# 732A90: Computational Statistics

Computer lab5 - Group11

*Sofie Jörgensen, Oriol Garrobé Guilera, David Hrabovszki*

*26 February 2020*

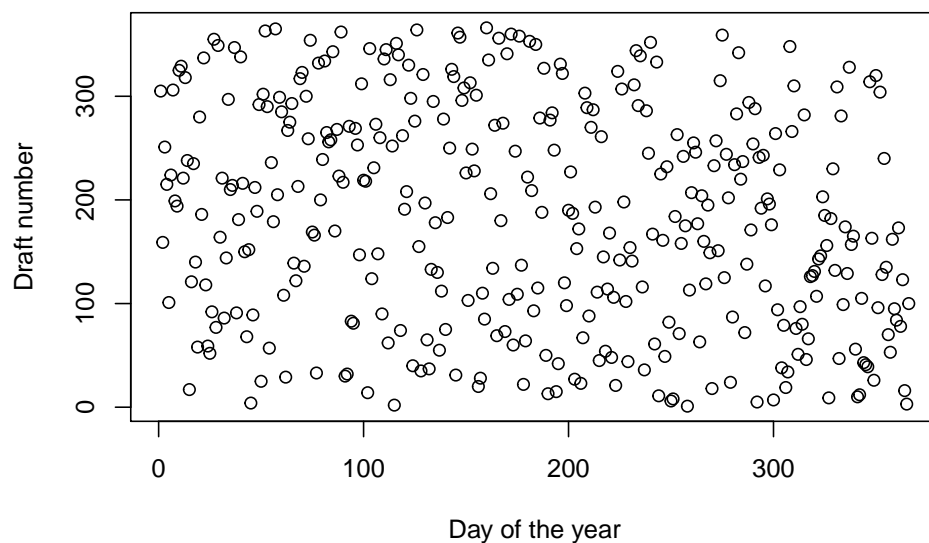## Question 1: Hypothesis testing

**1.**



Figure 1: Scatterplot of Draft number versus Day of the year

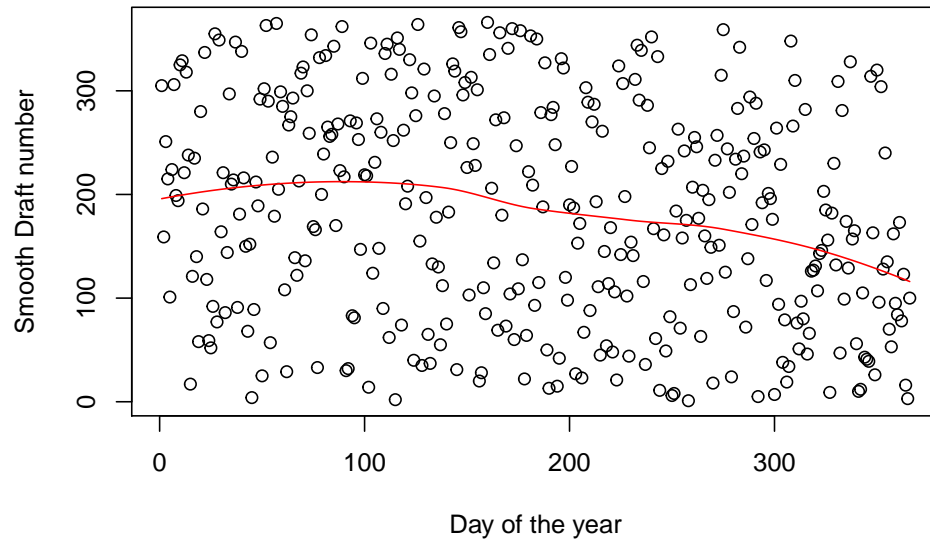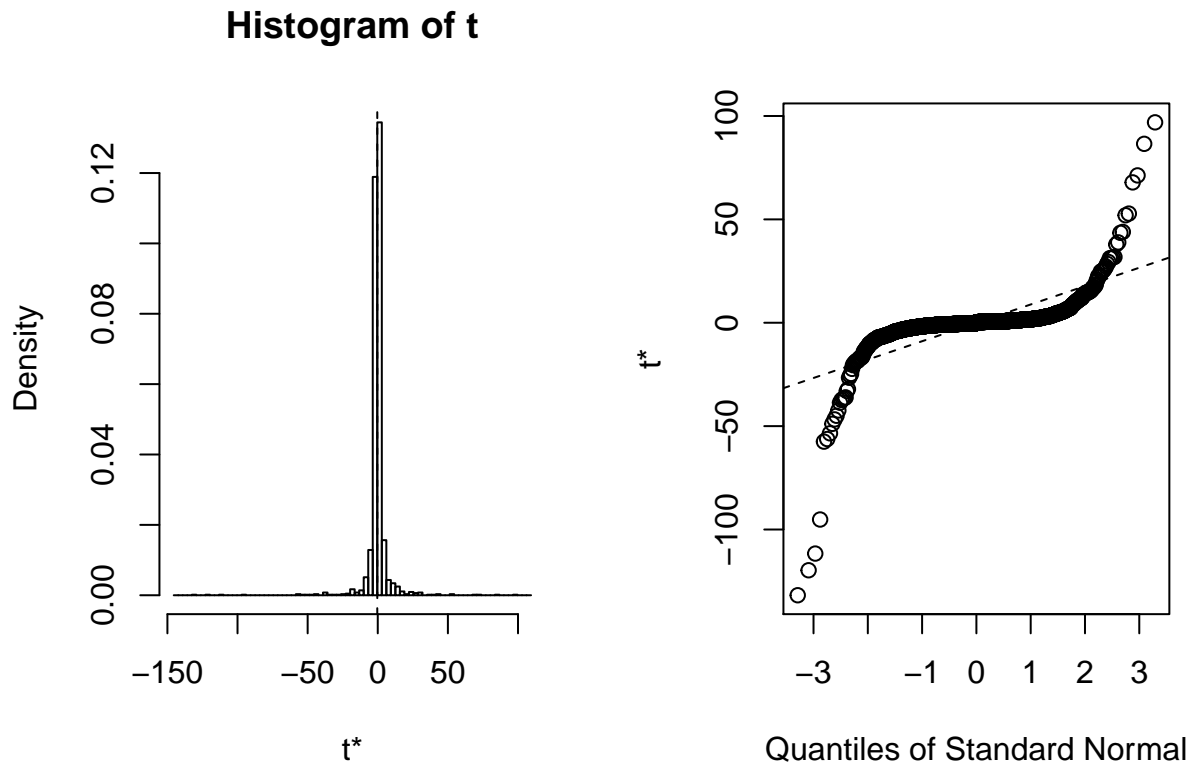The lottery looks random as there is no clear correlation between the variables.

**2.**



Figure 2: Scatterplot of Smooth Darft number versus Day of the year

It does not look as random as before. The days at he beginning of the year have on average larger draft number that those at the end of the year.

**3.**

```
## [1] -0.3479163
```

There should not be a trend in the data as the value is not significantly greater than 0.

**Histogram of t**



```
## [1] 0.9425
```

**4.**

```
## [1] 0.9395
```

The p-value is large.

**5.**

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

```
