# Statistics for Computer Science

## Assignment 2

## Kanitha Chim

## 501453

Field of Study Software System and Service Management

Faculty of Informatics
Masaryk University

May 24, 2020

# Exercise 3

1. Write down the formula for likelihood function of Poisson distribution.

The formula for the Poisson probability mass function is:

$$P(x, \lambda) = \frac{e^{-\lambda}\lambda^x}{x!}$$

In Poisson distribution, the parameter of interest is $\lambda$. Having sequence of $X_n$, the probability of observing the sequence $X_n$ will be the product of probabilities of each of them.

**Therefore, the kernel of likelihood function of Poisson distribution is:**

$$L(\lambda|X) = \prod_{i=1}^{N} \frac{\lambda^{X_i} e^{-\lambda}}{X_i!}$$

2. Write down the formula for log-likelihood function of Poisson distribution.

The formula for log-likelihood function of Poisson distribution is obtained by using natural logarithm on the likelihood function of Poisson distribution.

**Therefore, the kernel of log-likelihood function of Poisson distribution is:**

$$l(\lambda|X) = ln\left(\prod_{i=1}^{N} \frac{\lambda^{X_i} e^{-\lambda}}{X_i!}\right)$$

$$l(\lambda|X) = \sum_{i=1}^{N} X_i ln\lambda - N\lambda$$

3. Write down the likelihood equation and work out the exact formula for $\hat{\lambda}$.

$$L(\lambda) = \prod_{i=1}^{N} \frac{e^{-\lambda}\lambda^{x_i}}{x_i!} = e^{-N\lambda} \frac{\lambda^{\sum_1^N x_i}}{\prod_{i=1}^{N} x_i}$$

$$lnL(\lambda) = -N\lambda + \sum_1^N x_i ln(\lambda) - ln\left(\prod_{i=1}^{N} x_i\right)$$

$$\frac{dlnL(\lambda)}{dp} = -N + \sum_1^N x_i \frac{1}{\lambda}$$

$$\hat{\lambda} = \frac{\sum_{i=1}^{N} x_i}{N}$$

4. Create your own R-function for calculating the value of log-likelihood function of Poisson distribution.

```
1   #set up x according to exercise description
2   x <-c(117, 109, 109, 89, 120, 88, 99, 103, 109, 91, 107, 101, 109, 117,
        96, 95, 129, 96, 105, 98)
3
4   #n number of observation
```

```
 5  n <- 20
 6
 7  #the function will take 2 parameters lambda and x
 8  #lambda is mean
 9  #x is the sequence of observed values
10  #finally it will return the value of log-likelihood of Poisson
       distribution
11  poi.log.likelihood <- function(lambda, x){
12    n <- length(x) #n is number of obervations
13    log.like <- sum(x) * log(lambda) - n * lambda
14    return(-log.like)
15  }
16
17  #call function poi.log.likelihood
18  #passing 2 variable lambda(mean) and x
19  #store return value in variable ans.poi.log.like
20  ans.poi.log.like <- poi.log.likelihood(mean(x), x)
```

The value of log-likelihood function of Poisson distribution is **-7612.856** .

5. Using function optimize() find $\hat{\lambda}$. Compare it to the estimate you get from the exact formula.

```
21  #this function will take 2 parameters x and n
22  #x is the sequence of observed values
23  #n is number of obervation
24  #finally it will return lambda hat
25  lambda.hat <- function(x, n){
26    return(sum(x)/n)
27  }
28
29  #call function lambda.hat
30  #passing 2 variables x and n
31  #store value of estimate lambda in ans.lambda.hat
32  ans.lambda.hat <- lambda.hat(x, n)
33
34  #using optimize function to obtain lambda hat
35  #the optimize function take the log-likelihood of Poisson distribution
       to optimize with the given interval
36  #optimize function will return maximum and objective value
37  #in this case we interest in the value of maximum
38  lambda.hat.est <- optimize(f = poi.log.likelihood, interval = c(88, 129)
       , maximum = T, x = x)$maximum
```

