



# University of St.Gallen

## Classification

**University of St. Gallen**  
School of Management, Economics, Law,  
Social Sciences, International Affairs  
and Computer Science

### **Assignment 4**

Data Analytics I: Predictive Econometrics  
Prof. Jana Mareckova

submitted by

**Cyril Janak, 16-611-287**  
**Jonas Husmann, 16-610-917**  
**Niklas Kampe, 16-611-618**  
**Robin Scherrer, 18-617-969**

22.12.2021

# Contents

Requirements	1
Exercise 1	1
Exercise 2	1
Exercise 3	2
Exercise 4	2
Exercise 5	3
Exercise 6	5

## Requirements

To solve the following tasks, the required libraries and the data sets are loaded first.

```
library(rpart)
library(rpart.plot)
library(dplyr)
library(stringr)
library(ggplot2)
library(qpcR)

load("GHA/drugs.RData")
```

## Exercise 1

The share of males who consume soft drugs is ~29.18%

```
(m_s_drug <- (nrow(drugs[drugs$Gender=="male" & drugs$Soft_Drug==T,]) /
  nrow(drugs[drugs$Gender=="male",]) * 100) %>%
  round(., digits = 2) %>%
  paste0(., "%"))
```

```
## [1] "29.18%"
```

## Exercise 2

The difference between the share of male and female hard drug consumers is ~2.74%

```
m_h_drug <- nrow(drugs[drugs$Gender=="male" & drugs$Hard_Drug==T,]) /
  nrow(drugs[drugs$Gender=="male",])

f_h_drug <- nrow(drugs[drugs$Gender=="female" & drugs$Hard_Drug==T,]) /
  nrow(drugs[drugs$Gender=="female",])

(diff_h_drug <- ((m_h_drug - f_h_drug) * 100) %>%
  round(., digits = 2) %>%
  paste0(., "%"))
```

```
## [1] "2.74%"
```

### Exercise 3

From the shares of soft drug consumption for each age group, one can observe that only 16-17 year-olds consume soft drugs. Therefore, the consumption of soft drugs is decreasing in age, but not strictly as the groups of 18-19 and 20-24 year-olds are not consuming any soft drugs at all.

```
share_softdrugs_16_17 <- round((nrow(drugs[drugs$Age=="16-17 years" &
                                         drugs$Soft_Drug==T,]) /
  nrow(drugs[drugs$Age=="16-17 years",]))*100, digits = 2)
share_softdrugs_18_19 <- round((nrow(drugs[drugs$Age=="18-19 years" &
                                         drugs$Soft_Drug==T,]) /
  nrow(drugs[drugs$Age=="18-19 years",]))*100, digits = 2)
share_softdrugs_20_24 <- round((nrow(drugs[drugs$Age=="20-24 years" &
                                         drugs$Soft_Drug==T,]) /
  nrow(drugs[drugs$Age=="20-24 years",]))*100, digits = 2)

(shares_softdrugs <- data.frame(
  age = c("16-17 Years", "18-19 Years", "20-24 Years"),
  share = c(share_softdrugs_16_17, share_softdrugs_18_19, share_softdrugs_20_24)))
```

```
##           age share
## 1 16-17 Years  48.5
## 2 18-19 Years   0.0
## 3 20-24 Years   0.0
```

### Exercise 4

The chi-squared test results in a X-squared statistic of 9.40 at a p-value of 0.025. Hence, the hypothesis of independence is rejected ( $0.025 < 0.05$ ) and the earnings range and soft drug consumption are indeed dependent at a confidence interval of 5%.

```
drugs_table <- table(drugs$Earning, drugs$Soft_Drug)
chi_squared <- chisq.test(drugs_table)
(statistics <- chi_squared$statistic)
```

```
## X-squared
##  9.401385
```

```
(p_value <- chi_squared$p.value)
```

```
## [1] 0.02440394
```

