

## Mini project #4

**Group Member:** Chaoran Li, Wenting Wang

**Contribution of each member:**

Firstly, we discussed the mathematical models and code details together. Then, we divided the project into two part and finished our respective work. Wenting Wang mainly worked on Q1 while Chaoran Li worked on Q2. Then, we merged our code and solution into one report. Each member makes contribution to each sub task of this project and combines all to finish this project, as the details shown in table 1.

	Question1	Question2
Chaoran li	20%	20%
Wenting wang	80%	80%

Table 1: Member contribution table

**Question 1:**

(a) First, we set up the Null Hypothesis and Alternate Hypothesis:

$H_0$ : No difference of mean temperature between male and female:  $Tu_M - Tu_F = 0$

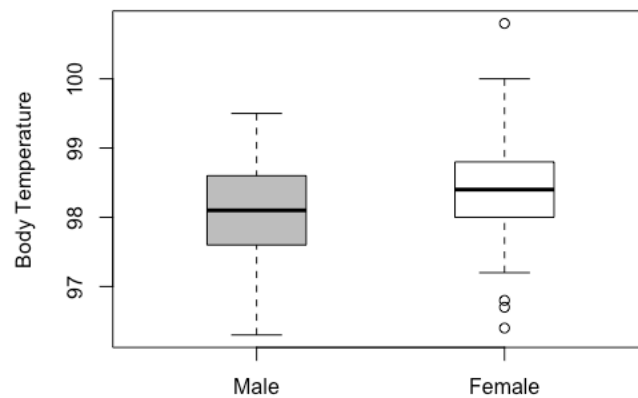
$H_1$ : There' s difference of mean temperature between male and female:  $Tu_M - Tu_F \neq$

0

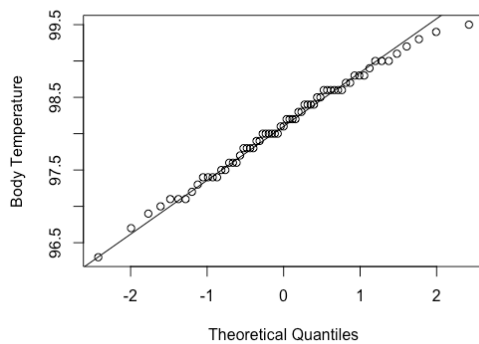
Next, we analyze the data:

The body temperatures of male and female are independent sample data, and the population standard deviation unknown. Based on the boxplots and QQ-plots below we can see the variances are different and the datasets are approximately normal.

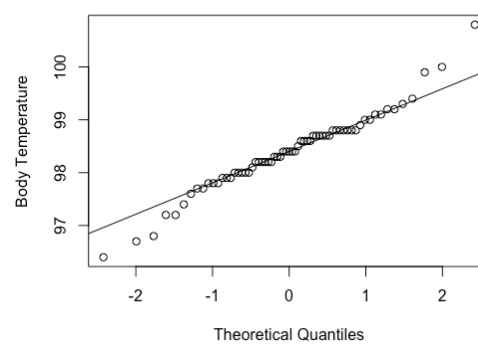
### Boxplots of Body Temperature



Q-Q Plot of Male



Q-Q Plot of Female



Thus, we choose to use t-test and get the results below:

```
welch Two sample t-test

data: male$body_temperature and female$body_temperature
t = -2.2854, df = 127.51, p-value = 0.02394
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.53964856 -0.03881298
sample estimates:
mean of x mean of y
 98.10462  98.39385
```

Then we can conclude that because the 95% confidence interval of difference between male and female temperature is  $[-0.53964856, -0.03881298]$ , which does not include 0. The P-value is 0.02394 which is less than 0.05. Thus, we reject the null hypothesis, which means there exists difference of mean temperature between male and female. And from the sample means, we can notice that the mean temperature of female is a little higher than male.

