.spline(age, mate_age), col
    = 3,lwd=5))

i <- 0
with(data[data$sex == i,], plot(age, mate_age, cex = 1.2, pch =
    21,bg=2,main="Female"))
with(data[data$sex == i,], lines(smooth.spline(age, mate_age), col
    = "lightgreen",lwd=5))
lmFemale <- lm(mate_age ~ age, data = data[data$sex== i,])
summary(lmFemale)
lmFemaleQuadratic <- lm(mate_age ~ age + I(age^2), data = data[
    data$sex == i,])
```

```r
39 summary(lmFemaleQuadratic)
40 #quadratic term not significant.
41 #Let's consider the fit of the simple linear regression:
42 x <- seq(min(data[data$sex == i,]$age), max(data[data$sex == i,]$
      age), length.out=400)
43 betahat <- lmFemale$coefficients
44 points(x, betahat[1]+betahat[2]*x, lwd=5, lty=2, col="black",type=
      "l")
45 abline(a=0,b=1,col="orange",lwd=5)
46 legend("topleft",c("Smooth Spline","Simple Linear","mate_age=age")
      ,lty = c(1,2,1),col=c("lightgreen","black","orange"),lwd=c
      (3,3,3))
47 #perfect match
48
49 i <- 1
50 lmMale <- lm(mate_age ~ age, data = data[data$sex== i,])
51 summary(lmMale)
52 lmMaleQuadratic <- lm(mate_age ~ age + I(age^2), data = data[data$
      sex == i,])
53 summary(lmMaleQuadratic) #quadratic term seems quite significant.
54 with(data[data$sex == i,],cor(age,age^2)) # 0.9861473
55 # the correlation between the age and quadratic age is quite large
      , so we cannot trust the p-value easily
56 #Let's further try the cubic term
57 lmMaleCubic <- lm(mate_age ~ age + I(age^2) + I(age^3), data =
      data[data$sex == i,])
58 summary(lmMaleCubic)
59 #In this case the linear term is not significant anymore, to
      figure out whether we should use the quadratic term or the
      cubic term,
60 #we will visually assess which one is better
61 with(data[data$sex == i,], plot(age, mate_age, cex = 1.2, pch =
      21,bg=3,main="Male"))
62 with(data[data$sex == i,], lines(smooth.spline(age, mate_age), col
       = "red",lwd=5))
63 # plot the quadratic line
64 x <- seq(min(data[data$sex == i,]$age), max(data[data$sex == i,]$
      age), length.out=400)
65 betahat <- lmMaleQuadratic$coefficients
66 points(x, betahat[1]+betahat[2]*x+betahat[3]*x^2, lwd=5, col="
      black",type="l")
67 x <- seq(min(data[data$sex == i,]$age), max(data[data$sex == i,]$
      age), length.out=400)
68 betahat <- lmMaleCubic$coefficients
69 points(x, betahat[1]+betahat[2]*x+betahat[3]*x^2+betahat[4]*x^3,
      lwd=5, col="skyblue",type="l")
70 abline(a=0,b=1,col="orange",lwd=5)
71 abline(lmMale,col="purple",lwd=5)
72 legend("topleft",c("Smooth Spline","Quadratic Linear","Cubic
```