The exact value of $\hat{\lambda}$ is **104.35** and the estimate value of $\hat{\lambda}$ is **128.9999**.
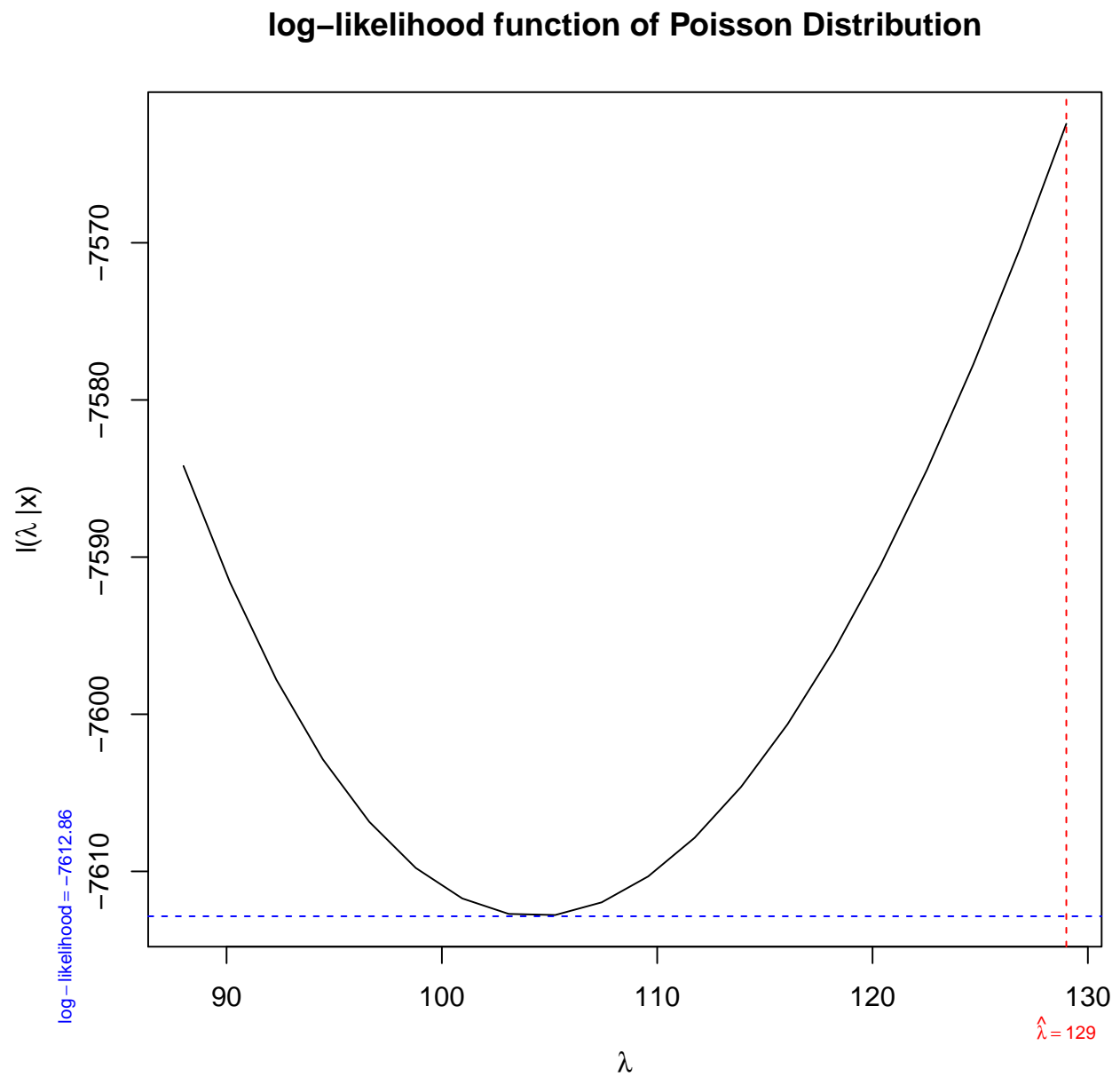Using function optimize the value of $\hat{\lambda}$ is greater than using the exact formula