Code:

```

3 # (a) Body temperature
4 bodytemp.heartrate <- read.csv(
5   file=file.path("./Mini Project 5/bodytemp-heartrate.csv"))
6 male <-bodytemp.heartrate[which(bodytemp.heartrate$gender==1),]
7 female <- bodytemp.heartrate[which(bodytemp.heartrate$gender==2),]
8 #plots body temperature
9 boxplot(male$body_temperature, female$body_temperature, boxwex=.4,
10         main = "Boxplots of Body Temperature", names = c('Male', 'Female'),
11         col=c("grey", "white"), ylab="Body Temperature")
12 qqnorm(male$body_temperature, main='Q-Q Plot of Male', ylab="Body Temperature")
13 qqline(male$body_temperature)
14 qqnorm(female$body_temperature, main='Q-Q Plot of Female',
15         ylab="Body Temperature")
16 qqline(female$body_temperature)
17 #t,test function for the body temperature difference
18 t.test(male$body_temperature, female$body_temperature,
19        alternative = 'two.sided', var.equal = F)

```

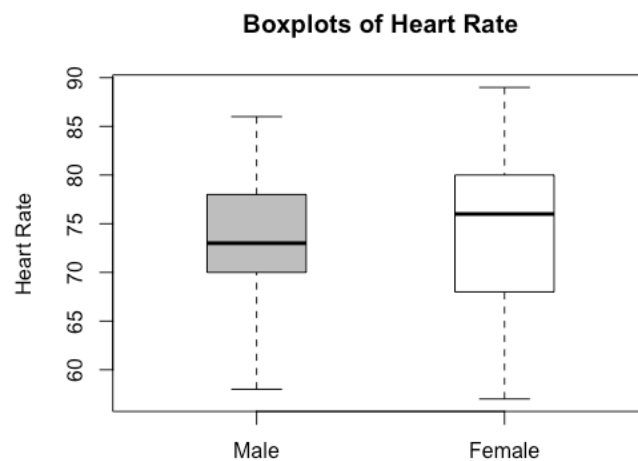
(b) First, we set up the Null Hypothesis and Alternate Hypothesis:

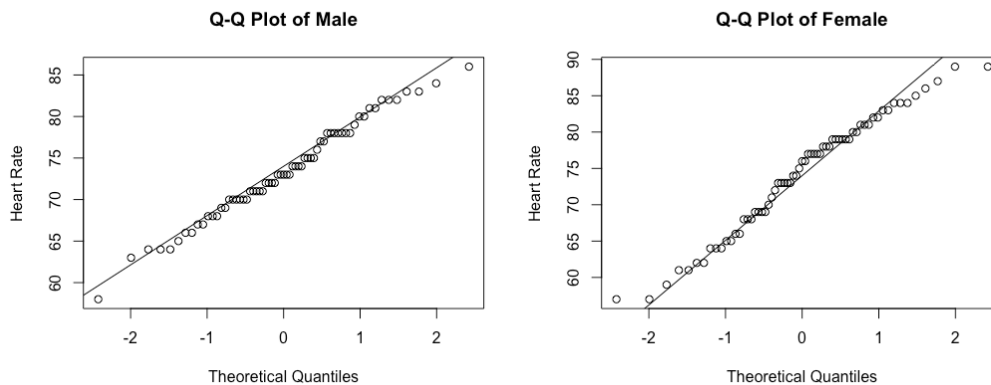
$H_0$ : No difference of mean heart rate between male and female:  $Hu_M - Hu_F = 0$

$H_1$ : There' s difference of mean heart rate between male and female:  $Hu_M - Hu_F \neq 0$

Next, we analyze the data:

The heart rates of male and female are independent sample data, and the population standard deviation unknown. Based on the boxplots and QQ-plots below we can see the variances are different and the datasets are approximately normal.





welch Two sample t-test

```
data: male$heart_rate and female$heart_rate
t = -0.63191, df = 116.7, p-value = 0.5287
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -3.243732  1.674501
sample estimates:
mean of x mean of y
 73.36923  74.15385
```

