

# Computational Statistics (732A90) Lab5

*Anubhav Dikshit(anudi287) and Thijs Quast(thiqu264)*

*12 Feburary 2019*

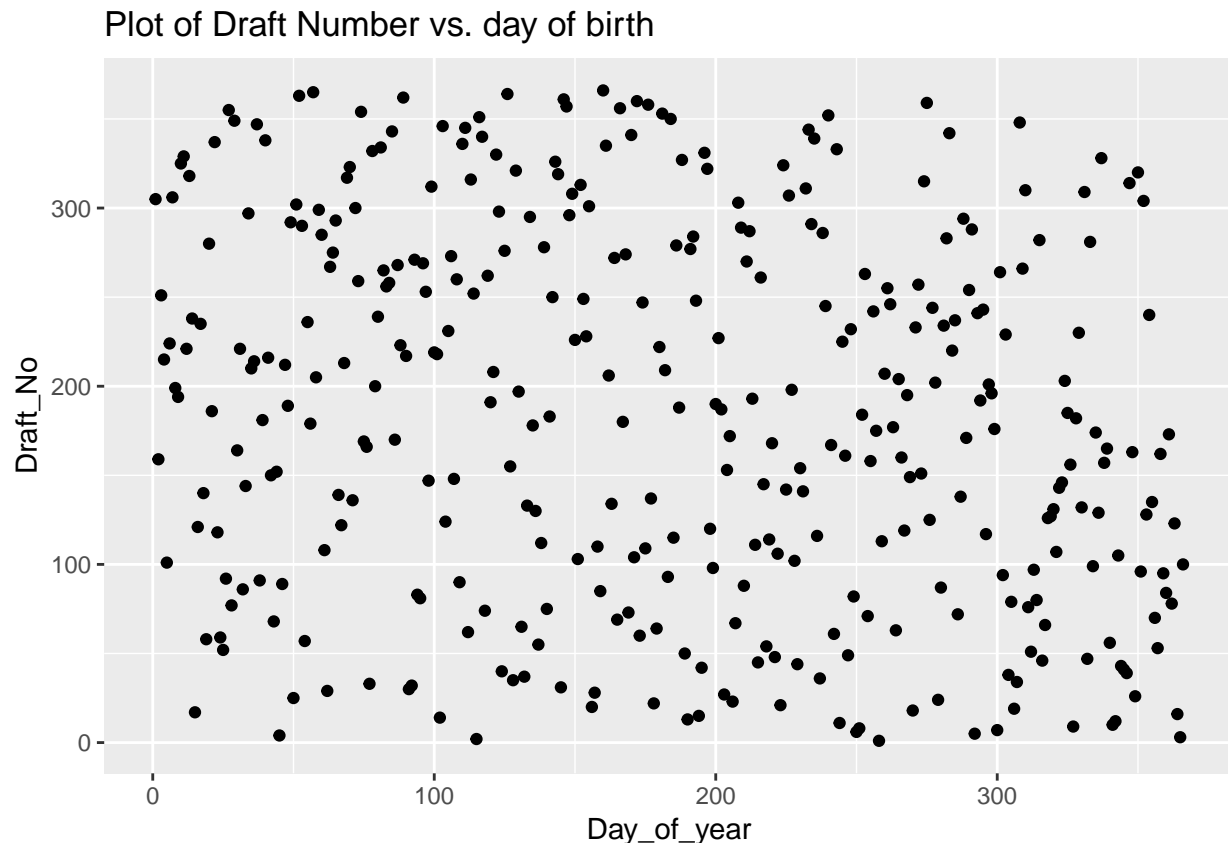
## Contents

<b>Question 1: Hypothesis testing</b>	<b>2</b>
1. Make a scatterplot of Y(draft_no) versus X(day_of_year) and conclude whether the lottery looks random. . . . .	2
2. Compute an estimate Y(hat) of the expected response as a function of X by using a loess smoother (use loess()), put the curve Y(hat) versus X in the previous graph and state again whether the lottery looks random. . . . .	2
3. To check whether the lottery is random, it is reasonable to use test statistics . . . . .	4
4. Implement a function depending on data and B that tests the hypothesis H0: Lottery is random versus H1: 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. . . . .	5
5 Make a crude estimate of the power of the test constructed in Step4 : . . . . .	7
<b>Question 2: Bootstrap, jackknife and confidence intervals</b>	<b>9</b>
1. Plot the histogram of Price. Does it remind any conventional distribution? Compute the mean price. . . . .	9
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	10
3. Estimate the variance of the mean price using the jackknife and compare it with the bootstrap estimate . . . . .	13
4. Compare the confidence intervals obtained with respect to their length and the location of the estimated mean in these intervals. . . . .	14
<b>Appendix</b>	<b>15</b>

## Question 1: Hypothesis testing

1. Make a scatterplot of  $Y(\text{draft\_no})$  versus  $X(\text{day\_of\_year})$  and conclude whether the lottery looks random.

```
lottery <- read.csv("lottery.csv", sep=";")  
  
ggplot(lottery, aes(x=Day_of_year, y = Draft_No)) + geom_point() +  
  ggtitle("Plot of Draft Number vs. day of birth")
```