```
        Linear","Simple Linear","mate_age=age"),col=c("red","black","
           skyblue","purple","orange"),lwd=c(3,3,3,3))
73  # Seems like the cubic estimation is not as good as the quadratic
           estimation, so let's go for the quadratic term
74
75  #Let's consider the fit of the quadratic linear regression:
76  with(data[data$sex == i,], plot(age, mate_age, cex = 1.2, pch =
           21,bg=3))
77  with(data[data$sex == i,], lines(smooth.spline(age, mate_age), col
            = "red",lwd=5))
78  x <- seq(min(data[data$sex == i,]$age), max(data[data$sex == i,]$
           age), length.out=400)
79  betahat <- lmMaleQuadratic$coefficients
80  points(x, betahat[1]+betahat[2]*x+betahat[3]*x^2, lwd=3, lty=2,
           col="black",type="l")
81  abline(a=0,b=1,col="skyblue",lwd=5)
82  legend("topleft",c("Smooth Spline","Quadratic Estimate","mate_age=
           age"),lty = c(1,2,1),col=c("red","black","skyblue"),lwd=c
           (3,3,3))
83  # perfect match
84
85  # For the full data, this motivates the following model
86  lm.joint <- with(data,lm(mate_age ~ age + sex*age + sex*I(age^2)-I
           (age^2)))
87  summary(lm.joint)
88
89  ### (II)
90  cuba.both <- data[data$country=="Cuba",]
91  cuba.male <- cuba.both[cuba.both$sex==1,]
92  cuba.female <- cuba.both[cuba.both$sex==0,]
93  lm.cuba <- with(cuba.both,lm(mate_age~age))
94  lm.cuba.male <- with(cuba.male,lm(mate_age~age))
95  lm.cuba.female <- with(cuba.female,lm(mate_age~age))
96  with(cuba.both,plot(age, mate_age, cex = 1.2, pch = 21,bg="orange"
           ,main="Cuba Male and Female"))
97  abline(lm.cuba,col="blue",lwd=5)
98  abline(a=0,b=1,col="red",lwd=5)
99  legend("topleft",c("Simple linear regression","mate_age=age"),lty=
           c(1,1),col=c("blue","red"),lwd=c(3,3))
100 with(cuba.male,plot(age, mate_age, cex = 1.2, pch = 21,bg="skyblue
           ",main="Cuba Male"))
101 legend("topleft",c("Simple linear regression","mate_age=age"),lty=
           c(1,1),col=c("blue","red"),lwd=c(3,3))
102 abline(lm.cuba.male,col="blue",lwd=5)
103 abline(a=0,b=1,col="red",lwd=5)
104 with(cuba.female,plot(age, mate_age, cex = 1.2, pch = 21,bg="
           yellow",main="Cuba Female"))
105 abline(lm.cuba.female,col="blue",lwd=5)
106 abline(a=0,b=1,col="red",lwd=5)
```

```r
107 legend("topleft",c("Simple linear regression","mate_age=age"),lty=
        c(1,1),col=c("blue","red"),lwd=c(3,3))
108
109 lm.cuba.null <- with(cuba.both, lm(mate_age ~ offset(1*age) - 1))
110 lm.cuba.null.male <- with(cuba.male, lm(mate_age ~ offset(1*age) -
        1))
111 lm.cuba.null.female <- with(cuba.female, lm(mate_age ~ offset(1*
        age) - 1))
112
113 anova1 <- anova(lm.cuba,lm.cuba.null)
114 anova2 <- anova(lm.cuba.male,lm.cuba.null.male)
115 anova3 <- anova(lm.cuba.female,lm.cuba.null.female)
116 max(residuals(lm.cuba.null)- with(cuba.both, mate_age - age)) == 0
117
118 ### (III)
119 five_country <- data[(data$country=="Hungary") | (data$country=="
        Pakistan")|(data$country=="Poland")|(data$country=="Slovenia")
        |(data$country=="Turkey"),]
120 five_country$Hungary <- (five_country$country=="Hungary")*1
121 five_country$Pakistan <- (five_country$country=="Pakistan")*1
122 five_country$Poland <- (five_country$country=="Poland")*1
123 five_country$Slovenia <- (five_country$country=="Slovenia")*1
124 five_country$Turkey <- (five_country$country=="Turkey")*1
125 with(five_country, plot(age,mate_age, cex = 1, pch = 20, col=(1*
        Hungary+2*Pakistan+3*Poland+4*Slovenia+7*Turkey),main="Scatter
        Plot of 5 Countries"))
126 with(five_country[five_country$Hungary == 1,], lines(smooth.spline
        (age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*
        Slovenia+7*Turkey),lwd=20))
127 with(five_country[five_country$Pakistan == 1,], lines(smooth.
        spline(age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*
        Slovenia+7*Turkey),lwd=5))
128 with(five_country[five_country$Poland == 1,], lines(smooth.spline(
        age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*Slovenia
        +7*Turkey),lwd=2))
129 with(five_country[five_country$Slovenia == 1,], lines(smooth.
        spline(age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*
        Slovenia+7*Turkey),lwd=4))
130 with(five_country[five_country$Turkey == 1,], lines(smooth.spline(
        age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*Slovenia
        +7*Turkey),lwd=2))
131 legend("topleft",c("Hungary","Pakistan","Poland","Slovenia","
        Turkey"),fill=c(1,2,3,4,7))
132
133 # For Hungary, Pakistan and Slovenia together
134 threeCountry = five_country[(five_country$country!="Poland")&(five
        _country$country!="Turkey"),]
135 with(threeCountry, plot(age, mate_age, cex = 1.2, pch = 21, main="
        Three Countries"))
```

