

# Final

Yi Chen (yc3356)

Section 003#

## Instructions (Read this completely first)

You should complete the exam by editing this file directly. Please knit the file often, so that if you make a mistake you catch it before the end of the exam. You will have exactly 160 minutes from your start time to complete the exam. **At the end you must turn in your knitted .pdf file and raw .Rmd file on Courseworks.**

**When the time is up, you must shut your computer immediately.** We will take off points from anyone whose computer is still open after time is up.

**You may use your class notes for the exam, but not the internet. You absolutely may not communicate with anyone else during the exam. Doing so will result in an F in this class and likely result in termination from the MA program.**

## Question 0 (5 points)

- a. (0.5 points) Place your section number as the date of the document. If you don't know your section number, you can determine it below based on when your lab meets.
  - Section 002 Lab meets TR 7:40pm-8:55pm
  - Section 003 Lab meets TR 11:40am-12:55pm
  - Section 004 Lab meets MW 8:40am-9:55am
  - Section 005 Lab meets TR 8:40am-9:55am
- b. (0.5 points) Write your name and UNI as the author of the document.
- c. (4 points) Please present your answers in a readable format. This includes things like indenting your code and generally presenting easy-to-read code. Presentation of the overall Markdown document will be considered as well.

## Question 1: Fitting Data (49 points)

- a. (4 points) You have an urn with 30 balls – 10 are red, 10 are blue, and 10 are green. Write a single line of code to simulate randomly picking 400 balls from the urn with replacement. Create a variable `num_green` that records the number of green balls selected in the 400 draws.

```
set.seed(1)
```

```
# Your answer to question 1.a here. Don't remove the set.seed(1) command.  
# Let 1 be the red, 2 be the blue, 3 be the green  
sample_1 <- sample(rep(1:3,each=10),size = 400,prob = c(rep(1/30,30)),replace = TRUE)  
num_green <- sum(sample_1 == 3)  
num_green
```

```
## [1] 124
```

- b. (4 points) Now repeat the above experiment 1000 times. Create a vector `data`, such that each element in `data` is the result (counting the number of green balls) from an independent trial like that described in

1.a.

```
set.seed(2)

# Your answer to question 1.b here. Don't remove the set.seed(2) command.
data <- rep(NA,1000)
for(i in 1:1000){
  sample_this <- sample(rep(1:3,each=10),size = 400,prob = c(rep(1/30,30)),replace = TRUE)
  num_green_this <- sum(sample_this == 3)
  data[i] <- num_green_this
}
head(data)
```

```
## [1] 129 142 128 127 123 136
```

# If you can't produce a data vector, uncomment the following line of code  
# and use it for the rest of the questions:

```
# data <- rnorm(1000, 133, 9)
```

- c. (6 points) Note that if a random variable  $X$  is the number of green balls selected in 400 draws with replacement from the urn, then  $X$  follows a binomial distribution, namely  $X \sim \text{bin}(n, p)$  where  $p$  is the probability of selecting a green ball from the urn in a single draw,  $n$  is the total number of draws, and

$$Pr(X = x) = \binom{n}{x} p^x (1 - p)^{n-x} \quad \text{for } x = 0, 1, \dots, 400,$$

with  $\mathbb{E}[X] = np$ . Recall that the binomial distribution is well-approximated by the normal distribution. To see that this is a good approximation, plot a histogram of your data from 1.b along with a normal density curve colored red having mean  $np = 400 * (1/3)$  and variance  $np(1 - p) = 400 * (1/3) * (2/3)$ .

```
# Your answer to question 1.c here.
library(ggplot2)
```