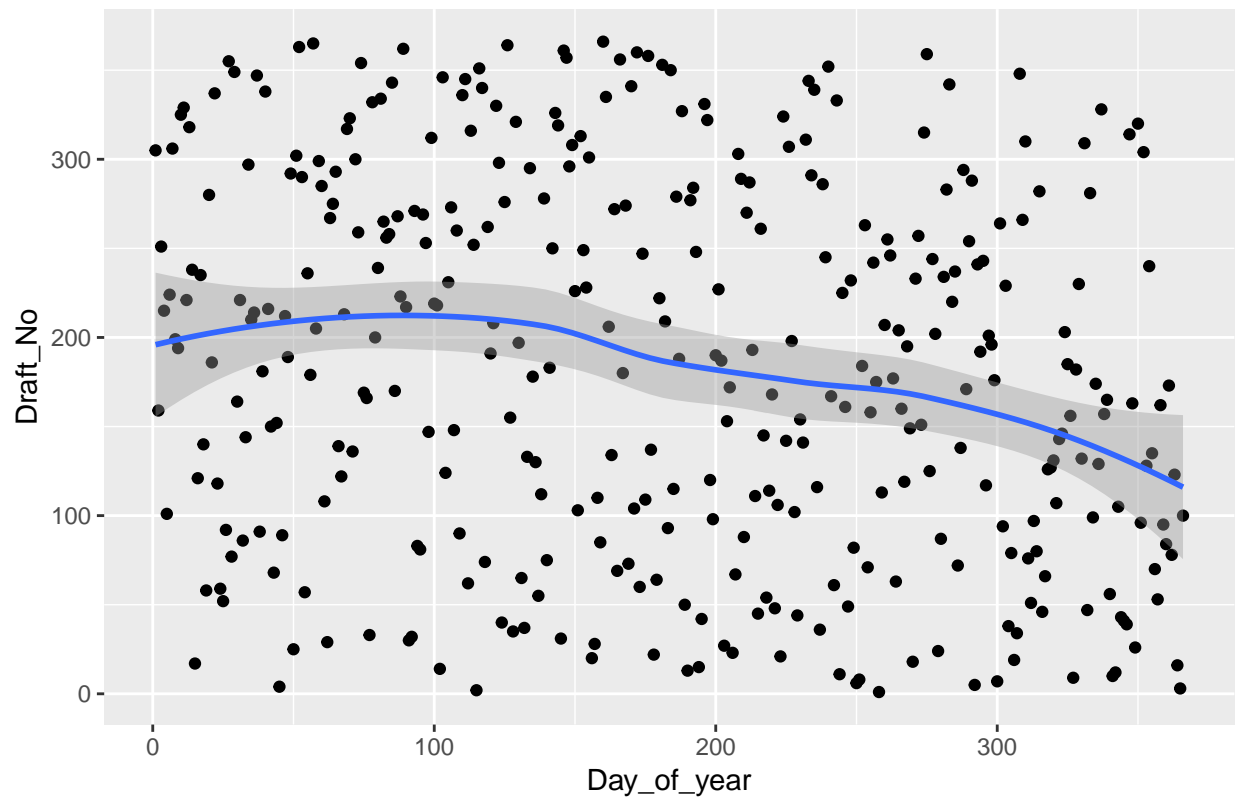


Analysis: The plot seems random and its very difficult to judge any sort of trend

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.

```
ggplot(lottery, aes(x=Day_of_year, y = Draft_No)) +  
  geom_point() +  
  geom_smooth(method = loess) +  
  ggtitle("Plot of Draft Number vs. Day of birth")
```

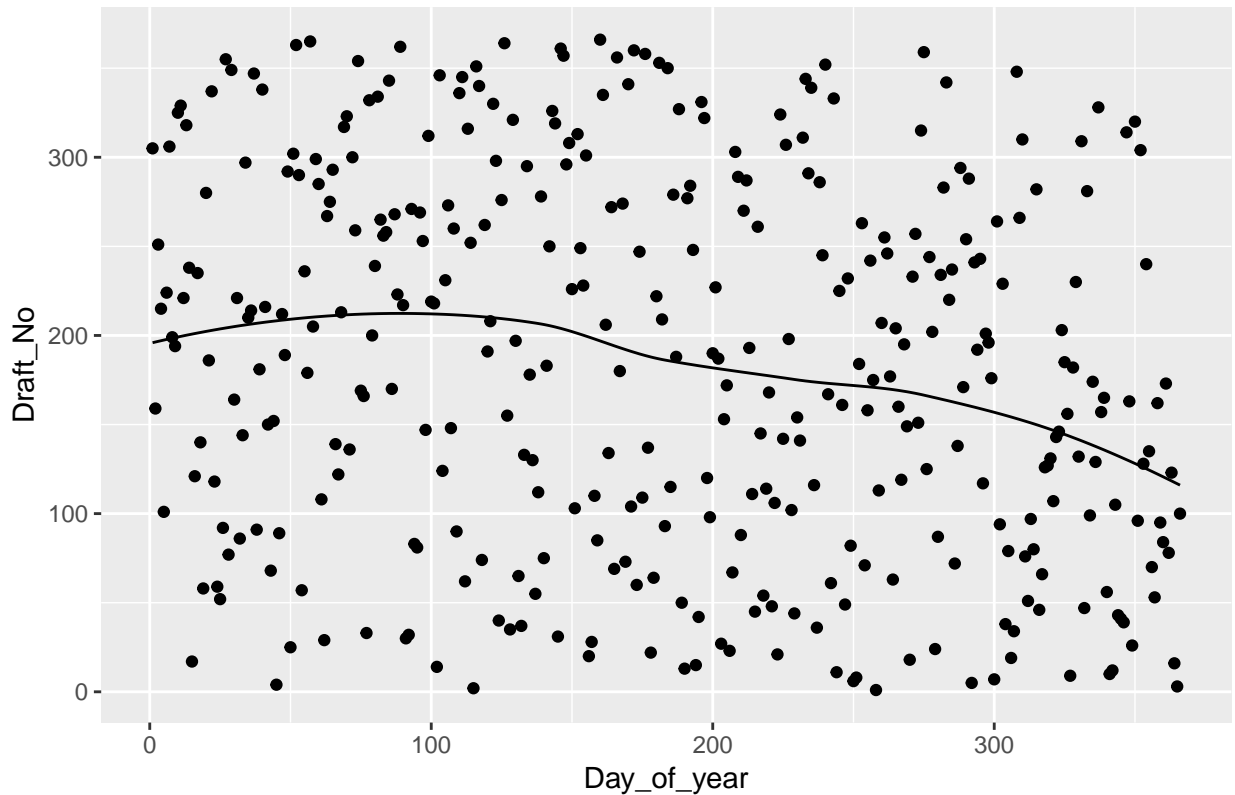
Plot of Draft Number vs. Day of birth



```
model <- loess(Draft_No ~ Day_of_year, lottery)
lottery$Y_hat <- predict(model, lottery)

ggplot(lottery, aes(x=Day_of_year, y = Draft_No)) +
  geom_point() +
  geom_line(aes(y = Y_hat)) +
  ggtitle("Plot of Draft Number vs. Day of birth without using ggplot loess")
```