6. Plot the log-likelihood function, highlight the maximum and denote the maximum likelihood estimate in plot margin.

```
39  #create sequence of lambda for x-axis from min(x) to max(x)
40  lambda.seq <- seq(from=min(x), to=max(x), length=20)
```

```
41
42  #prepare value for y-axis using function apply
43  l.lambda <- apply(X = as.matrix(lambda.seq), MARGIN = 1, FUN = poi.log.
        likelihood, x=x)
44
45  #generate plot for log-likelihood function of poisson distribution
46  plot(lambda.seq, l.lambda, type = 'l', main="log-likelihood function of
        Poisson Distribution", xlab = bquote(lambda), ylab =bquote(paste('l('
        , lambda, ' |x)')))
47
48  #set value of lambda estimation on vertical line
49  abline(v = lambda.hat.est, col = 'red', lty =2)
50
51  #set value of log-likelihood on horizontal line
52  abline(h = ans.poi.log.like, col = 'blue', lty = 2)
53
54  #use mtext to set symbol and value of lambda estimation on x-axis
55  mtext(bquote(hat(lambda) == .(round(lambda.hat.est, 2))), side = 1, line
        = 2,
56      at = lambda.hat.est, cex = 0.7, col = 'red')
57
58  #use mtext to set symbod and value of log-likelihood on y-axis
59  mtext(bquote(log-likelihood == .(round(ans.poi.log.like, 2))), side = 2,
        line = 2,
60      at = ans.poi.log.like, cex = 0.7, col = 'blue')
```

**log–likelihood function of Poisson Distribution**



## Exercise 4

1. Write down the null and the alternative hypotheses in mathematical form.

**Null hypothesis**— $H_0 : \rho = \rho_0$

**Alternative hypothesis**— $H_1 : \rho \neq \rho_0$,

where $\rho_0 = 0$

2. Calculate the value of test statistic.

```
61  #set working directory
62  setwd(getwd())
63
64  #read data from file txt
65  body <- read.table(file = 'body-measurements.txt', header = T)
66
67  #get only body height of female
68  body.f.height <- body[body$sex == 'f', 'body.H']
69
70  #get only neck of female
71  body.f.neck <- body[body$sex == 'f', 'neck.C']
72
73  #na.omit will clean data that have NA value
74  body.f.height <- na.omit(body.f.height)
75  body.f.neck <- na.omit(body.f.neck)
76
77  #get length of the data
78  n.x <- length(body.f.height)
79  n.y <- length(body.f.neck)
80  n <- n.x <- n.y
81
82  #value of rho from null hypothesis
83  rho0 <- 0
84
85  #significant level
86  alpha <- 0.05
87
88  #calculate the estimate of rho using function cor()
89  #passing vector of body.f.height and body.f.neck
90  rho.est <- cor(body.f.height, body.f.neck, method = c("pearson", "
        kendall", "spearman"))
91
92  #calculate ZR using Fisher Z -variable
93  ZR <- 1/2 * log((1+ rho.est)/(1- rho.est))
94
95  #calculate the value of xi
96  xi <- 1/2 * log((1+ rho0)/(1-rho0))
97
98  #calculate the value of test statistics
99  zW <- sqrt(n -3)*(ZR - xi)
100 # round(zW, 5) => 1.39395
```

**Result:** $zw = 1.39395$


3. Calculate the critical region and make your decision.

```
101 #calculate critical value using qnorm for lower bound
102 z.CR.l <-qnorm(alpha/2, lower.tail = T)
```

```
103  # round ( z . CR . l , 5)  => -1.95996
104
105  #calculate  critical  value  using  qnorm  for  upper  bound
106  z.CR.u <-qnorm(alpha/2, lower.tail = F)
107  # round ( z . CR . u , 5)  =>  1.95996
```

**Result:** $W = (-\infty, -1.95996) \cup (1.95996, \infty)$
The test statistic $zw = 1.39395$ does not belong to the critical region. Therefore $H_0$ is not rejected at significance level alpha $= 0.05$.