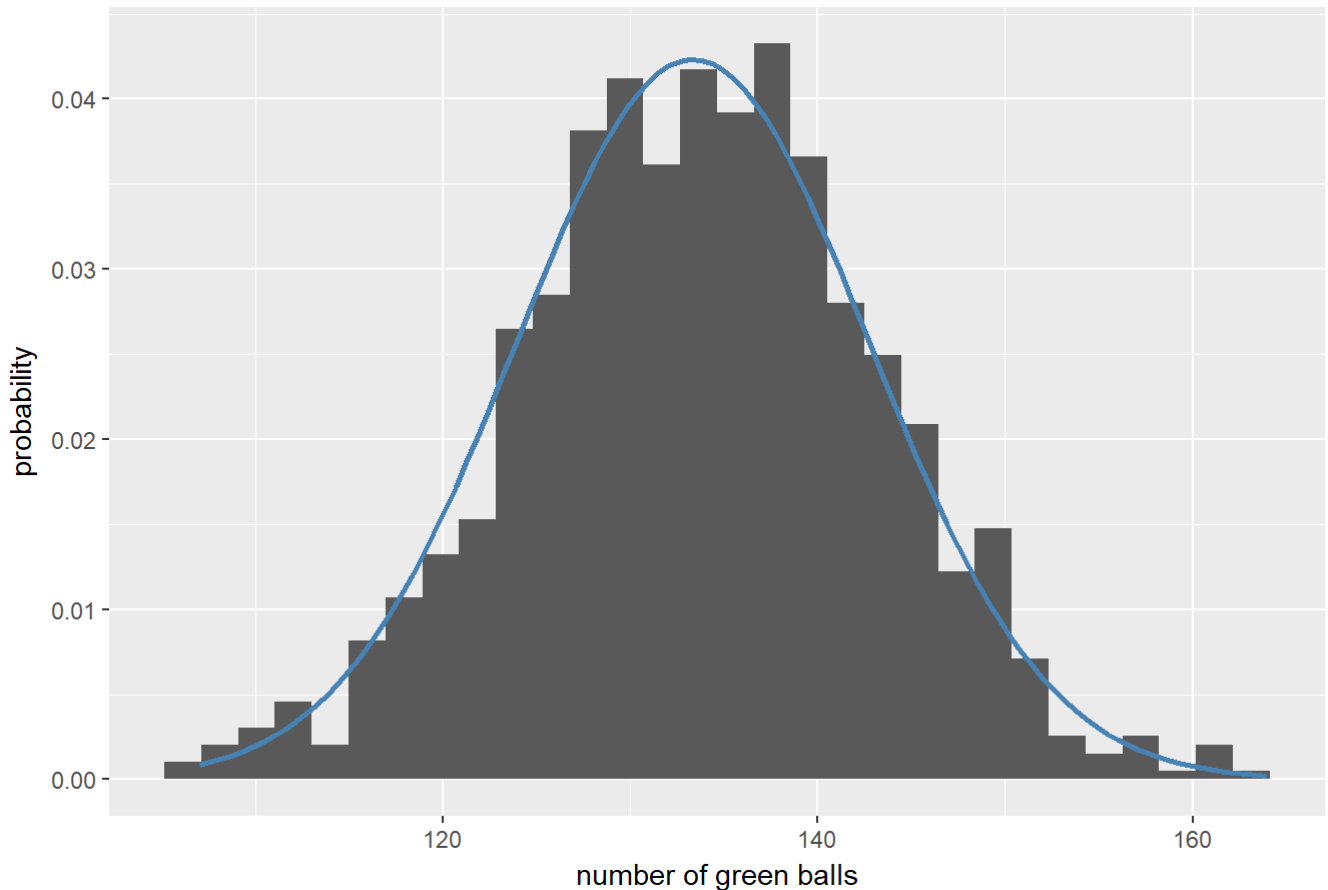
```
## Warning: package 'ggplot2' was built under R version 3.4.2
```

```
mean <- 400 * (1/3)
standard_deviation <- sqrt(400 * (1/3) * (2/3))
normal_args <- list(mean,standard_deviation)

ggplot()+
  geom_histogram(aes(x=data,y=..density..))+
  stat_function(aes(x=data), fun = dnorm,args = normal_args,col="steelblue",lwd=1)+
  labs(title="The estimated data v.s. the expected density",x="number of green balls",y="probability")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

The estimated data v.s. the expected density



- d. (5 points) Give the proportion of values in your data vector that are less than or equal to 100 or greater than or equal to 150. Using R functions for probability distributions, compare this proportion to the probability that a normal random variable having mean  $np = 400 * (1/3)$  and variance  $np(1 - p) = 400 * (1/3) * (2/3)$  is less than 100 or greater than 150.