##   [1] 1.000 1.000 1.000 1.000 1.000 0.990 0.960 0.955 0.920 0.905 0.825
##  [12] 0.800 0.825 0.790 0.750 0.680 0.655 0.705 0.700 0.645 0.670 0.695
##  [23] 0.660 0.655 0.615 0.645 0.615 0.625 0.615 0.625 0.670 0.655 0.575
##  [34] 0.535 0.535 0.545 0.540 0.575 0.565 0.575 0.575 0.580 0.550 0.525
##  [45] 0.540 0.510 0.510 0.480 0.480 0.490 0.480 0.500 0.510 0.520 0.525
##  [56] 0.530 0.555 0.545 0.535 0.535 0.525 0.520 0.530 0.535 0.545 0.525
##  [67] 0.510 0.510 0.500 0.495 0.490 0.490 0.495 0.500 0.495 0.495 0.480
##  [78] 0.470 0.465 0.470 0.475 0.485 0.480 0.485 0.475 0.475 0.490 0.490
##  [89] 0.500 0.510 0.510 0.510 0.510 0.500 0.500 0.495 0.515 0.530 0.525
```

---

# Question 2: Bootstrap, jackknife and confidence intervals

In this task we are going to estimate home prices in Albuquerque back in 1993, by using bootstrap, jackknife and confidence intervals. For this, we import the data set `prices1.csv` consisting of 110 observations, where

`Price` is the variable of interest.