Plot of Draft Number vs. Day of birth without using ggplot loess



Analysis: One can see that the overall trend is downward, the implies the more people were born in the first half of the year than the later half.

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

$$T = \frac{\hat{Y}(X_b) - \hat{Y}(X_a)}{X_b - X_a}$$

Where  $X_b = \operatorname{argmax}_x Y(X)$  and  $X_a = \operatorname{argmin}_x Y(X)$ .

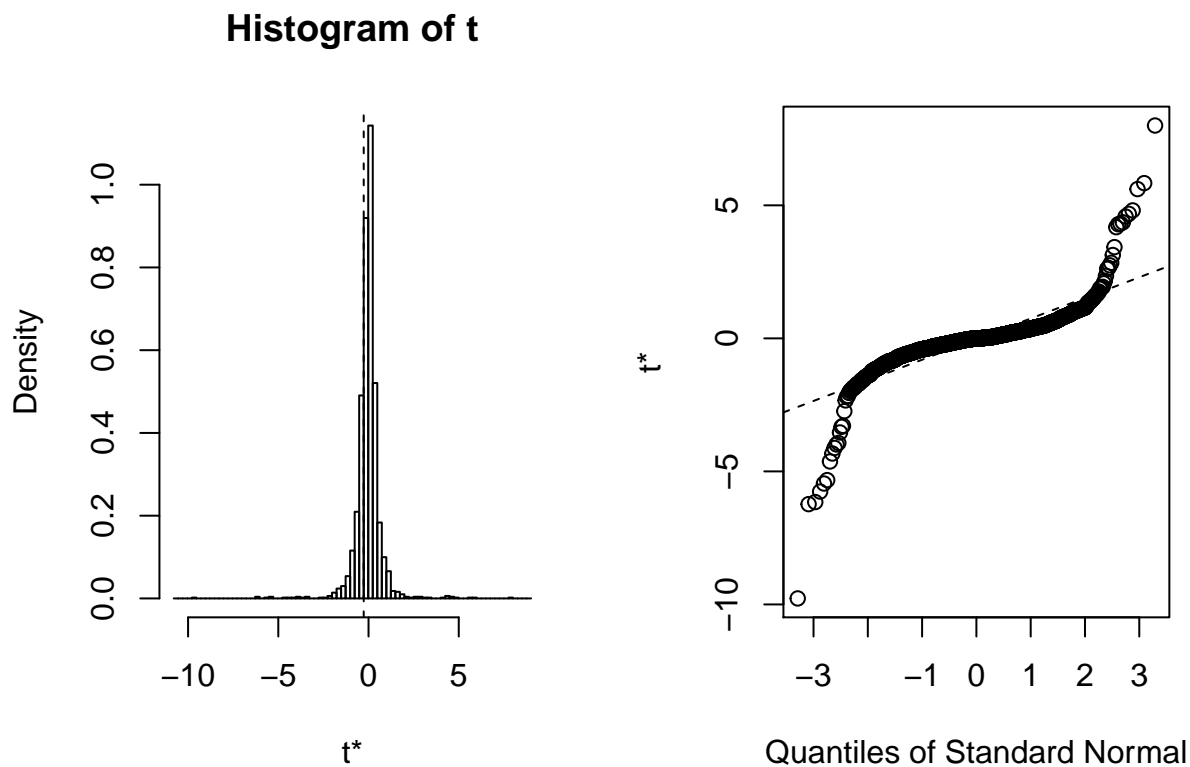
If this value is significantly greater than zero, then there should be a trend in the data and the lottery is not random. Estimate the distribution of T by using a non-parametric bootstrap with  $B = 2000$  and comment whether the lottery is random or not. What is the p-value of the test?

```
stat1 <- function(data, index){
  data <- data[index,]
  model <- loess(Draft_No ~ Day_of_year, data)
  res <- predict(model, data)
  X_a <- data$Day_of_year[which.min(data$Draft_No)]
  X_b <- data$Day_of_year[which.max(data$Draft_No)]
  Y_a <- res[X_a]
  Y_b <- res[X_b]
  answer <- ((Y_b - Y_a) / (X_b - X_a))
  return(answer)
}
```

```
res <- boot(data=lottery, statistic = stat1, R=2000)
pval <- length(which(res$t>=0))/2000
cat("The estimated p-value is",pval)
```

```
## The estimated p-value is 0.4865
```

```
plot(res)
```



Analysis: From the p-value obtained and plot of the histogram we can conclude that the lottery number is not random. The histogram is has its mean towards the right side and is skewed towards left. The p-value is 0.4865.