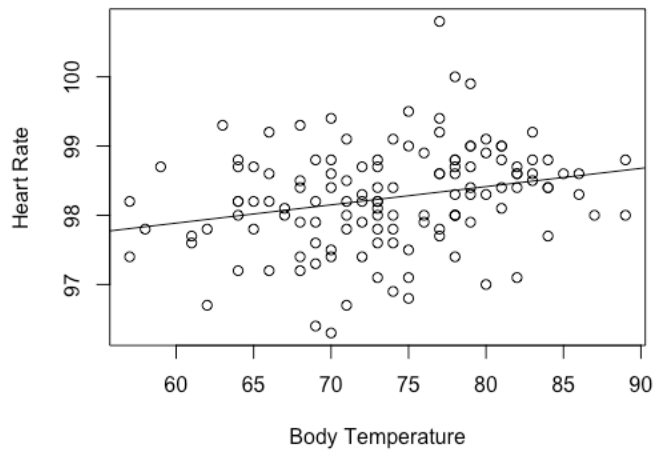
Then we can conclude that because the 95% confidence interval of difference between male and female heart rate is  $[-3.243732, 1.674501]$ , which does include 0. The P-value is 0.5287 which is greater than 0.05. Thus, we accept the null hypothesis, which means there is no difference of mean heart rates between male and female.

Code:

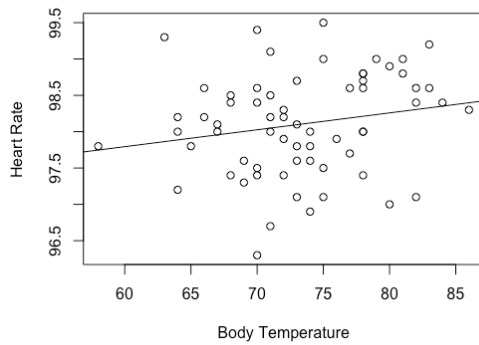
```
21 # (b) Heart rate
22 boxplot(male$heart_rate, female$heart_rate, boxwex=.4,
23         main = "Boxplots of Heart Rate", names = c('Male', 'Female'),
24         col=c("grey", "white"), ylab="Heart Rate")
25 qqnorm(male$heart_rate, main="Q-Q Plot of Male", ylab="Heart Rate")
26 qqline(male$heart_rate)
27 qqnorm(female$heart_rate, main="Q-Q Plot of Female", ylab="Heart Rate")
28 qqline(female$heart_rate)
29 #t.test function for the body temperature difference
30 t.test(male$heart_rate, female$heart_rate,
31        alternative='two.sided', var.equal=F)
```

(c) First, we draw the scatter plot and regression line below

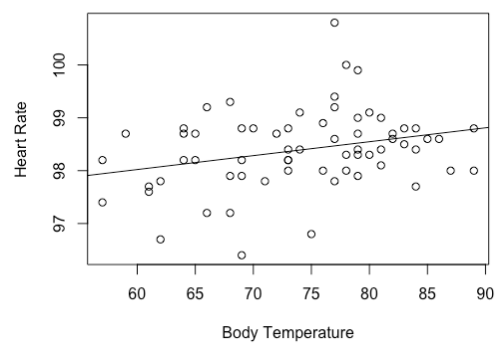
**Body Temperature and Heart Rate**



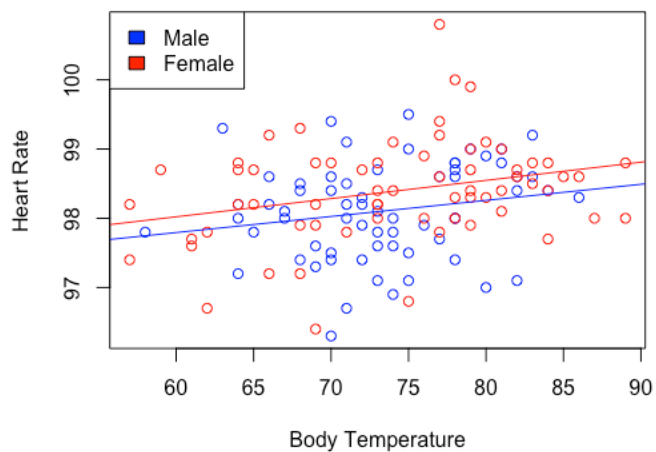
**Scatter Plot for Male**



**Scatter Plot for Female**



**Scatter Plot for Male and Female**



We can see that in both the plot of all people and the plots of male and female, there exists a weak positive trend between body temperature and heart rate. We can also see some outliers in the scatter plots. Then we can calculate the correlation between body temperature and heart rate in the 3 plots.