**1.**

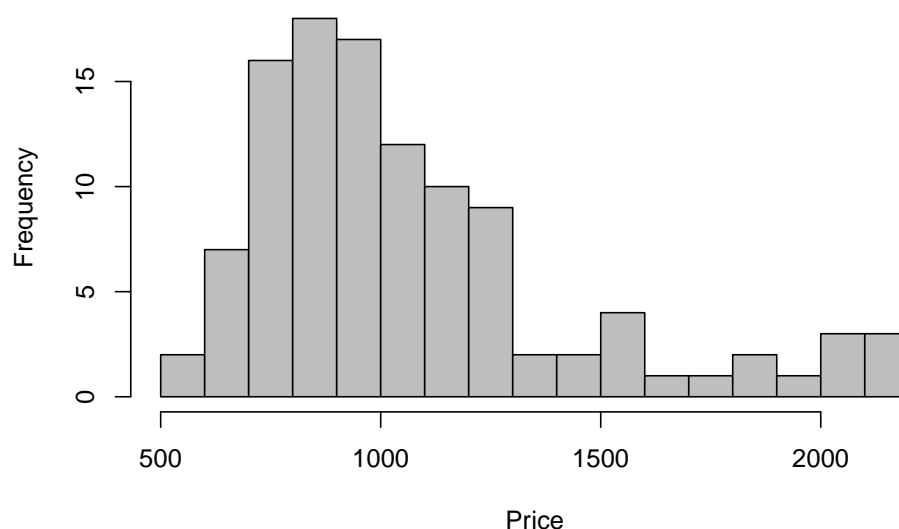First, we plot a histogram of the variable `Price`, to get an idea of its distribution.



Figure 3:   A histogram of the distribution of Price.

In Figure 3, we can observe that `Price` appears to be right-skewed, which reminds us of a Chi-squared distribution. The mean price of `Price` is computed to 1080.473.

**2.**

In the previous task, we presented a histogram of the price and computed the mean price. Now we are interested in estimating the distribution of the mean price by using bootstrap. The package `boot` provides appropriate functions for this purpose.

The following approach is being used: Draw $B = 2000$ bootstrap samples, i.e. resamples with replacement of size $n = 110$, and then compute the statistic, that is, the mean. Then we can form a bootstrap distribution that approximates the sampling distribution of the mean statistic.
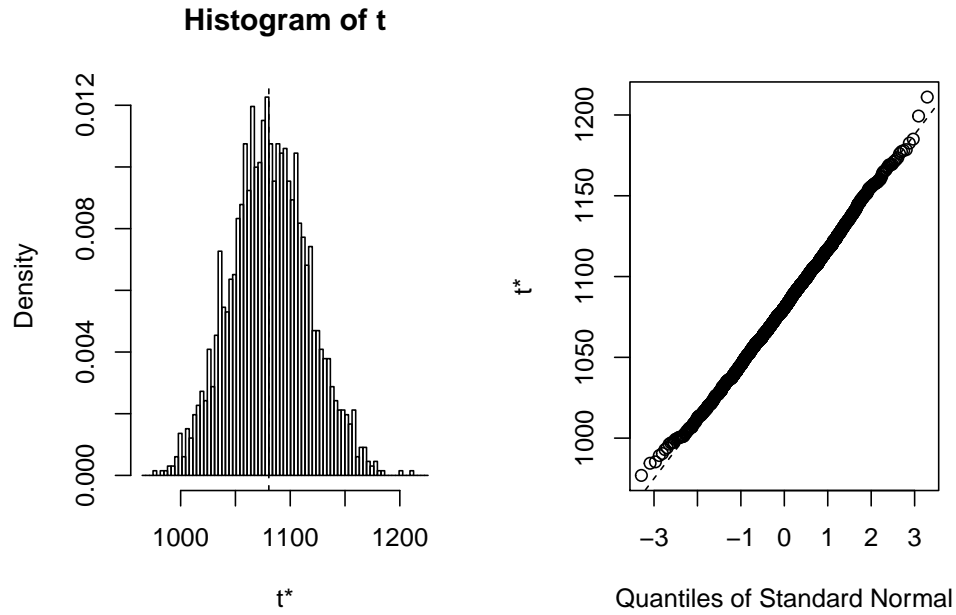
**Histogram of t**