```r
with(threeCountry, lines(smooth.spline(age, mate_age),col = (1*
    Hungary+2*Pakistan+4*Slovenia),lwd=5))
lmThree <- lm(mate_age ~ age, data = threeCountry)
summary(lmThree)
lmThreeQuadratic <- lm(mate_age ~ age + I(age^2), data =
    threeCountry)
summary(lmThreeQuadratic)
lmThreeCubic <- lm(mate_age ~ age + I(age^2) +I(age^3), data =
    threeCountry)
summary(lmThreeCubic)
with(threeCountry,cor(age,age^2))
with(threeCountry,cor(age,age^3))
with(threeCountry,cor(age^2,age^3))
betahat<-lmThree$coefficients
#quadratic and cubic term are significant.
#However, the correlations among these three terms are extremely
    high, so we should be carefully with the quadratic and cubic
    terms
#Let's plot the three linear lines to see which one fits the best.
x <- seq(min(threeCountry$age), max(threeCountry$age), length.out
    =400)
with(threeCountry,points(x, betahat[1]+betahat[2]*x, lwd=5, col="
    red",type="l"))
betahatQua <- lmThreeQuadratic$coefficients
with(threeCountry,points(x, betahatQua[1]+betahatQua[2]*x+
    betahatQua[3]*x^2, lwd=5, col="blue",type="l"))
betahatCub <- lmThreeCubic$coefficients
with(threeCountry,points(x, betahatCub[1]+betahatCub[2]*x+
    betahatCub[3]*x^2+betahatCub[4]*x^3, lwd=5, col="green",type="l
    "))
legend("topleft", legend=c("Smooth", "Simple Linear","Quadratic
    Linear","Cubic Linear"),lty = c(1,1,1,), col=c("black", "red","
    blue","green"), lwd=c(3,3,3,3))
# Seems like the simple linear fits the best, so let's go for it

# For Poland
with(five_country[five_country$Poland == 1,], plot(age, mate_age,
    cex = 1.2, pch = 21,bg=(1*Hungary+2*Pakistan+3*Poland+4*
    Slovenia+5*Turkey), main="Poland"))
with(five_country[five_country$Poland == 1,], lines(smooth.spline(
    age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*Slovenia
    +5*Turkey+1),lwd=5))
lmPoland <- lm(mate_age ~ age, data = five_country[five_country$
    Poland== 1,])
summary(lmPoland)
lmPolandQuadratic <- lm(mate_age ~ age + I(age^2), data = five_
    country[five_country$Poland == 1,])
summary(lmPolandQuadratic)
lmPolandCubic <- lm(mate_age ~ age + I(age^2) +I(age^3), data =
```