**4.Implement a function depending on data and B that tests the hypothesis H0: Lottery is random versus H1: 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.**

```
permu_function <- function(data, index){

  temp <- function(data, index){
    data <- data[index,]
    Y <- as.data.frame(data$Draft_No)
    X <- as.data.frame(data$Day_of_year)
    X <- X[sample(nrow(X), replace = FALSE),]
    data <- as.data.frame(cbind(X,Y))
```

```

colnames(data) <- c("Day_of_year", "Draft_No")
model <- loess(Draft_No ~ Day_of_year, data)
res <- predict(model, data)
X_a <- data$Day_of_year[which.min(data$Draft_No)]
X_b <- data$Day_of_year[which.max(data$Draft_No)]
Y_a <- res[X_a]
Y_b <- res[X_b]
answer <- ((Y_b - Y_a) / (X_b - X_a))
return(answer)
}

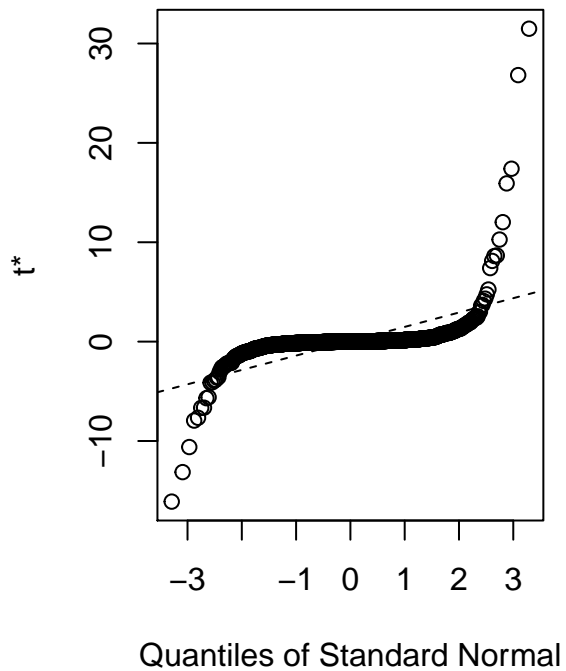
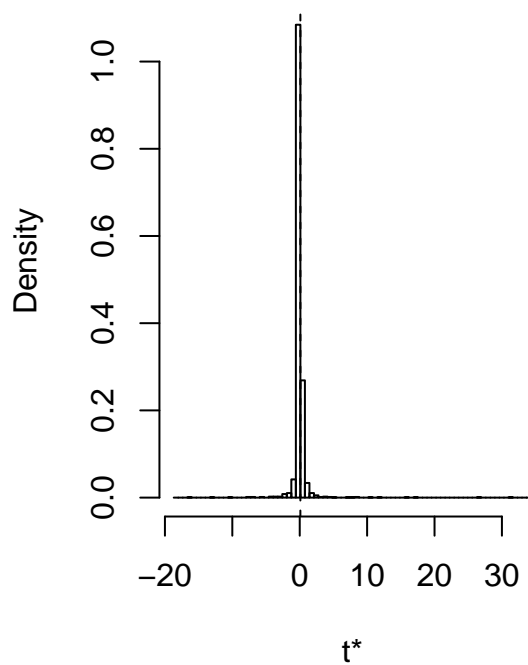
res2 <- boot(data=lottery, statistic = temp, R=2000)

plot(res2)
pval2 <- length(which(as.vector(abs(res2$t) > abs(res2$t0))))/2000
return(pval2)
}

pval <- permu_function(data=lottery, index=2000)

```

### Histogram of t



```
cat("The estimated p-value is",pval)
```

```
## The estimated p-value is 0.4335
```

Analysis: The plot suggests that lottery number is not random and the p-value is 0.4335.

## 5 Make a crude estimate of the power of the test constructed in Step4 :

(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  $\beta \sim N(183, sd = 10)$ .

```
permu_test <- function(B,X,Y){

  #compute test statistics from original data
  data <- cbind(X,Y)
  X_a <- data[which.min(data[,2])]
  X_b <- data[which.max(data[,2])]

  model <- loess(Y~X, data=as.data.frame(data), method="loess")

  fitted_X_a <- model$fitted[X_a]
  fitted_X_b <- model$fitted[X_b]

  test <- (fitted_X_b - fitted_X_a)/(X_b - X_a)

  #then compute the permuted data
  permu_T <- numeric()

  for (i in 1:B){
    permu_Y <- sample(1:length(Y), length(Y), replace=FALSE)
    permu_data <- cbind(X,permu_Y)

    #and do the same thing
    X_a <- permu_data[which.min(permu_data[,2])]
    X_b <- permu_data[which.max(permu_data[,2])]

    permu_model <- loess(permu_Y~X, data=as.data.frame(permu_data), method="loess")

    fitted_X_a <- permu_model$fitted[X_a]
    fitted_X_b <- permu_model$fitted[X_b]

    permu_test <- (fitted_X_b - fitted_X_a)/(X_b - X_a)
    permu_T[i] <- permu_test
  }

  #then compute the estimated p value
  estimated_pval <- length(which(abs(permu_T)>=abs(test)))/B
  return(estimated_pval)
}

new_Y <- function(alpha){

  X <- lottery$Day_of_year
  new_y <- numeric()

  for (i in 1:length(X)){
```

```

        beta <- rnorm(1,183,10)
        new_y[i] <- max(0,min(alpha*X[i]+beta,366))
    }

    return(new_y)
}

non_ran_Y <- new_Y(alpha=0.1)
head(non_ran_Y)

## [1] 205.1485 178.7821 181.2708 187.3799 189.2456 179.4140

```

(b) Plug this data into the permutation test with  $B=200$  and note whether it was rejected.

```
permu_test(B=200, X=lottery$Day_of_year, Y=non_ran_Y)
```

```
## [1] 0.45
```

Analysis: The obtained p-value is 0.45, so we can't reject null hypothesis.