```
# Your answer to question 1.d here.
propotion <- mean(data<=100 | data>=150)

expect_propotion <- pnorm(100,mean = mean,sd = standard_deviation) + pnorm(150,mean = mean,sd
= standard_deviation,lower.tail = FALSE)

propotion;expect_propotion
```

```
## [1] 0.045
```

```
## [1] 0.03875341
```

- e. (5 points) Write a function `MomentEstimator` that takes two input: `data`, a vector containing the number of green balls selected in each experiment, and `n` the total number of balls selected in each experiment (in 1.a,  $n = 400$  but we write the function where this could change) and returns a single output value `phat` that is the method of moments estimate of the probability  $p$ . After the function is written run the code `MomentEstimator(data, 400)` to see the method of moment estimator from your simulated data in 1.b and the code `MomentEstimator(80, 100)` to check the functionality.

```
# Your answer to question 1.e here.
```

```
MomentEstimator <- function(data,n){  
  mu <- mean(data)  
  phat <- mu/n  
  
  return(phat)  
}  
  
MomentEstimator(data,400)
```

```
## [1] 0.333905
```

```
MomentEstimator(80,100)
```

```
## [1] 0.8
```

f. (7 points) If  $num$  is the number of experiments run (i.e. in 1.b  $num$  is 1000) and  $x_i$  is the number of green balls selected in experiment  $i = 1, 2, \dots, num$ , then the log-likelihood for the binomial distribution is given by the following:

$$\ell(p) = \sum_{i=1}^{num} \left( \log \binom{n}{x_i} + x_i \log p + (n - x_i) \log(1 - p) \right).$$

Find the MLE estimate by writing a function that calculates the negative log-likelihood and then using `nlm()` to minimize it. Find the MLE estimate in this way on your data from part 1.b. Use an initial guess of  $p = 0.5$ .

```
# Your answer to question 1.f here.
```

```
Negloglikelihood <- function(data,p,num){  
  
  return(-sum(dbinom(data,prob = p,size = num,log = TRUE)))  
}  
  
nlm(Negloglikelihood,p=0.5,data=data,num=400)
```

```
## $minimum  
## [1] 3680.913  
##  
## $estimate  
## [1] 0.3339045  
##  
## $gradient  
## [1] -0.002743491  
##  
## $code  
## [1] 1  
##  
## $iterations  
## [1] 5
```

*# actually the mean of the model is indeed 400 rather than 1000. Otherwise the estimated value of this is too small.*

- g. (10 points) Use the bootstrap procedure to estimate the variance of your method of moments estimator. Use 5000 bootstrap resamples of the data (stored in vector `data`) you calculated in 1.b. The actual variance of the method of moments estimator is  $5.56e - 07$ .

```
set.seed(3)

# Your answer to question 1.g here. Don't remove the set.seed(3) command.
B <- 5000
param_ests <- rep(NA,nrow = B)

for(b in 1:B){
  resamp <- sample(data,1000,replace = TRUE)
  param_ests[b] <- MomentEstimator(resamp,n=400)
}

var(param_ests)
```

```
## [1] 5.974888e-07
```

- h. (8 points) Use simulation to provide evidence that the method of moments estimate is consistent (meaning that as the sample size increases  $n$ , the estimator converges to the population value).

```
set.seed(4)

