

MATH3426 - Midterm Project

Wilbert Caine

```
#change seed 20584260
set.seed(20387844)
data <- read.csv("Big_Five_Personality_small.csv")
head(data)

##   X EXT1 EXT2 EST1 EST2 AGR1 AGR2 CSN1 CSN2          dateload country
## 1 1     1     5     4     3     2     5     4     4 2016-03-03 02:11:06    AU
## 2 2     1     5     5     1     1     3     4     2 2016-03-03 02:22:56    AU
## 3 3     1     3     3     2     5     4     4     4 2016-03-03 02:26:37    AU
## 4 4     3     4     1     2     1     4     4     2 2016-03-03 02:43:01    AU
## 5 5     4     1     3     4     1     4     4     4 2016-03-03 03:11:20    AU
## 6 6     4     2     1     5     1     5     5     2 2016-03-03 03:26:32    AU
##   lat_appx_lots_of_err long_appx_lots_of_err
## 1           -37.9333            145.2333
## 2            -27.0000            133.0000
## 3           -31.6448            152.7946
## 4            -27.0000            133.0000
## 5            -27.0000            133.0000
## 6           -37.8139            144.9634
```

Part I. Simple Random Sampling

We first estimate the population mean of variable “EXT1”.

1. Find the population mean and population total of “EXT1”.

```
mu_ext1 <- mean(data$EXT1)
N <- nrow(data)
tau_ext1 <- mu_ext1*N
print(list("population mean of "EXT1""= mu_ext1, "population total of "EXT1""= tau_ext1))

## $`population mean of "EXT1"`
## [1] 2.64419
##
## $`population total of "EXT1"`
## [1] 886938
```

2. Use simple random sampling with $n=1500$ to estimate the population mean and total of “EXT1” respectively. Construct 96% confidence interval for the population mean and total of “EXT1”. What are the error bounds for these estimates?

```
n <- 15e2
N <- nrow(data)
alpha <- 0.04
sample_ext1 <- sample(data$EXT1, n)
```

```

ybar_ext1 <- mean(sample_ext1)
s_square_ext1 <- var(sample_ext1)
varhat_ybar_ext1 <- (s_square_ext1/n)*(1-n/N)
errorbound_muhat_ext1 = qnorm(1-alpha/2)*sqrt(varhat_ybar_ext1)

tauhat_ext1 <- N*ybar_ext1
varhat_tauhat_ext1 <- (N^2)*varhat_ybar_ext1
errorbound_tauhat_ext1 = qnorm(1-alpha/2)*sqrt(varhat_tauhat_ext1)

print(list("estimate of population mean"= ybar_ext1, "estimate of population mean"= tauhat_ext1,
          "error bound for estimate of population mean"= errorbound_muhat_ext1, "error bound for estimate of population mean"= errorbound_tauhat_ext1,
          "CI.lower for mu"= ybar_ext1-errorbound_muhat_ext1, "CI.lower for tau"= tauhat_ext1-errorbound_tauhat_ext1,
          "CI.upper for mu"= ybar_ext1+errorbound_muhat_ext1, "CI.upper for tau"= tauhat_ext1+errorbound_tauhat_ext1))

## $`estimate of population mean`
## [1] 2.702
##
## $`estimate of population mean`
## [1] 906329.2
##
## $`error bound for estimate of population mean`
## [1] 0.06672092
##
## $`error bound for estimate of population mean`
## [1] 22380.13
##
## $`CI.lower for mu`
## [1] 2.635279
##
## $`CI.lower for tau`
## [1] 883949
##
## $`CI.upper for mu`
## [1] 2.768721
##
## $`CI.upper for tau`
## [1] 928709.3

```

3. Repeat step 2 for 100 times and calculate the relative frequencies that the constructed confidence intervals contain the population mean and total, respectively. Are the relative frequencies close to 96%?

```

CI_cal_mu_ext1 <- function(sample, n, N){
  ybar_ext1 <- mean(sample)
  varhat_ybar_ext1 <- (var(sample)/n)*(1-n/N)
  errorbound_muhat_ext1 = qnorm(1-alpha/2)*sqrt(varhat_ybar_ext1)
  list("CI.lower for mu"= ybar_ext1-errorbound_muhat_ext1,
       "CI.upper for mu"= ybar_ext1+errorbound_muhat_ext1,
       "error bound for mu"= errorbound_muhat_ext1)
}

CI_cal_tau_ext1 <- function(sample, n, N){
  ybar_ext1 <- mean(sample)
  varhat_ybar_ext1 <- (var(sample)/n)*(1-n/N)

  tauhat_ext1 <- N*ybar_ext1
}