```r
          five_country[five_country$Poland == 1,])
167 summary(lmPolandCubic)
168 #quadratic and cubic term are not significant.
169 #Let's consider the fit of the simple linear regression:
170 x <- seq(min(five_country[five_country$Poland == 1,]$age), max(
        five_country[five_country$Poland == 1,]$age), length.out=400)
171 betahat <- lmPoland$coefficients
172 with(five_country,points(x, betahat[1]+betahat[2]*x, lwd=5, lty=2,
        col="red",type="l"))
173 #perfect match
174 legend("topleft", legend=c("Smooth", "Simple Linear"),lty = c(1,2)
        , col=c(4, "red"), lwd=c(3,3))
175
176 # For Turkey
177 with(five_country[five_country$Turkey == 1,], plot(age, mate_age,
        cex = 1.2, pch = 21,bg=(1*Hungary+2*Pakistan+3*Poland+4*
        Slovenia+5*Turkey), main="Turkey"))
178 with(five_country[five_country$Turkey == 1,], lines(smooth.spline(
        age, mate_age), col = (1*Hungary+2*Pakistan+3*Poland+4*Slovenia
        +5*Turkey+1),lwd=5))
179 lmTurkey <- lm(mate_age ~ age, data = five_country[five_country$
        Turkey== 1,])
180 summary(lmTurkey)
181 lmTurkeyQuadratic <- lm(mate_age ~ age + I(age^2), data = five_
        country[five_country$Turkey == 1,])
182 summary(lmTurkeyQuadratic)
183 lmTurkeyCubic <- lm(mate_age ~age+ I(age^2) +I(age^3), data = five
        _country[five_country$Turkey == 1,])
184 summary(lmTurkeyCubic)
185 #Let's consider the fit of the cubic linear regression:
186 x <- seq(min(five_country[five_country$Turkey == 1,]$age), max(
        five_country[five_country$Turkey == 1,]$age), length.out=400)
187 betahat.quadratic <- lmTurkeyQuadratic$coefficients
188 with(five_country[five_country$Turkey == 1,],points(x, betahat.
        quadratic[1]+betahat.quadratic[2]*x+betahat.quadratic[3]*x^2,
        lwd=5, lty=1, col="blue",type="l"))
189 betahat.cubic <- lmTurkeyCubic$coefficients
190 with(five_country[five_country$Turkey == 1,],points(x, betahat.
        cubic[1]+betahat.cubic[2]*x+betahat.cubic[3]*x^2+betahat.cubic
        [4]*x^3, lwd=5, lty=1, col="green",type="l"))
191 legend("topleft", legend=c("Smooth", "Quadratic Linear","Cubic
        Linear"), col=c(6,"blue","green"), lwd=c(3,3,3))
192
193 # add the interaction terms
194 # the linear age term of Turkey Cubic regression is not
        significant, here we don't include the simple linear
        interaction term for Turkey (Turkey*age)
195 # since it is not clear whether we need the interaction term for
        Poland (Poland*age)
```

```r
196  # we first run the model without the interaction term
197  lm1 <- with(five_country,lm(mate_age~age+Poland+Turkey+Turkey*I(
        age^2)+Turkey*I(age^3)-I(age^3)-I(age^2)))
198  summary(lm1)
199  lm2 <- with(five_country,lm(mate_age~age+Poland+Poland*age+Turkey+
        Turkey*I(age^2)+Turkey*I(age^3)-I(age^3)-I(age^2)))
200  summary(lm2)
201  stargazer(lm2,lm1)
202  anova(lm2,lm1)
203  # seems like perfect
204
205
206  # Assignment 2
207  ## (A) See the main body
208  ## (B) See the main body
209  ## (C) See the main body
210  ## (D)
211  set.seed(4110)
212  beta0 <- 1
213  beta1 <- 0.9
214  sigma.u <- 1
215  m <- 1000
216  n <- 2*m
217  iSeq <- (1:n)
218  u <- rnorm(n, mean=0, sd=sigma.u)
219  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-beta1
        ^2)
220  ### (I)
221  ### (II)
222  ### (III)
223  check <- NULL
224  y.calculated <- NULL
225  for (i in (1:n)) {
226    check[i] <- (beta0+beta1*y[i-(-1)^i]+u[i]-y[i])<1e-10
227    y.calculated[i] <- beta0+beta1*y[i-(-1)^i]+u[i]
228  }
229  sum(check)==n
230  plot(y.calculated,y,cex=1.5,pch=21,main = "Simulated y and
        calculated y",xlab="Y Calculated",ylab="Y Simulated")
231  abline(a=0,b=1,col="red",lwd=4)
232  # which verifies that every y satisfies the equation
233
234  ## (E)
235  ### (I)
236  numRep=1000
237  m <- 100
238  n <- 2*m
239  beta0 = 1
240  beta1 = 0.9
```

```r
sigma.u = 1
mu = beta0/(1-beta1)
mean.y.collect <- NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  mean.y.collect[i] <- mean(y)
}
hist(mean.y.collect,main = "Histogram of y means",col = "white",
    xlab = "mean of y")
abline(v=mu,col="red",lwd=4)
### (II)
numRep=1000
m <- 100
n <- 2*m
beta0 = 1
beta1 = 0.9
sigma.u = 1
counter <- NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  lo.lim <- mean(y) - 1.96*sqrt(var(y))/sqrt(n)
  up.lim <- mean(y) + 1.96*sqrt(var(y))/sqrt(n)
  counter[i] <- (mu>lo.lim) & (mu<up.lim)
}
sum(counter)/length(counter) #0.837