# Your answer to question 1.h here. Don't remove the set.seed(4) command.

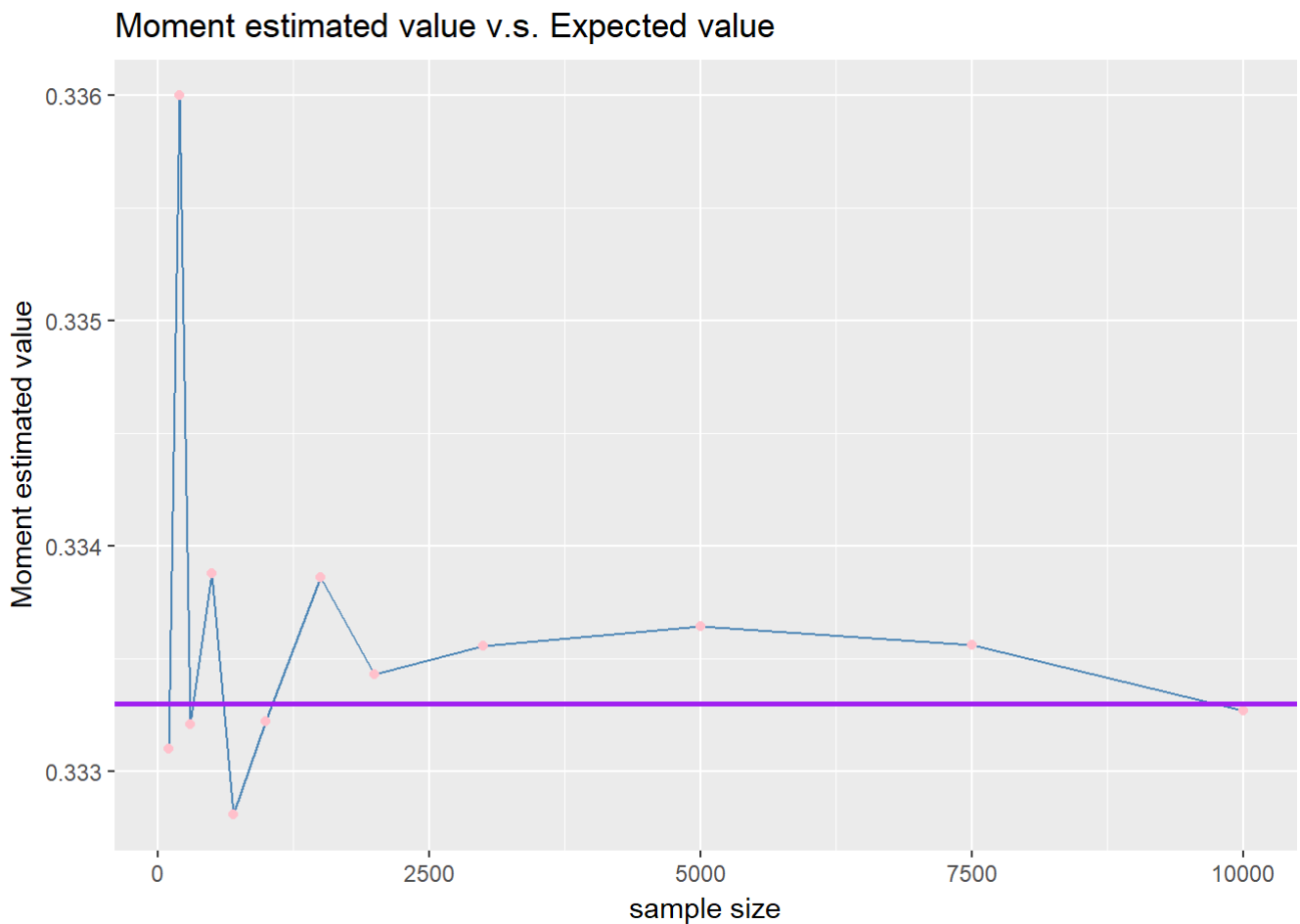
mu <- c(100,200,300,500,700,1000,1500,2000,3000,5000,7500,10000)

MM <- rep(NA,length(mu))
for(i in 1:length(mu)){
  data <- rep(NA,mu[i])
  for(j in 1:mu[i]){
    sample_this <- sample(rep(1:3,each=10),size = 400,prob = c(rep(1/30,30)),replace = TRUE)
    num_green_this <- sum(sample_this == 3)
    data[j] <- num_green_this
  }
  MMthis <- MomentEstimator(data,400)
  MM[i] <- MMthis
}

length(MM)
```

```
## [1] 12
```

```
ggplot()+
  geom_line(aes(x=mu,y=MM),col="steelblue")+
  geom_point(aes(x=mu,y=MM),col="pink")+
  geom_hline(yintercept = 0.3333,lwd=1,col='purple')+
  labs(title="Moment estimated value v.s. Expected value",x="sample size",y="Moment estimated value")
```



*# as we can see with the sample size becoming bigger and bigger, the estimated value of moment converge to the true value 0.33*

## Question 2: Transforming Data (46 points)

Gross domestic product (GDP) is a measure of the total market value of all goods and services produced in a given country in a given year. The percentage growth rate of GDP in year  $t$  is

$$100 \times \left( \frac{GDP_{t+1} - GDP_t}{GDP_t} \right) - 100$$

An important claim in economics is that the rate of GDP growth is closely related to the level of government debt, specifically with the ratio of the government's debt to the GDP. The file `debt.csv` contains measurements of GDP growth and of the debt-to-GDP ratio for twenty countries around the world, from the 1940s to 2010. Note that not every country has data for the same years, and some years in the middle of the period are missing data for some countries but not others.