Figure 4: Left plot: A histogram of the bootstrap distribution that approximates the sampling distribution of the mean statistic denoted by $t^*$. The dotted vertical line indicates the original mean value of the price. Right plot: QQ-plot of the bootstrap samples.
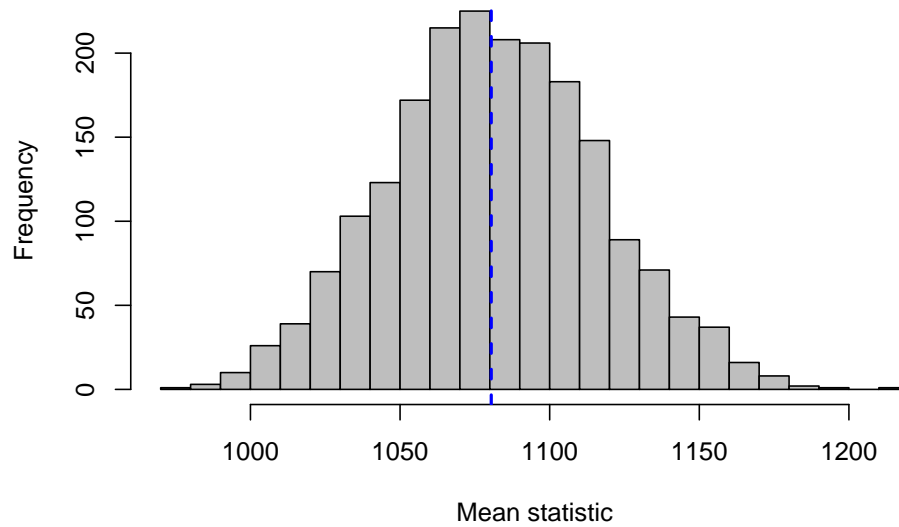


Figure 5: A histogram of the bootstrap distribution, with 2000 resamples, that approximates the sampling distribution of the mean statistic. The dotted vertical line indicates the original mean value of the price.

(Choose one of the histograms?)

The plot to the left in Figure 4 looks normally distributed. The right plot shows the QQ-plot of the bootstrap samples against standard normal quantiles, which follows a straight line, which also gives support for the distribution to be normal. The distribution of the mean price is estimated using bootstrap can be viewed in Figure 5. From the histogram we can notice that the distribution looks normally distributed.

Further, we will determine the bootstrap bias–correction and the variance of the mean price. The bias-corrected estimate is computed to 1079.523, by taking two times the original mean price minus the mean of the bootstrap sample. The bootstrap variance of the mean price equals 1349.957.

Thereafter we wish to compute a 95% confidence interval for the mean price using bootstrap percentile, bootstrap BCa, and first–order normal approximation, which are presented in Table 1.

| Type | 95% CI | Length |
|---|---|---|
| percentile | $(1011.761, 1155.064)$ | 143.3033 |
| BCa | $(1015.582, 1158.401)$ | 142.8193 |
| normal approx. | $(1008.088, 1152.495)$ | 144.407 |

Table 1: *The 95% confidence interval for the mean price using bootstrap percentile, bootstrap BCa, and first–order normal approximation.*

**3.**

Now let us consider an alternative way, the jackknive, of estimating the variance of the mean price. The jackknife estimation uses the leave-one-out method, which means that one observation of `Price` is omitted at each iteration. Since our data set consists of 110 observations, we will compute mean price of each sub-sample of size 109.

(Is this reasonable? Incorrect computations? But I try to explain the result anyway.) When estimating the variance of the mean price using the jackknife, we obtain a value of 12.23. This estimation is much lower compared to the estimated variance of the mean price, 1357.127, using bootstrap. Thus, it is clear that there is a big difference between the obtained estimations. One explanation is that only one observation differs between each sub-sample, and therefore the mean of each sub-sample will have a small variation among all the mean prices.