Cor(all) = 0.2536564, Cor(male) = 0.1955894, Cor(female) = 0.2869312

We can also get the simple linear regression as below:

```
> #get the fitted regression line
> lm(bodytemp.heartrate$body_temperature ~ bodytemp.heartrate$heart_rate)
```

Call:

```
lm(formula = bodytemp.heartrate$body_temperature ~ bodytemp.heartrate$heart_rate)
```

Coefficients:

(Intercept)	bodytemp.heartrate\$heart_rate
96.30675	0.02633

```
> lm(male$body_temperature ~ male$heart_rate)
```

Call:

```
lm(formula = male$body_temperature ~ male$heart_rate)
```

Coefficients:

(Intercept)	male\$heart_rate
96.39789	0.02326

```
> lm(female$body_temperature ~ female$heart_rate)
```

Call:

```
lm(formula = female$body_temperature ~ female$heart_rate)
```

Coefficients:

(Intercept)	female\$heart_rate
96.44211	0.02632

Above all, we can conclude that no matter male or female there is weak positive relationship between body temperature and heart rate. But the correlation of female is a little stronger than male.

Code:

```
33 # (c) Compare male and female
34 plot(x=bodytemp.heartrate$heart_rate, y=bodytemp.heartrate$body_temperature,
35      type="p", xlab="Body Temperature", ylab="Heart Rate",
36      main="Body Temperature and Heart Rate")
37 abline(lm(bodytemp.heartrate$body_temperature ~ bodytemp.heartrate$heart_rate))
38 cor(bodytemp.heartrate$body_temperature, bodytemp.heartrate$heart_rate)
39 #get the fitted regression line
40 lm(bodytemp.heartrate$body_temperature ~ bodytemp.heartrate$heart_rate)
41 #abline(temp)
42
43 #Scatter plots for body temperature and heart rate for males
44 plot(male$heart_rate, male$body_temperature,
45      type="p", xlab="Body Temperature", ylab="Heart Rate",
46      main="Scatter Plot for Male")
47 abline(lm(male$body_temperature ~ male$heart_rate))
48 cor(male$body_temperature, male$heart_rate)
49 lm(male$body_temperature ~ male$heart_rate)
```

```

50 # for females
51 plot(female$heart_rate, female$body_temperature, pch=1,
52      type="p", xlab="Body Temperature", ylab="Heart Rate",
53      main="Scatter Plot for Female")
54 abline(lm(female$body_temperature ~ female$heart_rate))
55 cor(female$body_temperature, female$heart_rate)
56 lm(female$body_temperature ~ female$heart_rate)
57 # mixed scatter plots
58 plot(bodytemp.heartrate$heart_rate, bodytemp.heartrate$body_temperature,
59      col=c('blue', 'red')[unclass(bodytemp.heartrate$gender)],
60      type="p", xlab="Body Temperature", ylab="Heart Rate",
61      main="Scatter Plot for Male and Female")
62 legend("bottomright", c("Male", "Female"), fill=c('blue', 'red'))
63 abline(lm(male$body_temperature ~ male$heart_rate), col='blue')
64 abline(lm(female$body_temperature ~ female$heart_rate), col='red')

```

## Question 2:

Compare large-sample z-interval (interval 1) and bootstrap percentile method interval (interval 2) from an exponential distribution.

(a) Given  $\lambda = 0.1$  and  $n = 30$ , compute two accuracies of two intervals which confidence interval include the mean.

Firstly, we design two functions to return the result whether the confidence interval include the mean in one investigation.