```
debt <- read.csv("debt.csv", as.is = TRUE)
dim(debt)
```

```
## [1] 1171    4
```

```
head(debt)
```

```
##      Country Year    growth    ratio
## 1 Australia 1946 -3.557951 190.41908
## 2 Australia 1947  2.459475 177.32137
## 3 Australia 1948  6.437534 148.92981
## 4 Australia 1949  6.611994 125.82870
## 5 Australia 1950  6.920201 109.80940
## 6 Australia 1951  4.272612  87.09448
```

- a. (5 points) Calculate the average GDP growth rate for each year (averaging over countries). This is a classic split/apply/combine problem, and you should use `split()` and a function from the apply family of functions to solve it. You should not need to use a loop to do this. (The average growth rates for 1972 and 1989 should be 5.63 and 3.19, respectively. Print these values in your output.)

```
# Your answer to question 2.a here.

split_year <- split(debt,debt$Year)

average_function <- function(df){
  return(mean(df$growth,na.rm=TRUE))
}

year.average1 <- sapply(split_year,average_function)
year.average1["1972"];year.average1["1989"];
```

```
##      1972
## 5.629986
```

```
##      1989
## 3.186842
```

- b. (5 points) Calculate the average GDP growth rate for each year (averaging over countries). This is a classic split/apply/combine problem, and you should use `ddply()` to solve it. Save your output as `year.avgs` and change the column names to be `Year` and `AverageGrowth`. You should not need to use a loop to do this. (The average growth rates for 1972 and 1989 should be 5.63 and 3.19, respectively. Print these values in your output.)

```
# Your answer to question 2.b here.
library(plyr)
year.avgs <- ddply(debt,.(debt$Year),average_function)
colnames(year.avgs) <- c("Year","AverageGrowth")
year.avgs$AverageGrowth[year.avgs$Year==1972];
```

```
## [1] 5.629986
```

```
year.avgs$AverageGrowth[year.avgs$Year==1989];
```

```
## [1] 3.186842
```

- c. (4 points) The `year.avgs` dataframe from 2.b will be sorted by Year, meaning row 1 corresponds to 1946, row 2 to 1947, and so on with row 64 corresponding to 2009. Produce a dataframe that instead has row 1 corresponding to the year with the largest average growth, row 2 to the year with the second largest average growth, and so on with row 64 corresponding to the year with the smallest average growth.

```
# Your answer to question 2.c here.  
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:plyr':  
##  
##   arrange, count, desc, failwith, id, mutate, rename, summarise,  
##   summarize
```

```
## The following objects are masked from 'package:stats':  
##  
##   filter, lag
```