**4.**

Lastly, we are going to compare the confidence intervals obtained with respect to their length and the location of the estimated mean in these intervals. In Table 1, we can see that the intervals differ slightly depending on which type is being used, but the length of the interval is more similar. To illustrate the difference we present a histogram of the bootstrap distribution with its corresponding confidence interval.
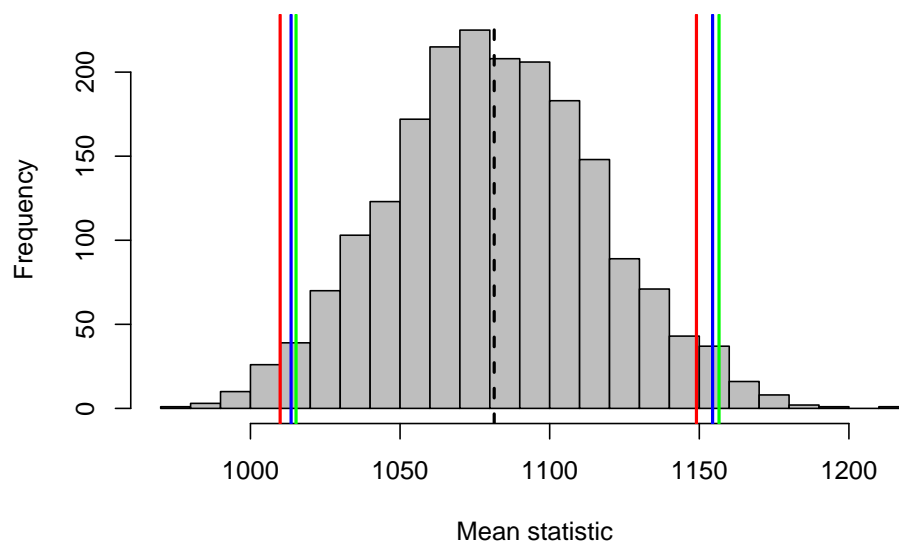
Figure 6: A histogram of the bootstrap distribution, with 2000 resamples, that approximates the sampling distribution of the mean statistic, with its corresponding 95 percent confidence interval. The dotted vertical line indicates the estimated mean price. The blue, green and red vertical lines corresponds to the confidence interval using bootstrap percentile, bootstrap BCa and first-order normal approximation, respectively.

The estimated mean is located around the center in these intervals. When we investigate the location even closer, the estimated mean is closest to the median of the interval using the first-order normal approximation, then comes the percentile, followed up by BCa.