(c) Repeat Steps 5a-5b for  $\alpha = 0.2, 0.3, \dots, 10$ .

```

seq <- seq(0.2,10,by=0.1)
pvals <- numeric()

for (i in 1:length(seq)){
    alpha <- seq[i]
    newY <- new_Y(alpha)

    pvals[i] <- permu_test(B=200, X=lottery$Day_of_year, Y=newY)
}

signif_p <- which(pvals<0.05)
power <- 1-sum(pvals>0.05)/length(pvals)

list("total_number_of_hypothesis_test"=length(seq),
     "number_of_test_which_reject_H0"=length(signif_p),
     "power"=power)

## $total_number_of_hypothesis_test
## [1] 99
##
## $number_of_test_which_reject_H0
## [1] 97
##
## $power
## [1] 0.979798

```

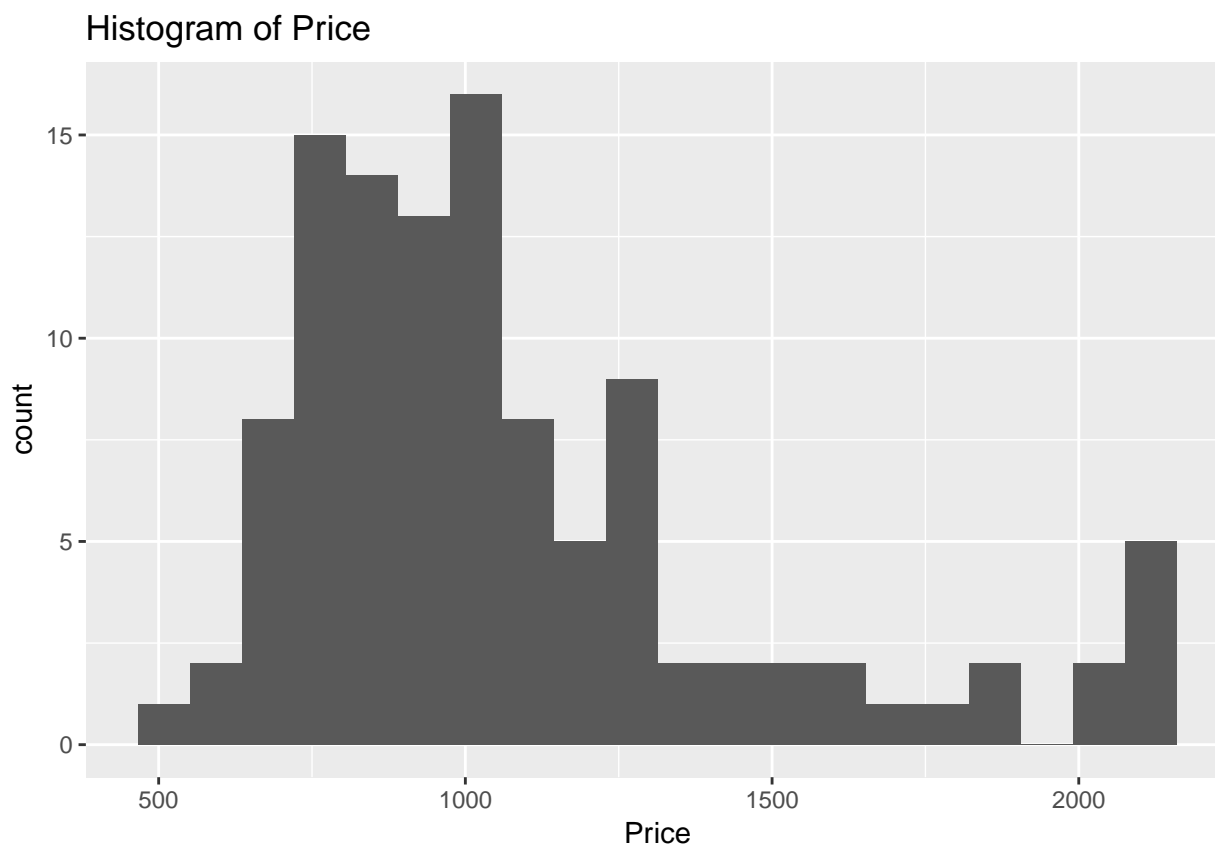
Analysis: From the result above, we can see that 98% of the cases reject the null hypothesis, which indicates that the lottery is not random. Since we have generated obviously non-random dataset, we can say that the test statistics performs well enough. The obtained power also supports that the quality of statistics is good enough.



## Question 2: Bootstrap, jackknife and confidence intervals

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

```
price_data <- read.csv("prices1.csv", sep=";")  
  
ggplot(data=price_data,aes(Price)) +  
  geom_histogram(bins=20) +  
  ggtitle("Histogram of Price")
```



```
cat("The mean price is",mean(price_data$Price))
```

```
## The mean price is 1080.473
```

Analysis: The distribution reminds us of the 'Beta distribution'