4. Calculate the p-value and make your decision.

```
108  #calculate  the  p-value  using  test  statistic  value  and  pnorm ()
109  p.val.zW <- 2 * (1 - pnorm(abs(zW)))
110  # round ( p . val . zW , 5)  =>  0.16333
```
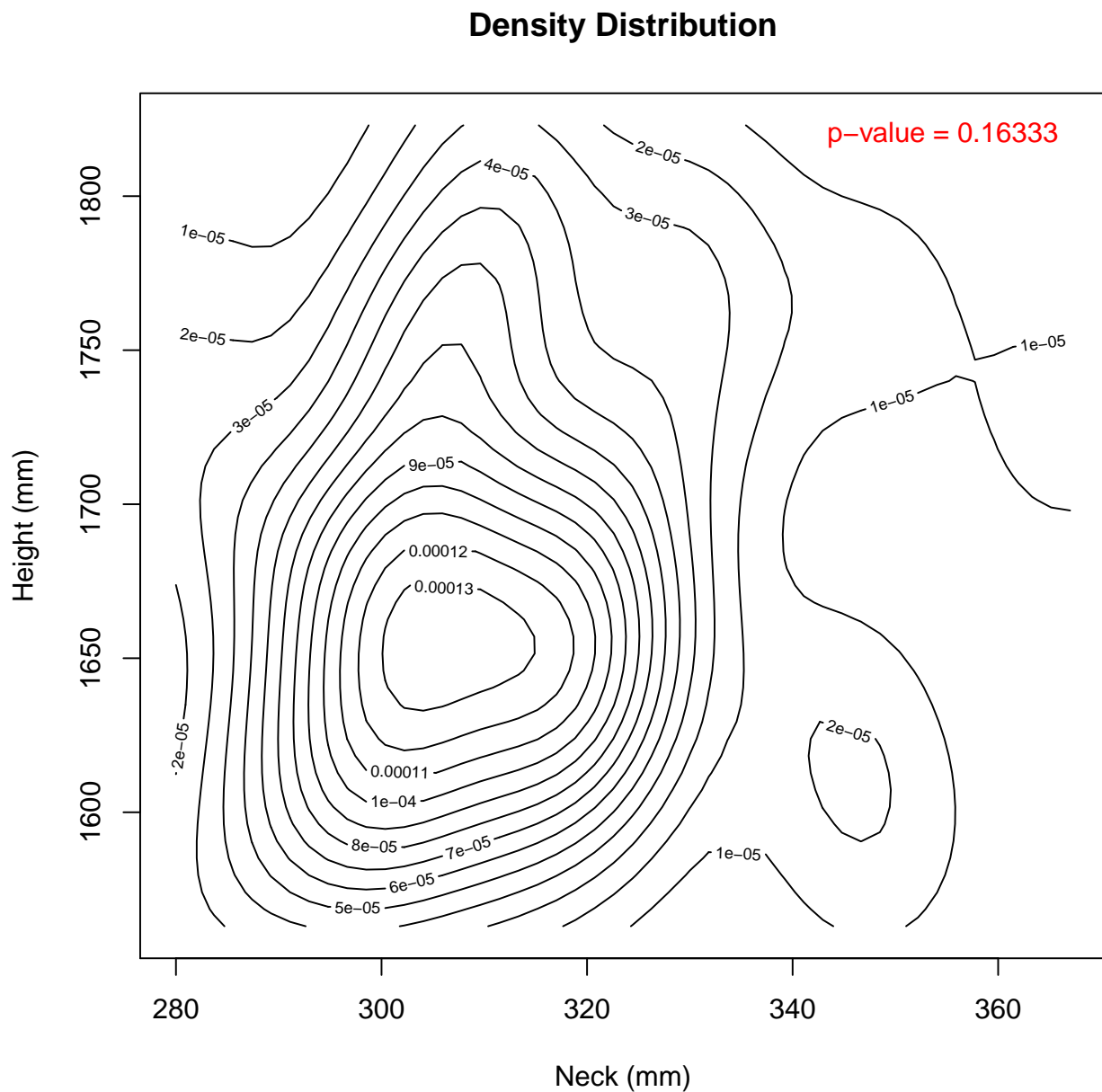
**Result:** $p - value = 0.16333$
The p-value is **0.16333** which is greater than significance level $\alpha = 0.05$. Therefore $H_0$ is not rejected at significance level alpha $= 0.05$.