# Assignment 3
## (A)
### (I) See the main body
### (II)
m <- 1000
n <- 2*m
numRep <- 1000
beta0 <- 1
beta1 <- 0.9
sigma.u <- 1
mu <- beta0/(1-beta1)
xy.collect <- NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
```

```r
286   x <- mu + 1/(1-beta1)*u
287   xy.collect[i] <- (mean(x)-mean(y))<1e-10
288 }
289 sum(xy.collect) #10000 which equals to the length of xy.collect,
        so for all the 1000 simulations, mean(x)=mean(y)
290 ## (B)
291 ### (I)
292 m <- 1000
293 n <- 2*m
294 numRep <- 5000
295 beta0 <- 1
296 beta1 <- 0.9
297 sigma.u <- 1
298 mu <- beta0/(1-beta1)
299 mean_collect=NULL
300 for (i in (1:numRep)) {
301   iSeq <- (1:n)
302   u <- rnorm(n, mean=0, sd=sigma.u)
303   y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
        beta1^2)
304   mean_collect[i] <- mean(y)
305 }
306
307 hist(mean_collect,freq = F,main = "Histogram of the means of y",
        col = "white",xlab = "mean of y")
308 x.plot<-seq(mu -3*sqrt((sigma.u/(1-beta1))^2/n),mu + 3*sqrt((sigma
        .u/(1-beta1))^2/n),length.out = 400)
309 y.plot<-(1/(sqrt(2*pi)*(sqrt((sigma.u/(1-beta1))^2/n))))*exp(-0.5*
        (x_plot - mu)^2/((sigma.u/(1-beta1))^2/n))
310 points(x.plot,y.plot, type="l", col="red",lwd=4)
311
312 ### (II) See the main body
313 ### (III)
314 m <- 1000
315 n <- 2*m
316 numRep <- 1000
317 beta0 <- 1
318 beta1 <- 0.9
319 sigma.u <- 1
320 mu <- beta0/(1-beta1)
321 interval_collect=NULL
322 for (i in (1:numRep)) {
323   iSeq <- (1:n)
324   u <- rnorm(n, mean=0, sd=sigma.u)
325   y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
        beta1^2)
326   up.bound <- mean(y)+1.96/sqrt(n)*sigma.u/(1-beta1)
327   low.bound <- mean(y)-1.96/sqrt(n)*sigma.u/(1-beta1)
328   interval_collect[i] <- (mu>low.bound)&(mu<up.bound)
```

```r
329 }
330 sum(interval_collect)/length(interval_collect) # 0.942, the result
        is around 0.95
331 ## (C)
332 ### (I)
333 set.seed(4110)
334 m <- 5000000
335 n <- 2*m
336 iSeq <- (1:n)
337 beta0 <- 1
338 beta1 <- 0.9
339 sigma.u <- 1
340 u <- rnorm(n, mean=0, sd=sigma.u)
341 y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-beta1
        ^2)
342 term1 <- (1+beta1^2)/(1-beta1^2)^2*(1/(n-1))*sum(u)^2
343 term2 <- 0
344 for (i in (1:m)) {
345   term2 = term2+u[2*i]*u[(2*i-1)]
346 }
347 term2<-4*(1/(n-1))*beta1/((1-beta1)^2)*term2
348 term3 <- 1/((1-beta1)^2)*n/(n-1)*(mean(u))^2
349 var(y) #50.16381
350 (term1+term2-term3) #50.19136
351
352 # From the two values above we can see that they are not exactly
        the same, but pretty close.(mean(u)^2)
353 ### (II)
354 set.seed(4110)
355 m <- 10000000
356 n <- 2*m
357 beta0 <- 1
358 beta1 <- 0.9
359 sigma.u <- 1
360 mu <- beta0/(1-beta1)
361 iSeq <- (1:n)
362 u <- rnorm(n, mean=0, sd=sigma.u)
363 y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-beta1
        ^2)
364 var(y) #50.14622
365 (1+beta1^2)/(1-beta1^2)^2*sigma.u^2 #50.1385
366 sigma.u^2/(1-beta1)^2 #100
367 # From the three values above we can see that var(y) is close to
        (1+beta1^2)/(1-beta1^2)^2*sigma.u^2 which is 50.1385 but not
        even close to
368 # sigma.u^2/(1-beta1)^2 which is 100.
369
370 ## (D)
371 ### (I)
```