## Exercise 5

```
#create random vector
v5 <- c(1:500) #Vector, 1-500

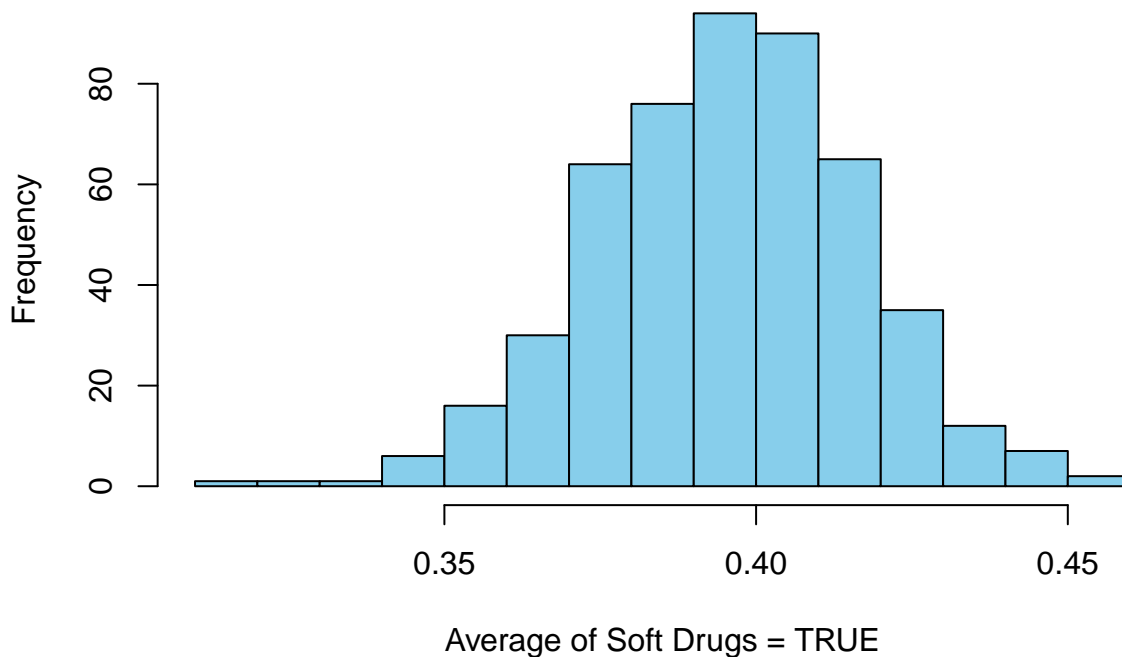
sd_no <-vector() #create empty vector for Soft Drugs = FALSE
sd_yes <-vector() #create empty vector for Soft Drugs = TRUE

for (value in v5) {
  #create 500 subsamples with 500 observations
  subsample <-drugs[sample(nrow(drugs), 500), ]
  #count the number of Soft Drugs = TRUE
  soft_drugs_yes <- sum(str_count(subsample$Soft_Drug, "TRUE"))
  #count the number of Soft Drugs = FALSE
  soft_drugs_no <- sum(str_count(subsample$Soft_Drug, "FALSE"))
  #Calculate the percentage of Soft Drugs = TRUE
  average_yes = (soft_drugs_yes / (soft_drugs_yes + soft_drugs_no))
  #Calculate the percentage of Soft Drugs = FALSE
  average_no = (soft_drugs_no / (soft_drugs_yes + soft_drugs_no))
  #Append value of TRUE to empty vector
  sd_yes = append(sd_yes, average_yes)
  #Append value of FALSE to empty vector
  sd_no = append(sd_no, average_no)
}

#Since True=1 and False=0, the Histogram showing sd_yes shows the average of the
#logical variable

hist(sd_yes,
      main ="Histogram of Soft Drugs Used",
      xlab = "Average of Soft Drugs = TRUE",
      col ="skyblue",
)
```