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

Bias correction

$$T1 = 2.T(D) - \frac{1}{D} \sum_{i=1}^B T_i^*$$

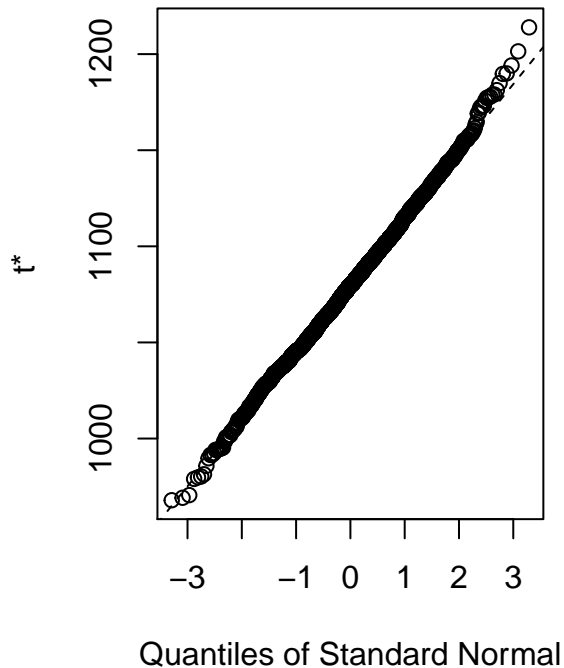
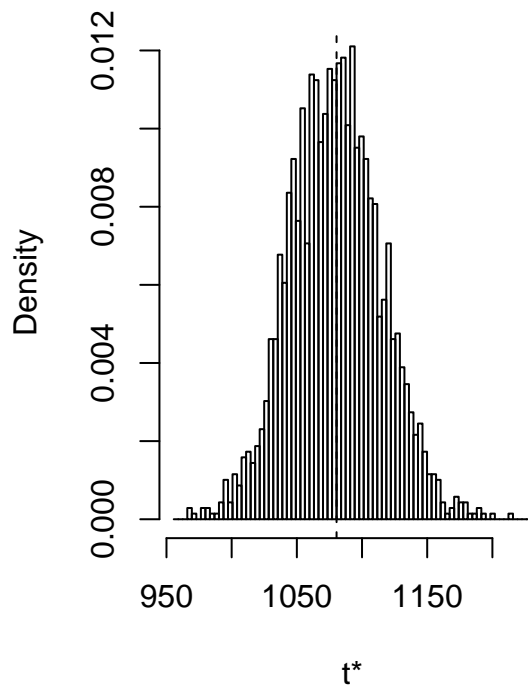
```
# Estimation of mean of Price
stat_mean <- function(data, index){
  data <- data[index,]
  answer <- mean(data$Price)
  return(answer)
}

res <- boot::boot(data=price_data, statistic = stat_mean, R=2000)
res

##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot::boot(data = price_data, statistic = stat_mean, R = 2000)
##
##
## Bootstrap Statistics :
##      original      bias    std. error
## t1* 1080.473 -0.6945636    34.86858

plot(res,index = 1)
```

### Histogram of t



```
#95% CI for mean using percentile
boot.ci(res, index=1, type=c('perc'))
```

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 2000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = res, type = c("perc"), index = 1)
##
## Intervals :
## Level      Percentile
## 95%      (1012, 1149 )
## Calculations and Intervals on Original Scale
```

```
#95% CI for mean using bca
boot.ci(res, index=1, type=c('bca'))
```

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 2000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = res, type = c("bca"), index = 1)
##
## Intervals :
## Level      BCa
## 95%      (1017, 1155 )
## Calculations and Intervals on Original Scale
```

```

#95% CI for mean using first order normal
boot.ci(res, index=1, type=c('norm'))

## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 2000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = res, type = c("norm"), index = 1)
##
## Intervals :
## Level      Normal
## 95%      (1013, 1150 )
## Calculations and Intervals on Original Scale

# Bias-correction and Variance of Price
boot.fn <- function(data,index){
  d <- data[index]
  res <- mean(d)
}

boot.result <- boot(data=price_data$Price, statistic=boot.fn, R=1000)
bias_cor <- 2*mean(price_data$Price)-mean(boot.result$t)

list("bias_correction"=bias_cor, "variance_of_the_mean_price"=35.93^2)

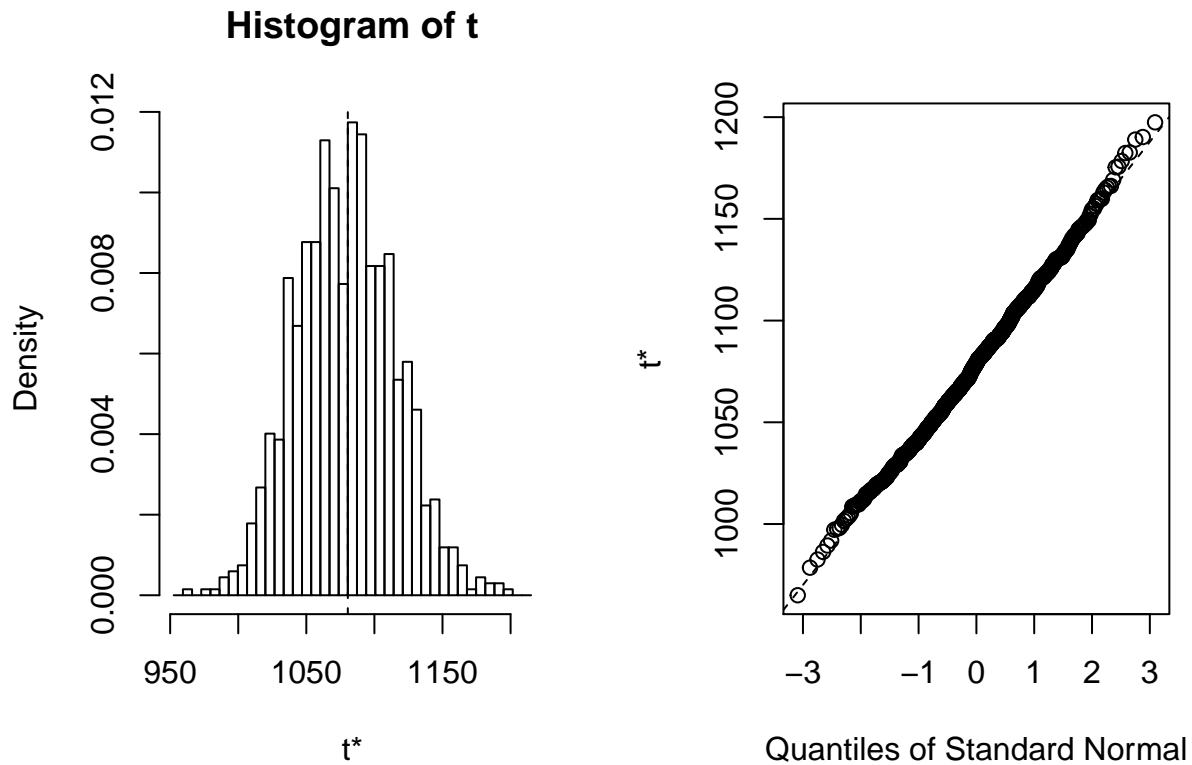
## $bias_correction
## [1] 1081.983
##
## $variance_of_the_mean_price
## [1] 1290.965

boot.result

##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = price_data$Price, statistic = boot.fn, R = 1000)
##
##
## Bootstrap Statistics :
##      original      bias    std. error
## t1* 1080.473 -1.510273    36.30546

plot(boot.result)