```
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union
```

```
year.avgs_new <-year.avgs %>%  
  arrange(desc(AverageGrowth))  
head(year.avgs_new)
```

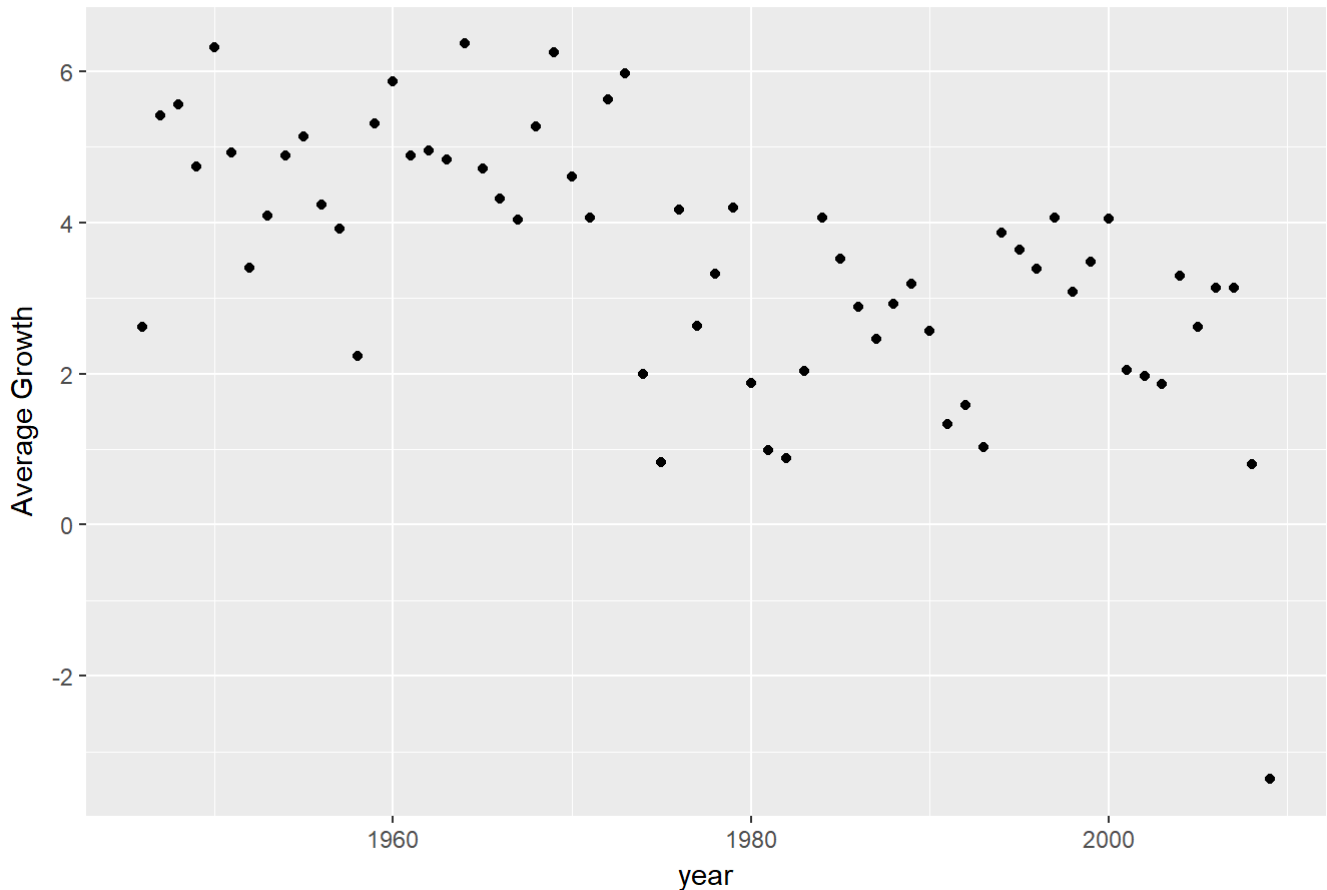
```
##   Year AverageGrowth  
## 1 1964      6.365472  
## 2 1950      6.321490  
## 3 1969      6.247051  
## 4 1973      5.971243  
## 5 1960      5.860439  
## 6 1972      5.629986
```

- d. (3 points) Make a plot of the growth rates (y-axis) versus the year (x-axis) using your results from either 2.a or 2.b (they should be the same). Make sure the axes are labeled appropriately.

```
# Your answer to question 2.d here.  
ggplot(data=year.avgs)+  
  geom_point(aes(x=Year,y=AverageGrowth))+  
  labs(title="Average Growth in every year",x="year",y="Average Growth")
```



Average Growth in every year



- e. (6 points) The function `cor(x,y)` calculates the correlation coefficient between two vectors `x` and `y`. First calculate the correlation coefficient between GDP growth and the debt ratio over the whole data set (all countries, all years). Your answer should be  $-0.1995$ . Second, compute the correlation coefficient separately for each year, and plot a histogram of these coefficients. The mean of these correlations should be  $-0.1906$ . Do not use a loop.