```

```

varhat_tauhat_ext1 <- (N^2)*varhat_ybar_ext1
errorbound_tauhat_ext1 = qnorm(1-alpha/2)*sqrt(varhat_tauhat_ext1)
list("CI.lower for tau"= tauhat_ext1-errorbound_tauhat_ext1,
     "CI.upper for tau"= tauhat_ext1+errorbound_tauhat_ext1,
     "error bound for tau"= errorbound_tauhat_ext1)
}
mrep <- 1e2
temp_mu_ext1 <- logical(mrep)
temp_tau_ext1 <- logical(mrep)
for (i in 1:mrep){
  sample_ext1 <- sample(data$EXT1, n)
  CI_mu <- CI_cal_mu_ext1(sample_ext1,n,N)
  temp_mu_ext1[i] <- (CI_mu[[1]]<=mu_ext1)*(CI_mu[[2]]>=mu_ext1)
  CI_tau <- CI_cal_tau_ext1(sample_ext1,n,N)
  temp_tau_ext1[i] <- (CI_tau[[1]]<=tau_ext1)*(CI_tau[[2]]>=tau_ext1)
}
print(list("relative frequencies that the constructed confidence intervals contain the population mean"=
           "relative frequencies that the constructed confidence intervals contain the population total"))

## $`relative frequencies that the constructed confidence intervals contain the population mean`#
## [1] 0.92
##
## $`relative frequencies that the constructed confidence intervals contain the population total`#
## [1] 0.92

```

4. Use the sample variance from step 2, calculate the the sample size lower bound for estimating the population mean of “EXT1” with an error bound 0.015, at the significance level 0.05.

```

alpha <-0.05
V_ext1 <- (0.015/qnorm(1-alpha/2))^2
sample_size_lower_ext1 <- (N*s_square_ext1)/(N*V_ext1+s_square_ext1)
sample_size_lower_ext1

```

```
## [1] 25117.59
```

5. Use the sample size lower bound from step 4, random sample this number of observations and calculate the error bound at the significance level 0.05, is it close to 0.015? Repeat this process 100 times and calculate the relative frequency that the error bound is smaller than 0.015.

```

n <- ceiling(sample_size_lower_ext1)
alpha <- 0.05
sample_ext1 <- sample(data$EXT1, n)
s_square_ext1 <- var(sample_ext1)
varhat_ybar_ext1 <- (s_square_ext1/n)*(1-n/N)
errorbound_muhat_ext1 = qnorm(1-alpha/2)*sqrt(varhat_ybar_ext1)
print(list("error bound for random sample with the sample size lower bound from step 4"= errorbound_muhat_ext1))

```

```
## $`error bound for random sample with the sample size lower bound from step 4`#
## [1] 0.01492682
```

```

mrep <- 1e2
temp_mu_ext1 <- logical(mrep)
for (i in 1:mrep){
  sample_ext1 <- sample(data$EXT1, n)
  s_square_ext1 <- var(sample_ext1)
  varhat_ybar_ext1 <- (s_square_ext1/n)*(1-n/N)
  errorbound_muhat_ext1 = qnorm(1-alpha/2)*sqrt(varhat_ybar_ext1)
}
```

```

    temp_mu_ext1[i] <- errorbound_muhat_ext1<0.015
}

print(list("relative frequency that the error bound is smaller than 0.015"= mean(temp_mu_ext1)))

## $`relative frequency that the error bound is smaller than 0.015`-
## [1] 0.88

```

We then estimate the population proportion of answer “5” for variable “EXT2”.

6. For the variable “EXT2”, what is the population proportion of answer “5”?

```

N <- nrow(data)
p_ext2.5 <- mean(data$EXT2== 5)
print(list("population proportion of answer "5" for variable "EXT2""= p_ext2.5))

```

```

## $`population proportion of answer "5" for variable "EXT2`-
## [1] 0.1170113

```

7. Use simple random sampling with $n=5000$ to estimate the population proportion of answer “5” for variable “EXT2”. What is the confidence interval and error bound of your estimate at significance level 0.05 ?

```

n <- 5e3
N <- nrow(data)
alpha <- 0.05
sample_ext2 <- sample(data$EXT2, n)

phat_ext2.5 <- mean(sample_ext2==5)
varhat_phat_ext2.5 <- (1-n/N)*(phat_ext2.5*(1-phat_ext2.5))/(n-1)
errorbound_phat_ext2.5 = qnorm(1-alpha/2)*sqrt(varhat_phat_ext2.5)

print(list("estimate the population proportion of answer "5" for variable "EXT2""= phat_ext2.5,
          "error bound for estimate the population proportion"= errorbound_phat_ext2.5,
          "CI.lower for p"= phat_ext2.5-errorbound_phat_ext2.5,
          "CI.upper for p"= phat_ext2.5+errorbound_phat_ext2.5))

```

```

## $`estimate the population proportion of answer "5" for variable "EXT2`-
## [1] 0.1132
##
## $`error bound for estimate the population proportion`-
## [1] 0.008717287
##
## $`CI.lower for p`-
## [1] 0.1044827
##
## $`CI.upper for p`-
## [1] 0.1219173