## Histogram of Soft Drugs Used



```
#Average of full sample
soft_drugs_yes_full <- sum(str_count(drugs$Soft_Drug, "TRUE"))
soft_drugs_no_full <-sum(str_count(drugs$Soft_Drug, "FALSE"))
average_yes_full <- (soft_drugs_yes_full / (soft_drugs_yes_full+soft_drugs_no_full))

mean_full = round(mean(sd_yes), 5)
mean_sub = round(average_yes_full, 5)
diff_mean = round(mean_sub - mean_full, 5)

print(paste0("The average in the subsample is: ", mean_sub))
```

```
## [1] "The average in the subsample is: 0.39716"
```

```
print(paste0("The average in the full sample is: ", mean_full))
```

```
## [1] "The average in the full sample is: 0.39641"
```

```
print(paste0("The difference is approx: ", diff_mean))
```

```
## [1] "The difference is approx: 0.00075"
```

## Exercise 6

When keeping the numbers of draws fixed at 500 one can see a higher density around the mean of the whole sample (see exercise 5) with an increasing sample size. If the sample size is kept fixed at 500 and the number of draws (100, 500, 2500) is varied one can observe that the curve is smoother for higher number of draws.

```
set.seed(11111)
#adjusting the soft_drug column to be a logical to make it easier work with it
drugs$Soft_Drug = as.logical(drugs$Soft_Drug)
#defining sample sizes and number of runs
sample_sizes = c(100, 500, 2500)
N_runs = c(100, 500, 2500)

#creating the function that we will use
# within the function first create the samples which will then be used for the
# mean calculations
subsample_mean = function(data, sample_size, draws) {
  samples = matrix(replicate(draws, sample(x = c(1:length(data)),
                                           size = sample_size, replace = FALSE)),
                  nrow = sample_size)
  means = rep(NA, draws)
  for(i in 1:draws){
    means[i] = mean(data[c(samples[,i])])
  }
  return(means)
}

#means_list is used to store the retrieved means
means_list = list()
#creating empty sample_size to use for the j in the for loop
sample_size = list()
#creating counter for the for loop
counter = 0
for (i in N_runs) {
  counter = counter + 1
  t = matrix(data = NA, nrow = N_runs[counter], ncol = 3)
  k=0
  for (j in sample_sizes) {
    k = k+1
    #storing the function results
    t[, k] = subsample_mean(data = drugs$Soft_Drug,
                           sample_size = j, draws = i)
  }
  means_list[[counter]] = t
}
```

*#fixing draws at 500*

```
draws500 = data.frame(means_list[[2]])
colnames(draws500) = c("subsample_a", "subsample_b", "subsample_c")
tail(draws500)
```

```
##      subsample_a subsample_b subsample_c
## 495         0.33         0.426         0.3808
## 496         0.38         0.396         0.3856
## 497         0.34         0.388         0.4096
## 498         0.41         0.378         0.3960
## 499         0.33         0.382         0.3976
## 500         0.36         0.428         0.3956
```

