

# Computational Statistics - Lab 05

*Annalena Erhard (anner218) and Maximilian Pfundstein (maxpf364)*

*2019-02-11*

## Contents

<b>1</b>	<b>Question 1: Hypothesis testing</b>	<b>1</b>
1.1	Task 1.1 . . . . .	1
1.2	Task 1.2 . . . . .	2
1.3	Task 1.3 . . . . .	3
1.4	Task 1.4 . . . . .	5
1.5	Task 2.5 . . . . .	5
<b>2</b>	<b>Question 2: Bootstrap, jackknife and confidence intervals</b>	<b>5</b>
2.1	Task 2.1 . . . . .	6
2.2	Task 2.2 . . . . .	6
2.3	Task 2.3 . . . . .	6
2.4	Task 2.4 . . . . .	6
<b>3</b>	<b>Source Code</b>	<b>6</b>

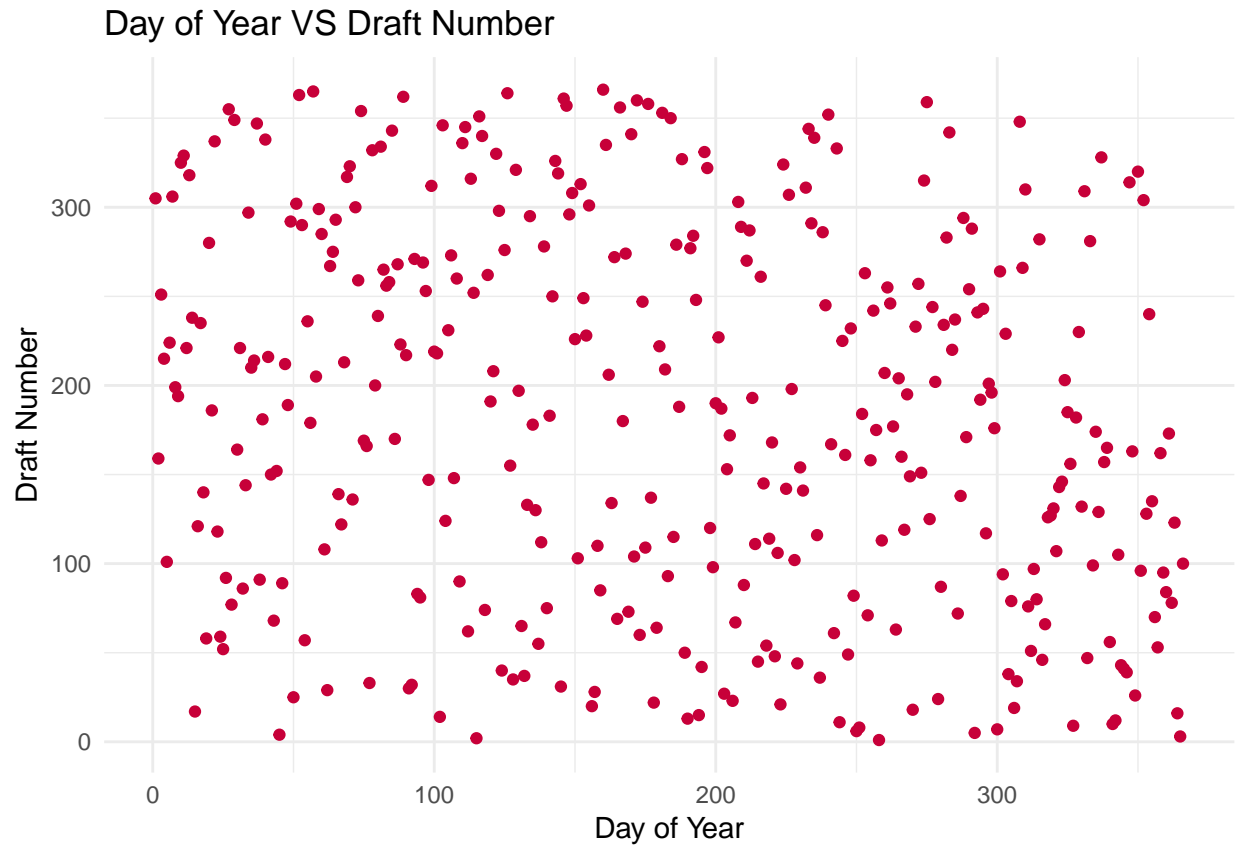
## 1 Question 1: Hypothesis testing

In 1970, the US Congress instituted a random selection process for the military draft. All 366 possible birth dates were placed in plastic capsules in a rotating drum and were selected one by one. The first date drawn from the drum received draft number one, the second date drawn received draft number two, etc. Then, eligible men were drafted in the order given by the draft number of their birth date. In a truly random lottery there should be no relationship between the date and the draft number. Your task is to investigate whether or not the draft numbers were randomly selected. The draft numbers ( $Y = DraftNo$ ) sorted by day of year ( $X = Day\ of\ year$ ) are given in the file `lottery.xls`.