```

8. Repeat step 7 for 100 times and calculate the relative frequencies that the confidence intervals contain the population proportion of answer “5” for the variable “EXT2”.

```

CI_cal_p_ext2.5 <- function(sample, n, N){
  phat_ext2.5 <- mean(sample==5)
  varhat_phat_ext2.5 <- (1-n/N)*(phat_ext2.5*(1-phat_ext2.5))/(n-1)
  errorbound_phat_ext2.5 = qnorm(1-alpha/2)*sqrt(varhat_phat_ext2.5)
}

```

```

list("CI.lower for p"= phat_ext2.5-errorbound_phat_ext2.5,
     "CI.upper for p"= phat_ext2.5+errorbound_phat_ext2.5,
     "error bound for p"= errorbound_phat_ext2.5)
}
mrep <- 1e2
temp_p_ext2.5 <- logical(mrep)
for (i in 1:mrep){
  sample_ext2.5 <- sample(data$EXT2, n)
  CI_p <- CI_cal_p_ext2.5(sample_ext2.5,n,N)
  temp_p_ext2.5[i] <- (CI_p[[1]]<=p_ext2.5)*(CI_p[[2]]>=p_ext2.5)
}
print(list("relative frequencies that the constructed confidence intervals contain the population proportion`

## $`relative frequencies that the constructed confidence intervals contain the population proportion`

## [1] 0.93

9. Use the sample variance from step 7, calculate the sample size lower bound for estimating the population proportion of answer “5” of variable “EXT2” with an error bound of 0.02, at the significance level 0.04.

alpha <- 0.04
V_ext2 <- (0.02/qnorm(1-alpha/2))^2
sample_size_lower_ext2.5 <- (N*phat_ext2.5*(1-phat_ext2.5))/(N*V_ext2+phat_ext2.5*(1-phat_ext2.5))
sample_size_lower_ext2.5

## [1] 1055.209

```

Part II. Stratified Simple Random Sampling

By looking at the country code of all participants in the provided dataset, the participants come from 6 different countries “AU”, “CA”, “DE”, “GB”, “US”, “PH”. In this part, we implement stratified SRS to estimate population parameters for the variables “EST1” and “EST2”.

We first estimate the population mean of variable “EST1”.

1. Find the population sizes of the 6 strata and population mean of “EST1”.
2. To estimate the population mean of “EST1”, use simple random sampling to sample 400 from ‘AU’, 600 from ‘CA’, 140 from ‘DE’, 660 from ‘GB’, 200 from ‘PH’ and 1200 from ‘US’. Compute your point estimate and error bound at significance level 0.05. Construct the confidence interval.
3. Repeat step 2 for 100 times and calculate the relative frequencies that the constructed confidence intervals contain the population mean of “EST1”. Is the relative frequency close to 95%?
4. Use the sample variances and allocation weights from step 2, calculate the the sample size lower bound for estimating the population mean of “EST1” with an error bound 0.02, at the significance level 0.05.
5. Use the sample size lower bound from step 4 and allocating weights as step 2, randomly sample these numbers of observations and calculate the error bound at the significance level 0.05, is it close to 0.02? Repeat this process 100 times and calculate the relative frequency that the error bound is smaller than 0.02.

We then estimate the population proportion of answer “5” for variable “EST2”.

6. Find population proportion of answer “5” for variable “EST2”.

7. Use stratified simple random sampling with $n=5000$ and Neyman allocation to estimate the population proportion of answer “5” for the variable “EST2”. What is the confidence interval and error bound of your estimate at significance level 0.05 ?
8. Repeat step 7 for 100 times and calculate the relative frequencies that the confidence intervals contain the population proportion of answer “5” for the variable “EST2”.
9. Use the sample variances from step 7 with Neyman allocation, calculate the sample size lower bound for estimating the population proportion of answer “5” of variable “EST2” with an error bound of 0.02, at the significance level 0.05.
10. Repeat step 9, but use proportional allocation.
11. Use the sample size obtained from step 10 and with proportional allocation, randomly sample this number of observations from the population by stratified SRS and calculate the error bound. Repeat the process 100 times and calculate the relative frequencies that the error bound is less than 0.02.

Part III. Systematic Random Sampling

Now, we apply systematic random sampling to estimate the population mean of variable “AGR1”.

We estimate the population mean of variable “AGR1”.

1. Find the population mean of variable “AGR1”.
2. Choose $k=100$ and do a systematic random sampling. Calculate the point estimate and its error bound at significance level 0.05. Use half-sample estimator to estimate the required variance.
3. Choose $k=100$ and do a systematic random sampling. Calculate the point estimate and its error bound at significance level 0.05. Use successive differences estimator to estimate the required variance.
4. Choose $k=100$ and do a systematic random sampling. Calculate the point estimate and its error bound at significance level 0.05. Use 10 repeated system samples to estimate the required variance.
5. Use step 4 to construct 95% confidence interval of population mean of “AGR1”. Repeat 100 times and calculate the relative frequency of those intervals containing the population mean.