```



```
boot.ci(boot.result)
```

```
## Warning in boot.ci(boot.result): bootstrap variances needed for studentized
## intervals
```

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
```

```
## Based on 1000 bootstrap replicates
```

```
##
```

```
## CALL :
```

```
## boot.ci(boot.out = boot.result)
```

```
##
```

```
## Intervals :
```

```
## Level      Normal      Basic
```

```
## 95%   (1011, 1153 )   (1009, 1149 )
```

```
##
```

```
## Level      Percentile      BCa
```

```
## 95%   (1012, 1152 )   (1016, 1159 )
```

```
## Calculations and Intervals on Original Scale
```

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

Jackknife(n=B):

$$\widehat{Var}[T(\cdot)] = \frac{1}{n(n-1)} \sum_{i=1}^n (T_i^* - J(T))^2$$

where,

$$T_i^* = nT(D) - (n-1)T(D_i^*), \quad J(T) = \frac{1}{n} \sum_{i=1}^n T_i^*$$

When you compute the equation given above, you got

$$\frac{n-1}{n} \sum_{i=1}^n (T_i^* - J(T))^2$$

Reference: The Jackknife Estimation Method, Avery I. McIntosh (<http://people.bu.edu/aimcinto/jackknife.pdf>)

```
result <- numeric()
n <- NROW(price_data)

for (i in 1:n){
  updated_price <- price_data$Price[-i]
  result[i] <- mean(updated_price)
}

var_T <- (n-1)/n*sum((result-mean(result))^2)
mean_T <- mean(result)

cat("The variance from jackknife method is:", var_T)
```

```
## The variance from jackknife method is: 1320.911
```

Analysis: The obtained variance using Jackknife method is 1320.911 while using bootstrapping the obtained value was 1290.965. Considering the fact that Jackknife overestimate variance, the answer seems reasonable.

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

```
confidence_interval_jackknife <- c((mean_T - 1.96*var_T), (mean_T + 1.96*var_T))

confidence_interval_jackknife

## [1] -1508.513 3669.458

intervals <- c("(1150-1011)", "(1148-1011)", "(1149-1013)", "(1146-1007)")
length <- c(1150-1011, 1148-1011, 1149-1013, 1146-1007)
Center_of_interval <- c((1150+1011)/2, (1148+1011)/2, (1149+1013)/2, (1146+1007)/2)

dt <- data.frame(intervals, length, Center_of_interval, row.names = c("Normal", "Basic", "Percentile", "BCa"))

dt %>% kable(col.names = c("Confidence interval", "Length of interval", "Center of interval"))
```

	Confidence interval	Length of interval	Center of interval
Normal	(1150-1011)	139	1080.5
Basic	(1148-1011)	137	1079.5
Percentile	(1149-1013)	136	1081.0
BCa	(1146-1007)	139	1076.5

## Appendix

```
knitr::opts_chunk$set(echo = TRUE)
options(scipen=999)
library(dplyr)
library(ggplot2)
library(knitr)
library("boot")
set.seed(12345)
lottery <- read.csv("lottery.csv", sep=";")

ggplot(lottery, aes(x=Day_of_year, y = Draft_No)) + geom_point() +
  ggtitle("Plot of Draft Number vs. day of birth")
ggplot(lottery, aes(x=Day_of_year, y = Draft_No)) +
  geom_point() +
  geom_smooth(method = loess) +
  ggtitle("Plot of Draft Number vs. Day of birth")

model <- loess(Draft_No ~ Day_of_year, lottery)
lottery$Y_hat <- predict(model, lottery)

ggplot(lottery, aes(x=Day_of_year, y = Draft_No)) +
  geom_point() +
  geom_line(aes(y = Y_hat)) +
  ggtitle("Plot of Draft Number vs. Day of birth without using ggplot loess")

stat1 <- function(data, index){
  data <- data[index,]
  model <- loess(Draft_No ~ Day_of_year, data)
  res <- predict(model, data)
  X_a <- data$Day_of_year[which.min(data$Draft_No)]
  X_b <- data$Day_of_year[which.max(data$Draft_No)]
  Y_a <- res[X_a]
  Y_b <- res[X_b]
  answer <- ((Y_b - Y_a) / (X_b - X_a))
  return(answer)
}

res <- boot(data=lottery, statistic = stat1, R=2000)
pval <- length(which(res$t>=0))/2000
cat("The estimated p-value is", pval)
plot(res)

permu_function <- function(data, index){
  temp <- function(data, index){
    data <- data[index,]
    Y <- as.data.frame(data$Draft_No)
    X <- as.data.frame(data$Day_of_year)
```