5. Plot the density of the distribution, that the test statistic follows, and visualise the p-value.

```
111  #create  plot  for  bivariate  distribution
112
113  #create  data  frame  for  body  neck  and  height
114  df <- data.frame(body.f.neck, body.f.height)
115
116  #load  library  called  MASS  to  use  function  kde2d
117  library(MASS)
118
119  #call  function  kde2d  to  generate  density  of  body  neck  and  height
120  body_density <- kde2d(df$body.f.neck, df$body.f.height, n=n)
121
122  #plot  bivariate  distribution
123  plot1 <- contour(body_density, xlab='Neck (mm)', ylab= 'Height (mm)',
         main='Density Distribution')
124
125  #add  text  showing  p-value  on  top  right  of  the  plot
126  text(x = 355, y = 1820, labels = "p-value = 0.16333 ", xpd = NA, col = '
         red')
```
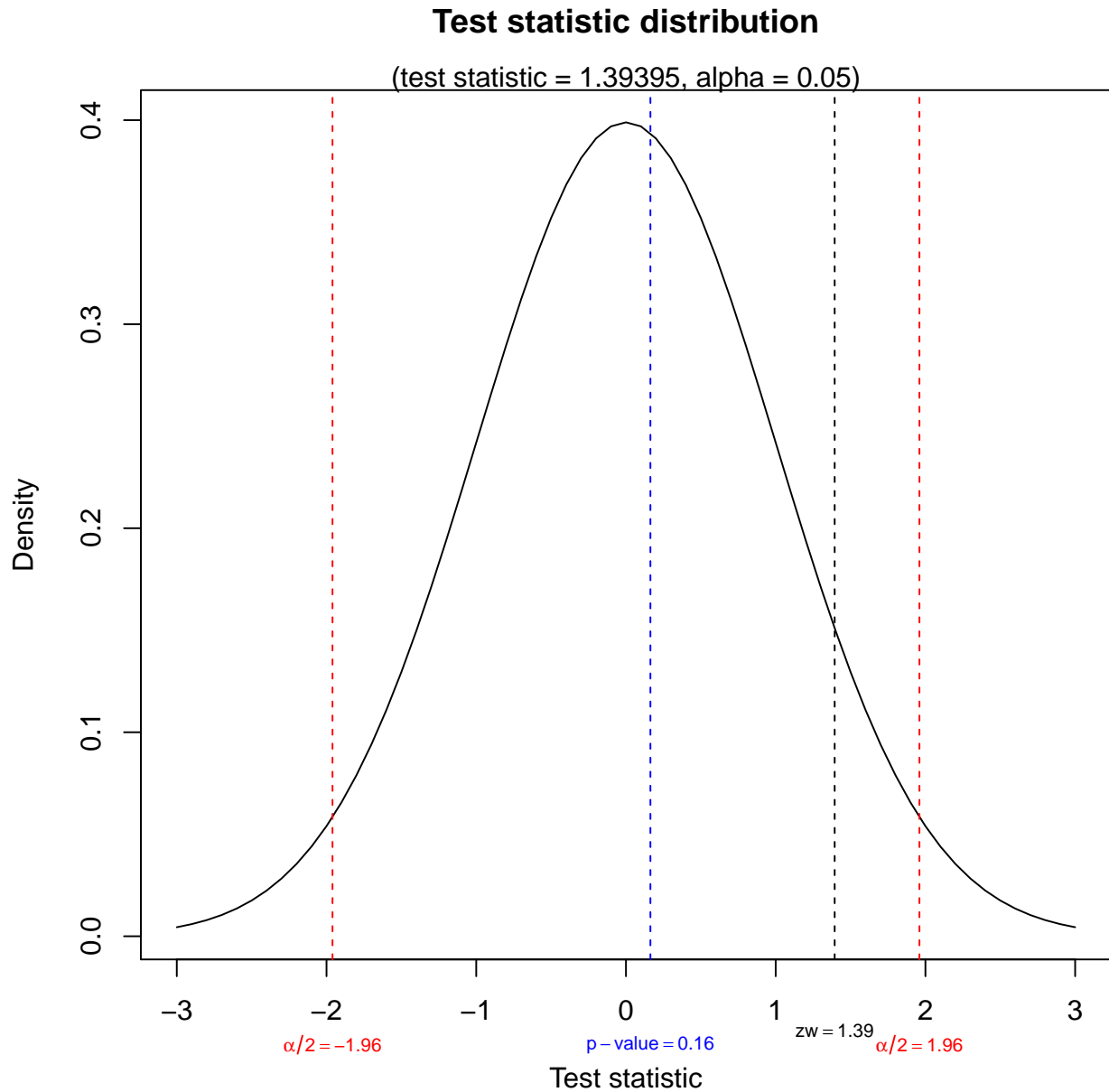
## Density Distribution



```
127  #create another plot based on test statistics density distribution
128
129  #set sequence of test statistic from -3 to 3
130  zw.seq <- seq(-3,3, by = .1)
131
132  #using dnorm to create density distribution
133  dvalues <- dnorm(zw.seq)
134
135  #plot the test statistic