```
# Your answer to question 2.e here.
```

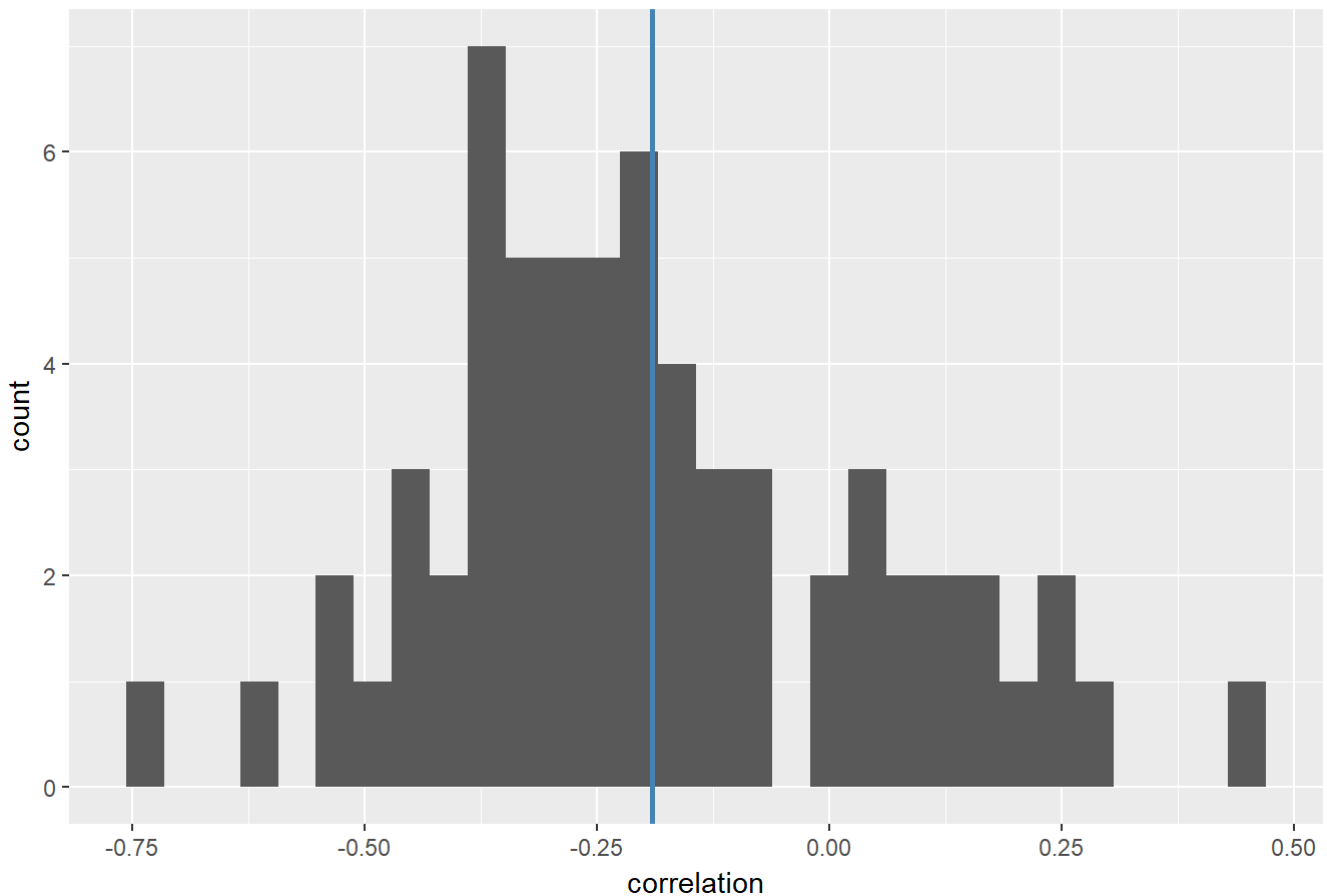
```
whole_ratio <- cor(debt$ratio,debt$growth)
whole_ratio
```

```
## [1] -0.199468
```

```
cor_function <- function(df){
  return(cor(df$ratio,df$growth))
}
year_cor <- ddply(debt,.(debt$Year),cor_function)
colnames(year_cor) <- c("year","correlation")
ggplot(year_cor)+
  geom_histogram(aes(x=correlation))+
  geom_vline(xintercept = -0.1906,lwd=1,col="steelblue")+
  labs(title="histogram of the correlation")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

histogram of the correlation



- f. (3 points) Some economists claim that high levels of government debt cause slower growth. Other economists claim that low economic growth leads to higher levels of government debt. The data file, as given, lets us relate this year's debt to this year's growth rate; to check these claims, we need to relate current debt to future growth. Create a new dataframe that contains all the rows of `debt` for France. It should have 54 rows and 4 columns. Note that some years are missing from the middle of this data set.

```
# Your answer to question 2.f here.
debt_france <- debt%>%
  filter(Country=='France')
dim(debt_france)
```

```
## [1] 54 4
```