---

# Appendix

```r
knitr::opts_chunk$set(echo = FALSE)
# R version
RNGversion('3.5.1')
# Packages
library("boot")
# 1.1

# Import the data set
data <- read.csv2("lottery.csv")

Y <- data$Draft_No
X <- data$Day_of_year

plot(X, Y, xlab = "Day of the year", ylab = "Draft number")

# 1.2
Y_loess <- loess(Draft_No ~ Day_of_year,data = data)
```

```r
smoothed <- predict(Y_loess)

plot(X, Y, xlab = "Day of the year", ylab = "Smooth Draft number")
lines(smoothed, x=X, col = "red")

# 1.3
Xb <- which.max(smoothed)
Xa <- which.min(smoothed)

T_test <- (smoothed[Xb] - smoothed[Xa])/(Xb - Xa)
T_test
# computing bootstrap samples
stat1 <- function(data, ind){
  data1 <- data[ind,]# extract bootstrap sample
  # Y <- data1$Draft_No
  # X <- data1$Day_of_year
  Y_loess <- loess(Draft_No ~ Day_of_year, data = data1)
  smoothed <- predict(Y_loess)

  Xb <- which.max(smoothed)
  Xa <- which.min(smoothed)
  T_test <- (smoothed[Xb] - smoothed[Xa])/(Xb - Xa)
  return(T_test)
}
# Seed
set.seed(1234567890)
res <- boot(data, stat1, R=2000) #make bootstrap
plot(res)
p_value <- sum(abs(res$t) >= abs(res$t0))/(res$R) # should it be (res$R + 1)
p_value # 0.9425

#1.4

permutation <- function(data, B) {
  # Seed
  set.seed(1234567890)

  stat <- numeric(B)
  n <- length(data[,1])
  for (b in 1:B) {
    id <- sample(n, replace = TRUE)
    Gb <- data[id,]
    stat[b] <- t_test(Gb)
  }
  stat0 <- t_test(data)
  return(sum(abs(stat)>abs(stat0))/B)
}

t_test <- function(data) {
  loess <- loess(Draft_No ~ Day_of_year,data = data)
  Y_hat <- predict(loess, data)

  Xb <- which.max(unname(Y_hat))
```

```r
  Xa <- which.min(unname(Y_hat))

  t_stat <- (unname(Y_hat[Xb]) - unname(Y_hat[Xa])) / (Xb - Xa)

  return(t_stat)
}

permutation(data, 2000)
#1.5

# a)

create_dataset <- function(X, alpha) {
  # Seed
  set.seed(1234567890)

  Y_set <- alpha*X+rnorm(366,183,10)
  Y_set[which(Y_set>366)] <- 366
  Y_set[which(Y_set<0)] <- 0
  dataset <- data.frame(Day_of_year=X, Draft_No=Y_set)
}


# b)
dataset <- create_dataset(X = X, alpha = 0.1)

permutation(data = dataset, B = 200)
# c)

alphas <- seq(from=0.2, to=10, by=0.1)
p_values <- numeric(length(alphas))
for (e in alphas) {
  dataset <- create_dataset(X = X, alpha = e)
  p_values[which(alphas==e)] <- permutation(data=dataset, B=200)
}
p_values
# 2.1

# Import the data set
data <- read.csv2("prices1.csv")

# Compute the mean price
price_mean <- mean(data$Price)

# Histogram of Price
hist(data$Price, breaks = 20, col = "grey", xlab = "Price", main = "")

# 2.2

# Seed
set.seed(123456789)

# Function that returns the mean statistic, that will be bootstrapped.
```

```r
# 1st argument: data set to bootstrap
# 2nd argument: index vector of the data set
mean_statistic <- function(d, i){
  data <- d[i,] # use all observations
  return(mean(data$Price))
}

# Bootstrap
price_boot <- boot(data, mean_statistic, R = 2000)

# Mean of bootstrap samples

boot_mean <- mean(price_boot$t)
plot(price_boot)

# Histogram of the bootstrap distribution
hist(price_boot$t[,1], main = "", xlab = "Mean statistic", col = 'grey', breaks = 20)
abline(v = price_mean, col = "blue", lwd = 2, lty = 2) # Original mean price

# Compute bootstrap estimated bias
bias <- boot_mean - price_mean

# Compute bias-correction, alt 1
price_mean - bias

# Compute bias-correction, alt 2
2*price_mean - boot_mean

#estimated variance
var(price_boot$t)
# 95% confidence interval for the mean price,
#   using bootstrap percentile, bootstrap BCa, and first-order normal approximation
ci_95 <- boot.ci(boot.out = price_boot, type = c("perc", "bca", "norm"))

# CI and their length
ci_perc <- ci_95$perc[ , c(4, 5)]
length_perc <- diff(ci_perc)

ci_bca <- ci_95$bca[ , c(4, 5)]
length_bca <- diff(ci_bca)


ci_norm <- ci_95$normal[ , c(2, 3)]
length_norm <- diff(ci_norm)
# Jackknife the mean
jack <- c()
for (i in seq_along(data$Price)){
  jack[i] <- mean(data$Price[-i])
}

# Variance of mean price
var(jack)
```

```r
# Median of CI
sum(ci_perc)/2
sum(ci_bca)/2
sum(ci_norm)/2
boot_mean
hist(price_boot$t[,1], main = "", xlab = "Mean statistic", col = 'grey', breaks = 20)
abline(v = boot_mean, lwd = 2, lty = 2) # estimated mean price
abline(v = ci_perc, col = 'blue',lwd = 2) # percentile CI
abline(v = ci_bca, col = 'green',lwd = 2) # BCa CI
abline(v = ci_norm, col = 'red',lwd = 2) # normal approx CI
```