136  plot2 <- plot(zw.seq, dvalues, type = 'l', xlab='Test statistic', ylab =
         'Density', main = 'Test statistic distribution')
137
```

```
138  #add text under main title
139  text(x = 0, y = 0.42, labels = '(test statistic = 1.39395, alpha = 0.05)
         ', col = 'black', xpd= NA)
140
141  #add vertical line for test statistic value and text
142  abline(v = zW, col = 'black', lty =2)
143  mtext(bquote(zw == .(round(zW, 2))), side = 1, line = 1.5,
            at = zW, cex = 0.7, col = 'black')
144
145
146  #add vertival line for lower bound critical region and text
147  abline(v = z.CR.l, col = 'red', lty =2)
148  mtext(bquote(alpha/2 == .(round(z.CR.l, 2))), side = 1, line = 2,
            at = z.CR.l, cex = 0.7, col = 'red')
149
150
151  #add vertical line for upper bound critical region and text
152  abline(v = z.CR.u, col = 'red', lty =2)
153  mtext(bquote(alpha/2 == .(round(z.CR.u, 2))), side = 1, line = 2,
            at = z.CR.u, cex = 0.7, col = 'red')
154
155
156  #add vertical line for p-value and text
157  abline(v = p.val.zW, col = 'blue', lty =2)
158  mtext(bquote(p-value == .(round(p.val.zW, 2))), side = 1, line = 2,
            at = p.val.zW, cex = 0.7, col = 'blue')
159
```