Day	Month	Mo.Number	Day_of_year	Draft_No
1	Jan	1	1	305
2	Jan	1	2	159
3	Jan	1	3	251
4	Jan	1	4	215
5	Jan	1	5	101
6	Jan	1	6	224

### 1.1 Task 1.1

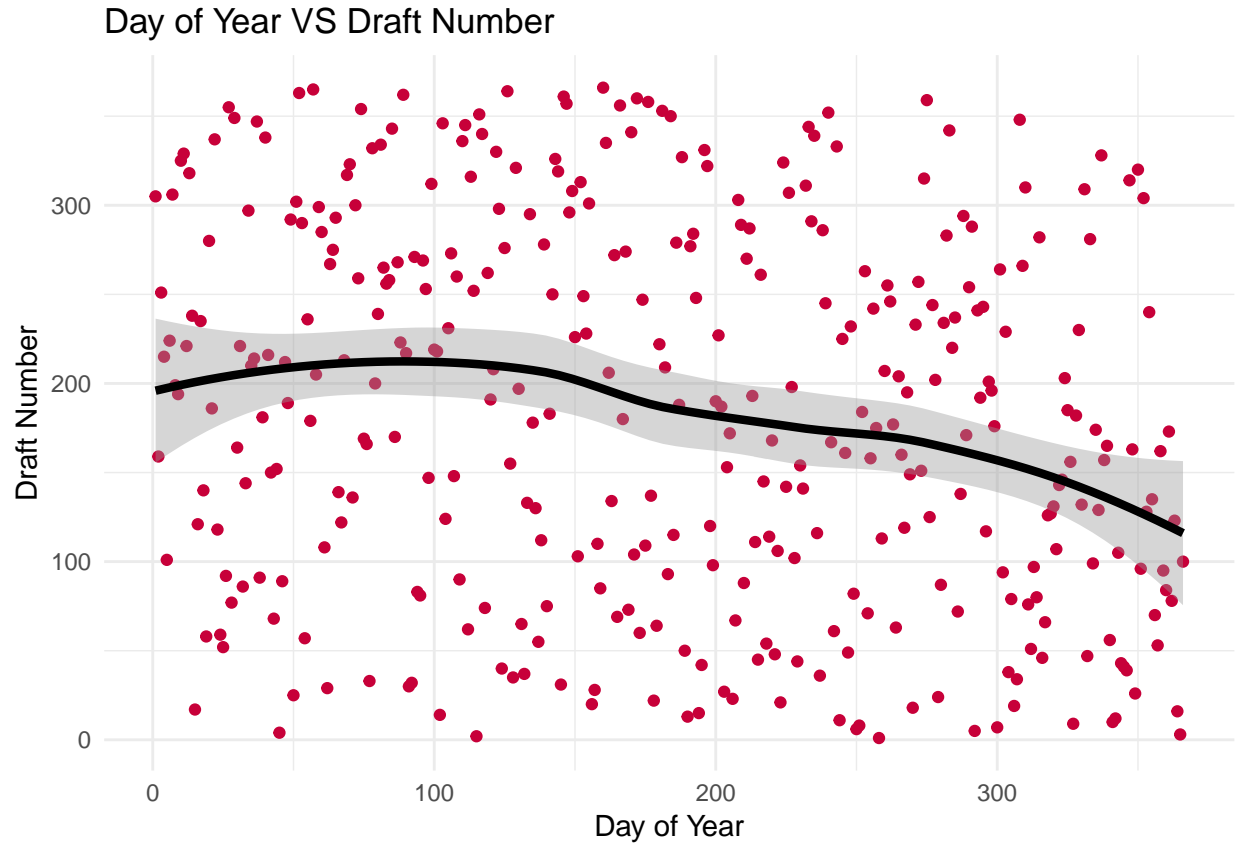
Make a scatterplot of Y versus X and conclude whether the lottery looks random.



**Answer:** In the plot, it the lottery appears to be random, because of the relatively even distribution of the points over the area.

## 1.2 Task 1.2

Compute an estimate  $\hat{Y}$  of the expected response as a function of  $X$  by using a loess smoother (use `loess()`), put the curve  $\hat{Y}$  versus  $X$  in the previous graph and state again whether the lottery looks random.



**Answer:** Including the insights from this smoothing we see, that the lottery actually doesn't look random. There are two reasons for this:

- The confidence interval of the curve is relatively small, especially compared to the spread of the data points.
- The data points follow the pattern of the curve slightly, especially above the line.