```r
numRep=1000
m <- 100
n <- 2*m
beta0 <- 1
beta1 <- 0.9
sigma.u <- 1
mu <- beta0/(1-beta1)
beta0hat_collect=NULL
beta1hat_collect=NULL
sd.u.hat_collect=NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
  lmfit <- lm(y[iSeq] ~ y[iSeq-(-1)^iSeq])
  beta0hat <- lmfit$coefficients[1]
  beta1hat <- lmfit$coefficients[2]
  sd.u.hat <- sigma(lmfit)
  beta1hat_collect[i] <- beta1hat
  beta0hat_collect[i] <- beta0hat
  sd.u.hat_collect[i] <- sd.u.hat
}
hist(beta0hat_collect,xlim = c(min(beta0hat_collect),1),main = "
    Histogram of Estimated beta0hat",xlab = "beta0hat")
abline(v=beta0,col="red",lwd=4)
hist(beta1hat_collect,xlim = c(0.9,max(beta1hat_collect)),main = "
    Histogram of Estimated beta1hat",xlab = "beta1hat")
abline(v=0.9,col="red",lwd=4)
hist(sd.u.hat_collect,xlim = c(min(sd.u.hat_collect),1),main = "
    Histogram of Estimated sigmaHat",xlab = "sigmaHat")
abline(v=sigma.u,col="red",lwd=4)

### (II)
numRep=1000
m <- 100
n <- 2*m
beta0 <- 1
beta1 <- 0.9
sigma.u <- 1
mu <- beta0/(1-beta1)
beta1hat_collect=NULL
beta0hat_collect=NULL
sd.u.hat_collect=NULL
for (i in (1:numRep)) {
  iSeq <- (1:n)
  u <- rnorm(n, mean=0, sd=sigma.u)
  y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
    beta1^2)
```

```r
    ratio <- var(y)/var(y[2*(1:m)] - y[2*(1:m)-1])
    beta1hat <- (2*ratio-sqrt(4*ratio-1))/(2*ratio-1)
    beta0hat <- mean(y)*(1-beta1hat)
    #extract residuals:
    u.hat <- y[iSeq] - (beta0hat + beta1hat*y[iSeq-(-1)^iSeq])
    sd.u.hat <- sd(u.hat)
    beta0hat_collect[i] <- beta0hat
    beta1hat_collect[i] <- beta1hat
    sd.u.hat_collect[i] <- sd.u.hat
}
hist(beta0hat_collect,freq = F,main = "Histogram of Beta0Hat",xlab
    = "beta0hat")
abline(v=beta0,col="red",lwd=4)
hist(beta1hat_collect,freq = F,main = "Histogram of Beta1Hat",xlab
    = "beta1hat")
abline(v=beta1,col="red",lwd=4)
hist(sd.u.hat_collect,main = "Histogram of Sd.u.Hat",xlab = "Sd.u.
    Hat")
abline(v=sigma.u,col="red",lwd=4)

### (III)
numRep=10000
m <- 100
beta0 = 1
n <- 2*m
beta1 = 0.9
sigma.u = 1
mu <- beta0/(1-beta1)
confinterval_collect=NULL
for (i in (1:numRep)) {
    iSeq <- (1:n)
    u <- rnorm(n, mean=0, sd=sigma.u)
    y <- beta0/(1-beta1) + (u[iSeq]+beta1*u[iSeq-(-1)^iSeq])/(1-
      beta1^2)
    ratio <- var(y)/var(y[2*(1:m)] - y[2*(1:m)-1])
    beta1hat <- (2*ratio-sqrt(4*ratio-1))/(2*ratio-1)
    beta0hat <- mean(y)*(1-beta1hat)
    #extract residuals:
    u.hat <- y[iSeq] - (beta0hat + beta1hat*y[iSeq-(-1)^iSeq])
    sd.u.hat <- sd(u.hat)
    up.bound <- mean(y)+1.96/sqrt(n)*sd.u.hat/(1-beta1hat)
    low.bound <- mean(y)-1.96/sqrt(n)*sd.u.hat/(1-beta1hat)
    confinterval_collect[i] <- (mu>low.bound)&(mu<up.bound)
}
sum(confinterval_collect)/length(confinterval_collect) #0.9438
    which is around 95%
```