```

X <- X[sample(nrow(X), replace = FALSE),]
data <- as.data.frame(cbind(X,Y))
colnames(data) <- c("Day_of_year", "Draft_No")
model <- loess(Draft_No ~ Day_of_year, data)
res <- predict(model, data)
X_a <- data$Day_of_year[which.min(data$Draft_No)]
X_b <- data$Day_of_year[which.max(data$Draft_No)]
Y_a <- res[X_a]
Y_b <- res[X_b]
answer <- ((Y_b - Y_a) / (X_b - X_a))
return(answer)
}

res2 <- boot(data=lottery, statistic = temp, R=2000)

plot(res2)
pval2 <- length(which(as.vector(abs(res2$t) > abs(res2$t0))))/2000
return(pval2)
}

pval <- permu_function(data=lottery, index=2000)

cat("The estimated p-value is",pval)

permu_test <- function(B,X,Y){

  #compute test statistics from original data
  data <- cbind(X,Y)
  X_a <- data[which.min(data[,2])]
  X_b <- data[which.max(data[,2])]

  model <- loess(Y~X, data=as.data.frame(data), method="loess")

  fitted_X_a <- model$fitted[X_a]
  fitted_X_b <- model$fitted[X_b]

  test <- (fitted_X_b - fitted_X_a)/(X_b - X_a)

  #then compute the permuted data
  permu_T <- numeric()

  for (i in 1:B){
    permu_Y <- sample(1:length(Y), length(Y), replace=FALSE)
    permu_data <- cbind(X,permu_Y)

    #and do the same thing
    X_a <- permu_data[which.min(permu_data[,2])]
    X_b <- permu_data[which.max(permu_data[,2])]

    permu_model <- loess(permu_Y~X, data=as.data.frame(permu_data), method="loess")

    fitted_X_a <- permu_model$fitted[X_a]
    fitted_X_b <- permu_model$fitted[X_b]
  }
}

```



```

        permu_test <- (fitted_X_b - fitted_X_a)/(X_b - X_a)
        permu_T[i] <- permu_test
    }

    #then compute the estimated p value
    estimated_pval <- length(which(abs(permu_T)>=abs(test)))/B
    return(estimated_pval)
}

new_Y <- function(alpha){

    X <- lottery$Day_of_year
    new_y <- numeric()

    for (i in 1:length(X)){
        beta <- rnorm(1,183,10)
        new_y[i] <- max(0,min(alpha*X[i]+beta,366))
    }

    return(new_y)
}

non_ran_Y <- new_Y(alpha=0.1)
head(non_ran_Y)
permu_test(B=200, X=lottery$Day_of_year, Y=non_ran_Y)

seq <- seq(0.2,10,by=0.1)
pvals <- numeric()

for (i in 1:length(seq)){
    alpha <- seq[i]
    newY <- new_Y(alpha)

    pvals[i] <- permu_test(B=200, X=lottery$Day_of_year, Y=newY)
}

signif_p <- which(pvals<0.05)
power <- 1-sum(pvals>0.05)/length(pvals)

list("total_number_of_hypothesis_test"=length(seq),
     "number_of_test_which_reject_H0"=length(signif_p),
     "power"=power)

price_data <- read.csv("prices1.csv", sep=";")

ggplot(data=price_data,aes(Price)) +
  geom_histogram(bins=20) +
  ggtitle("Histogram of Price")

cat("The mean price is",mean(price_data$Price))

```

```

# Estimation of mean of Price
stat_mean <- function(data, index){
  data <- data[index,]
  answer <- mean(data$Price)
  return(answer)
}

res <- boot::boot(data=price_data, statistic = stat_mean, R=2000)
res
plot(res,index = 1)

#95% CI for mean using percentile
boot.ci(res, index=1, type=c('perc'))

#95% CI for mean using bca
boot.ci(res, index=1, type=c('bca'))

#95% CI for mean using first order normal
boot.ci(res, index=1, type=c('norm'))
# Bias-correction and Variance of Price
boot.fn <- function(data,index){
  d <- data[index]
  res <- mean(d)
}

boot.result <- boot(data=price_data$Price, statistic=boot.fn, R=1000)
bias_cor <- 2*mean(price_data$Price)-mean(boot.result$t)

list("bias_correction"=bias_cor, "variance_of_the_mean_price"=35.93^2)
boot.result
plot(boot.result)
boot.ci(boot.result)

result <- numeric()
n <- NROW(price_data)

for (i in 1:n){
  updated_price <- price_data$Price[-i]
  result[i] <- mean(updated_price)
}

var_T <- (n-1)/n*sum((result-mean(result))^2)
mean_T <- mean(result)

cat("The variance from jackknife method is:", var_T)

confidence_interval_jackknife <- c((mean_T - 1.96*var_T), (mean_T + 1.96*var_T))

confidence_interval_jackknife

intervals <- c("(1150-1011)", "(1148-1011)", "(1149-1013)", "(1146-1007)")

```

```

length <- c(1150-1011, 1148-1011,1149-1013,1146-1007)
Center_of_interval <- c((1150+1011)/2,(1148+1011)/2,(1149+1013)/2,(1146+1007)/2)

dt <- data.frame(intervals, length, Center_of_interval,row.names = c("Normal", "Basic", "Percentile", "I

dt %>% kable(col.names = c("Confidence interval", "Length of interval","Center of interval"))

```