### 1.3 Task 1.3

To check whether the lottery is random, it is reasonable to use test statistics

$$T = \frac{\hat{Y}(X_b) - \hat{Y}(X_a)}{X_b - X_a}, \text{ where } X_b = \operatorname{argmax}_X Y(X), X_a = \operatorname{argmin}_X Y(X)$$

If this value is significantly greater than zero, then there should be a trend in the data and the lottery is not random. Estimate the distribution of  $T$  by using a non-parametric bootstrap with  $B = 2000$  and comment whether the lottery is random or not. What is the p-value of the test?

**Answer:**

```
data = data.frame(X = lottery$Day_of_year, Y = lottery$Draft_No)

test_statistics = function(X, Y, Y_hat) {

  b_index = which.max(Y)
  a_index = which.min(Y)
```

```

    return((Y_hat[b_index] - Y_hat[a_index]) / (X[b_index] - X[a_index]))
}

f = function(data, ind) {
  data1 = data[ind,]
  model = loess(Draft_No ~ Day_of_year, data1)

  T_value =
    test_statistics(data1$Day_of_year, data1$Draft_No, Y_hat = model$fitted)

  return(T_value)
}

# T(D) for the original data
data$Y_hat = loess(Draft_No ~ Day_of_year, lottery)$fitted
T_value_original = test_statistics(data$X, data$Y, data$Y_hat)

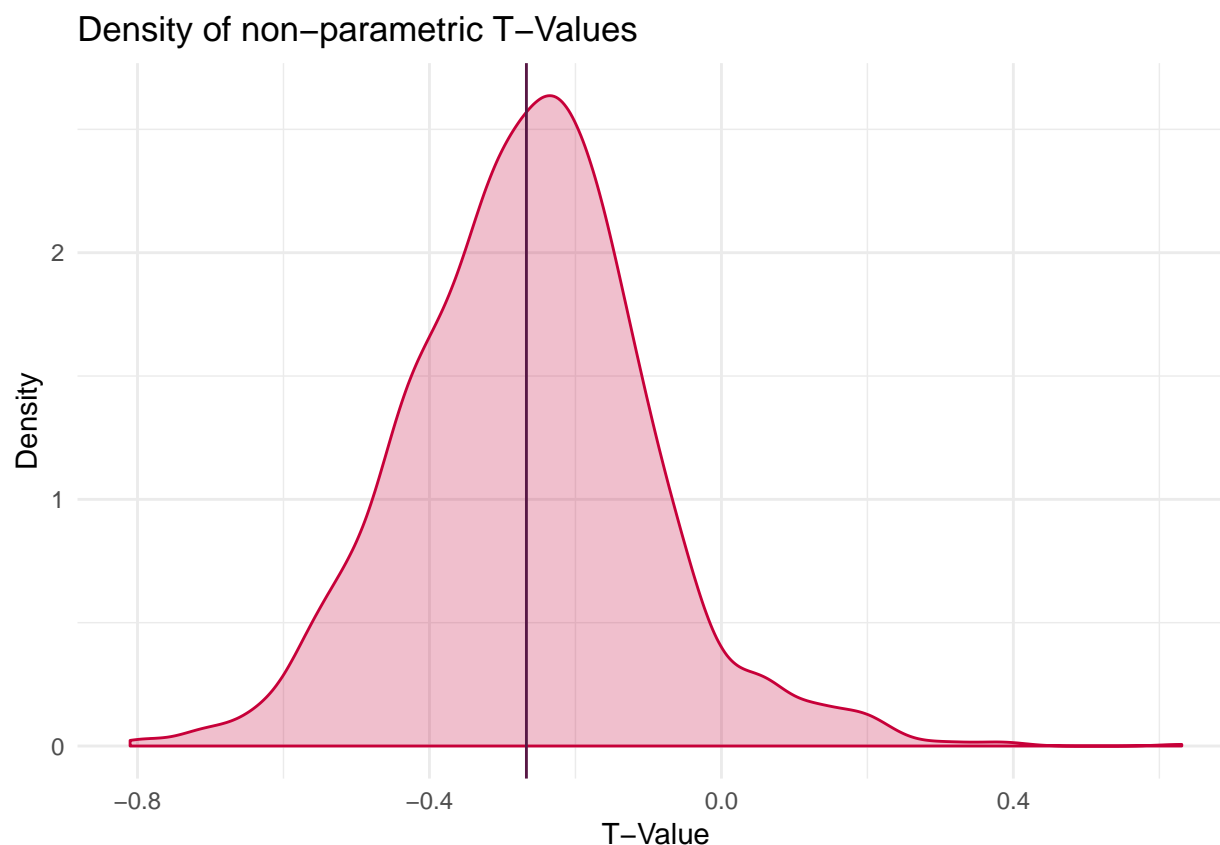
# T for the bootstrapped samples
nonparam_bootstrap =
  boot(lottery, statistic = f, R = 2000, parallel = "multicore")
p_value_original = mean(nonparam_bootstrap$t > T_value_original)

print(T_value_original)

## [1] -0.2671794
print(p_value_original)

## [1] 0.5155

```



## 1.4 Task 1.4

Implement a function depending on data and B that tests the hypothesis

$H_0$ : Lottery is random versus  $H_1$ : Lottery is non-random

by using a permutation test with statistics T. The function is to return the p-value of this test. Test this function on our data with  $B = 2000$ .