*#fixing subsample size at 500*

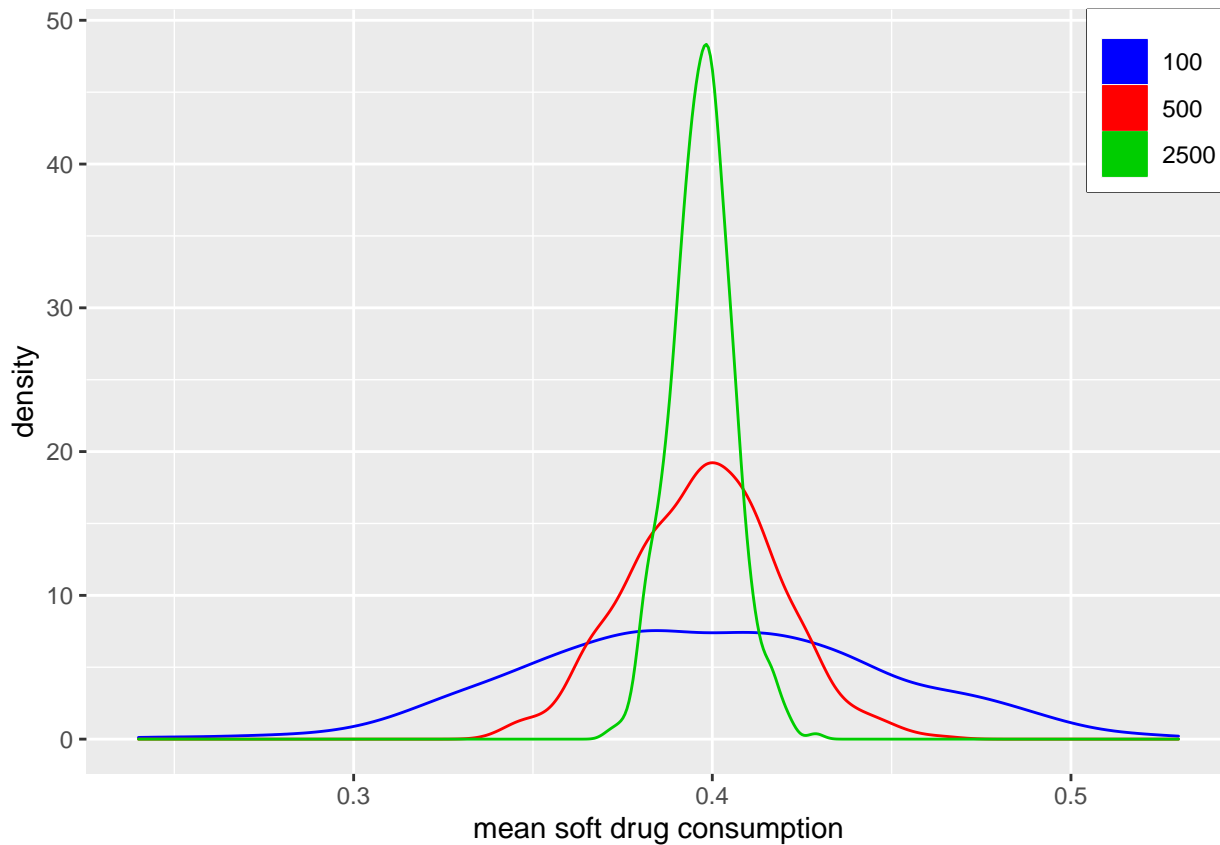
```
subsample500 = data.frame(qpcR::cbind.na(means_list[[1]][,2],
                                          means_list[[2]][,2],
                                          means_list[[3]][,2]))
colnames(subsample500) = c("draws_a", "draws_b", "draws_c")
tail(subsample500)
```

```
##      draws_a draws_b draws_c
## 2495        NA        NA  0.434
## 2496        NA        NA  0.430
## 2497        NA        NA  0.426
## 2498        NA        NA  0.406
## 2499        NA        NA  0.430
## 2500        NA        NA  0.414
```

*#plotting for the 500 draws with the 3 different subsample sizes*

```
ggplot(data = draws500) +
  geom_density(aes(x = subsample_a, color = "100")) +
  geom_density(aes(x = subsample_b, color = "500")) +
  geom_density(aes(x = subsample_c, color = "2500")) +
  scale_color_manual("", values = c("100" = "blue",
                                    "500" = "red",
                                    "2500" = "green3")) +
  scale_fill_manual("", values = c("100" = "blue",
                                    "500" = "red",
                                    "2500" = "green3")) +
  xlab("mean soft drug consumption") +
  ylab("density") +
  theme(legend.title = element_blank(),
        legend.position = c(1, 1),
        legend.justification = c("right", "top"),
        legend.box.background = element_rect())
```





*#plotting for the 500 subsample size with 3 different draw numbers*

```
ggplot(data = subsample500) +
  geom_density(aes(x = draws_a, color = "100")) +
  geom_density(aes(x = draws_b, color = "500")) +
  geom_density(aes(x = draws_c, color = "2500")) +
  scale_color_manual("", values = c("100" = "blue" ,
                                    "500" = "red",
                                    "2500" = "green3")) +
  scale_fill_manual("", values = c("100" = "blue",
                                   "500" = "red",
                                   "2500" = "green3")) +
  xlab("mean soft drug consumption") +
  ylab("density") +
  theme(legend.title = element_blank(),
        legend.position = c(1, 1),
        legend.justification = c("right", "top"),
        legend.box.background = element_rect())
```