- g. Create a new column in your dataframe for France created in 2.f, labeled `next.growth`, which gives next year's growth *if* the next year is in the data frame, or `NA` if the next year is missing. Do this in two steps.
1. (7 points) First write a function `n.growth()` that takes in two arguments, a year and a dataframe (that has the same columns as `debt` but rows only corresponding to a single country), and outputs the proper next growth value for that year and that dataframe (i.e. it gives next year's growth if the next year is in the input dataframe, or `NA` if the next year is missing).

```
# Your answer to question 2.g.1 here.
n.growth <- function(year,df){
  growth_next_year <- c()

  if(as.numeric(year+1) %in% df$Year){
    growth_next_year <- df$growth[df$Year==(year+1)]
  }else{
    growth_next_year <- NA
  }
  return(growth_next_year)
}
```

2. (5 points) Next use `n.growth()` and one of the functions from the `apply` family to create the `next.growth` column in the dataframe. ( `next.growth` for 1971 should be 5.886, but for 1972 it should be NA .)

```
# Your answer to question 2.g.2 here.
n.growth(year = 1971,debt_france)
```

```
## [1] 5.885827
```

```
n.growth(year=1972,debt_france)
```

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

```
next_year_growth <- sapply(debt_france$Year,n.growth,df=debt_france)
names(next_year_growth) <- debt_france$Year
next_year_growth
```

```
##      1950      1951      1952      1953      1954      1955
## 6.1349693 2.6274304 2.9185868 4.8258706 5.7902231 5.0246747
##      1956      1957      1958      1959      1960      1961
## 6.0230671 2.5382756 2.8290766 7.0691632 5.4960742 6.6305819
##      1962      1963      1964      1965      1966      1967
## 5.4885787 6.4962406 4.6879413 5.3142703 4.5850410 4.3595396
##      1968      1969      1970      1971      1972      1979
## 6.9467261 5.7932851 5.3723294 5.8858268      NA 1.0716383
##      1980      1981      1982      1983      1984      1985
## 1.5475053 0.1249935 3.8668937 5.8941037 5.3540295 4.0375523
##      1986      1987      1988      1989      1990      1991
## 1.7800135 -0.1728488 0.9975585 1.9275140 3.1046290 3.5231050
##      1992      1993      1994      1995      1996      1997
## 2.7866371 5.0513843 4.1863491 5.0997663 5.3926403 2.6827835
##      1998      1999      2000      2001      2002      2003
## 2.0257706 3.2535680 1.9901541 1.5020556 1.0135240 3.8641264
##      2004      2005      2006      2007      2008      2009
## 2.7392074 2.2808885 3.1340918 2.1321303 -1.9066761      NA
```

- h. (8 points) Add a `next.growth` column, as in 2.g, to the whole of the `debt` data frame. Make sure that you do not accidentally put the first growth value for one country as the `next.growth` value for another. (The `next.growth` for France in 2009 should be NA , not 9.167 .) Hints: Write a function to encapsulate what you did in 2.f, and apply it using `ddply()` .

```

# Your answer to question 2.g here.
newfunction <- function(df){
  next_year_growth <- c()
  for(i in df$Year){
    growth_next_year <- n.growth(i,df)
    next_year_growth <- c(next_year_growth,growth_next_year)
  }
  return(next_year_growth)
}

split_country <- split(debt,debt$Country)
result <- unlist(sapply(split_country,FUN = newfunction))
debt$next.growth <- result
head(debt,5)

```

```

##      Country Year    growth    ratio next.growth
## 1 Australia 1946 -3.557951 190.4191    2.459475
## 2 Australia 1947  2.459475 177.3214    6.437534
## 3 Australia 1948  6.437534 148.9298    6.611994
## 4 Australia 1949  6.611994 125.8287    6.920201
## 5 Australia 1950  6.920201 109.8094    4.272612

```