## 1.5 Task 2.5

Make a crude estimate of the power of the test constructed in Step 4:

- a) Generate (an obviously non-random) dataset with  $n = 366$  observations by using same X as in the original data set and  $Y(x) = \max(0, \min(\alpha x + \beta, 366))$ , where  $\alpha = 0.1$  and  $\beta \sim N(183, sd = 10)$ .
- b) Plug these data into the permutation test with  $B = 200$  and note whether it was rejected.
- c) Repeat Steps 5a-5b for  $\alpha = 0.2, 0.3, \dots, 10$ .

What can you say about the quality of your test statistics considering the value of the power?

## 2 Question 2: Bootstrap, jackknife and confidence intervals

The data you are going to continue analyzing is the database of home prices in Albuquerque, 1993. The variables present are **Price**; **SqFt**: the area of a house; **FEATS**: number of features such as dishwasher,

refrigerator and so on; **Taxes**: annual taxes paid for the house. Explore the file `prices1.xls`

```
prices = read_xls("prices1.xls")
```

## 2.1 Task 2.1

Plot the histogram of Price. Does it remind any conventional distribution? Compute the mean price.

## 2.2 Task 2.2

Estimate the distribution of the mean price of the house using bootstrap. Determine the bootstrap bias-correction and the variance of the mean price. Compute a 95% confidence interval for the mean price using bootstrap percentile, bootstrap BCa, and first-order normal approximation (Hint: use `boot()`, `boot.ci()`, `plot.boot()`, `print.bootci()`)

## 2.3 Task 2.3

Estimate the variance of the mean price using the jackknife and compare it with the bootstrap estimate

## 2.4 Task 2.4

Compare the confidence intervals obtained with respect to their length and the location of the estimated mean in these intervals.

# 3 Source Code

```
knitr::opts_chunk$set(echo = TRUE, cache = FALSE, include = TRUE, eval = TRUE)
library(knitr)
library(readxl)
library(ggplot2)
library(gridExtra)
library(boot)

set.seed(12345)

lottery = read_xls("lottery.xls")
kable(head(lottery))

ggplot(lottery)+
  geom_point(aes(x = Day_of_year, y = Draft_No), color = "#C70039") +
  labs(title = "Day of Year VS Draft Number",
       y = "Draft Number", x = "Day of Year", color = "Legend") +
  theme_minimal()

ggplot(lottery)+
  geom_point(aes(x = Day_of_year, y = Draft_No), color = "#C70039") +
```

```

geom_smooth(mapping = aes(x = Day_of_year, y = Draft_No),
             method = "loess", size = 1.5, color = "#000000") +
labs(title = "Day of Year VS Draft Number",
     y = "Draft Number", x = "Day of Year", color = "Legend") +
theme_minimal()

data = data.frame(X = lottery$Day_of_year, Y = lottery$Draft_No)

test_statistics = function(X, Y, Y_hat) {

  b_index = which.max(Y)
  a_index = which.min(Y)

  return((Y_hat[b_index] - Y_hat[a_index]) / (X[b_index] - X[a_index]))
}

f = function(data, ind) {
  data1 = data[ind,]
  model = loess(Draft_No ~ Day_of_year, data1)

  T_value =
    test_statistics(data1$Day_of_year, data1$Draft_No, Y_hat = model$fitted)

  return(T_value)
}

# T(D) for the original data
data$Y_hat = loess(Draft_No ~ Day_of_year, lottery)$fitted
T_value_original = test_statistics(data$X, data$Y, data$Y_hat)

# T for the bootstrapped samples
nonparam_bootstrap =
  boot(lottery, statistic = f, R = 2000, parallel = "multicore")
p_value_original = mean(nonparam_bootstrap$t > T_value_original)

print(T_value_original)
print(p_value_original)

df = data.frame(nonparam_bootstrap$t)

ggplot(df) +
  geom_density(aes(x = nonparam_bootstrap.t, color = "#c70039",
                  fill = "#c70039", alpha = 0.25)) +
  geom_vline(aes(xintercept = T_value_original), color = "#581845") +
  labs(title = "Density of non-parametric T-Values",
       y = "Density", x = "T-Value", color = "Legend") +
  theme_minimal()

test_hypothesis = function (data, statistics, B) {

```

```
}
```

```
prices = read_xls("prices1.xls")
```