## Test statistic distribution

(test statistic = 1.39395, alpha = 0.05)



6. Calculate the confidence interval for $\rho$ and make your decision.

```
160  #calculate the confidence interval Uaplha/2 = 1.95996 for lower bound
161  CI.zW.l <- tanh(ZR - 1.95996/sqrt(n-3))
162  # round(CI.zW.l, 5) => -0.08418
163
164  #calculate the confidence interval Uaplha/2 = 1.95996 for upper bound
165  CI.zW.u <- tanh(ZR + 1.95996/sqrt(n-3))
166  # round(CI.zW.u, 5) => 0.46209
```

**Result:** $CI : (-0.08418, 0.46209)$
$\rho_0 = 0$ belongs to the value of $CI : (-0.08418, 0.46209)$. Therefore $H_0$ is not rejected at significance level alpha $= 0.05$.

7. Interpret your conclusion.

**Statistical conclusion:**

$H_0$ is not rejected on a significant level $\alpha = 0.05$, because (1) test statistics does not belong to critical region, (2) $\rho_0 = 0$ belongs to the confidence interval, and (3) p-value is greater than 0.05.

**Verbal conclusion:**

We are not rejecting null hypothesis that is stated there is no correlation between body height and neck circumference of females.

# Exercise 5

1. Write down the null and the alternative hypotheses in mathematical form.

**Null hypothesis**— $H_0 : p_1 - p_2 = p_0$

**Alternative hypothesis**— $H_1 : p_1 - p_2 \neq p_0$,

where $p_0 = 0$

2. Calculate the value of test statistic.

```
167  #setup first group data n1: number of product from group 1, x1: number
          of faulty product from group 1
168  n1 <- 200
169  x1 <- 32
170
171  #calculate the estimate probability of group 1
172  p1.est <- x1 / n1
173
174  #setup second group data n2: number of product from group 2, x2: number
          of faulty product from group 2
175  n2 <- 230
176  x2 <- 21
177
178  #calculate the estimate probability of group 2
179  p2.est <- x2 /n2
180
181  #calculate the estimate standard deviation
182  sd.est <- sqrt((p1.est*(1-p1.est))/n1 + (p2.est*(1-p2.est))/n2)
183
184  #value of p from null hypothesis
185  p0 <- 0
186
187  #significant level
188  alpha <- 0.05
```

```
189
190  #calculate the value of test statistics using Wald Test
191  zW.obs <- (p1.est - p2.est - p0)/sd.est
192  # round(zW.obs, 5) => 2.13765
```

**Result:** $zw = 2.13765$

3. Calculate the critical region and make your decision.

```
193  #calculate the critical region for lower bound
194  z.CR.l <-qnorm(alpha/2, lower.tail = T)
195  # round(z.CR.l, 5) => -1.95996
196
197  #calculate the critical region for upper bound
198  z.CR.u <-qnorm(alpha/2, lower.tail = F)
199  # round(z.CR.u, 5) => 1.95996
```

**Result:** $W = (-\infty, -1.95996) \cup (1.95996, \infty)$
The test statistic $zw = 2.13765$ belongs to the upper bound critical region. Therefore $H_0$ is rejected at significance level alpha $= 0.05$.

4. Calculate the p-value and make your decision.

```
200  #calculate the p-value using value of test statistic and pnorm()
201  p.val.zW <- 2 * (1 - pnorm(abs(zW.obs)))
202  # round(p.val.zW, 5) => 0.03255
```

**Result:** $p - value = 0.03255$
The p-value is **0.03255** which is smaller than significance level $\alpha = 0.05$. Therefore $H_0$ is rejected at significance level alpha $= 0.05$.

5. Calculate the confidence interval and make your decision.

```
203  #calculate the confidence interval for lower bound Ualpha/2 = 1.95996
204  CI.zW.l <- p1.est - p2.est - (1.95996 * sd.est)
205  # round(CI.zW.l, 5) => 0.00571
206
207  #calculate the confidence interval for upper bound Ualpha/2 = 1.95996
208  CI.zW.u <- p1.est - p2.est + (1.95996 * sd.est)
209  # round(CI.zW.u, 5) => 0.13168
```

**Result:** $CI : (0.00571, 0.13168)$
$p_0 = 0$ does not belong to the value of $CI : (0.00571, 0.13168)$. Therefore $H_0$ is rejected at significance level alpha $= 0.05$.

6. Interpret your conclusion.

**Statistical conclusion:**
$H_0$ is rejected on a significant level $\alpha = 0.05$, because (1) test statistics belongs to critical region, (2) $p_0 = 0$ does not belong to the confidence interval, and (3) p-value is less than 0.05.

**Verbal conclusion:**
We are rejecting null hypothesis that the probability of getting a faulty product are the same in both factories.