For interval 1, even though we could know the distribution belongs to exponential distribution, the standard deviation or variance are unknown. Hence, we have to use the sample's standard deviation or variance. Besides, even though  $n = 5$  or  $10$  is not enough for a large-sample assumption, we will still use `qnorm()` for getting large-sample z-interval.

```

4 # Question 2
5 # Large-sample z-interval and bootstrap percentile method interval.
6
7 # (a) Given lambda = 0.1 and n = 30, compare two intervals' accuracy.
8 # 1 - alpha = 0.95 => [0.025, 0.975].
9
10 # Large-sample z-interval
11 getInterval1 <- function(n, lambda) {
12   population.mean = 1/lambda
13   sample = rexp(n, rate=lambda)
14   center = mean(sample)
15   margin = qnorm(0.975) * sd(sample) / sqrt(n)
16   if (center - margin > population.mean) {
17     return (0)
18   } else if (center + margin < population.mean) {
19     return (0)
20   } else {
21     return (1)
22   }
23 }

```

For interval 2, we know the distribution is exponential distribution. Hence, we would use parametric bootstrap here.

```

25 # Bootstrap percentile method interval
26 getInterval2 <- function(n, lambda) { #
27   population.mean = 1/lambda
28   sample = rexp(n, rate=lambda)
29   lambda.est = 1 / mean(sample)
30   sample.boot = c(mean(sample),
31     replicate(999, expr=mean(rexp(n, rate=lambda.est)))) # 1000 times
32   sample.boot = sort(sample.boot)
33   if (sample.boot[25] > population.mean) {
34     return (0)
35   } else if (sample.boot[975] < population.mean) {
36     return (0)
37   } else {
38     return (1)
39   }
40 }

```

Secondly, we use the given parameters to calculate the two accuracy of two intervals during repeating 5000 times.

```

42 q2.round = 5000
43 q2.a.lambda = 0.1
44 q2.a.n = 30
45 q2.a.interval1s = replicate(q2.round, expr=getInterval1(q2.a.n, q2.a.lambda))
46 q2.a.accuracy1 = sum(q2.a.interval1s) / q2.round
47 q2.a.interval2s = replicate(q2.round, expr=getInterval2(q2.a.n, q2.a.lambda))
48 q2.a.accuracy2 = sum(q2.a.interval2s) / q2.round
49 sprintf("Lambda = %f, n = %d. Accuracy:", q2.a.lambda, q2.a.n)
50 sprintf("interval 1: %f; interval 2 %f.", q2.a.accuracy1, q2.a.accuracy2)

```

Then, we get the results below:



```

> q2.round = 5000
> q2.a.lambda = 0.1
> q2.a.n = 30
> q2.a.interval1s = replicate(q2.round, expr=getInterval1(q2.a.n, q2.a.lambda))
> q2.a.accuracy1 = sum(q2.a.interval1s) / q2.round
> q2.a.interval2s = replicate(q2.round, expr=getInterval2(q2.a.n, q2.a.lambda))
> q2.a.accuracy2 = sum(q2.a.interval2s) / q2.round
> sprintf("Lambda = %f, n = %d. Accuracy:", q2.a.lambda, q2.a.n)
[1] "Lambda = 0.100000, n = 30. Accuracy:"
> sprintf("interval 1: %f; interval 2 %f.", q2.a.accuracy1, q2.a.accuracy2)
[1] "interval 1: 0.915200; interval 2 0.935600."

```

The accuracy of interval 2 which is 0.9356 is larger than that of interval 1 which is 0.9152.

(b) Repeat (a) for the remaining combinations of  $n \in \{5, 10, 30, 100\}$  and  $\lambda \in \{0.01, 0.1, 1, 10\}$ .

```

52 # (b) Traverse all lambdas and ns.
53 q2.b.lambdas = c(0.01, 0.1, 1, 10)
54 q2.b.ns = c(5, 10, 30, 100)
55 q2.intervals = c("interval 1", "interval 2")
56 q2.b accuracies = array(0,
57   dim=c(length(q2.b.lambdas), length(q2.b.ns), length(q2.intervals)))
58 for (i in 1:length(q2.b.lambdas)) {
59   lambda = q2.b.lambdas[i]
60   for (j in 1:length(q2.b.ns)) {
61     n = q2.b.ns[j]
62     interval1s = replicate(q2.round, expr=getInterval1(n, lambda))
63     accuracy1 = sum(interval1s) / q2.round
64     q2.b accuracies[i, j, 1] = sum(interval1s) / q2.round
65     interval2s = replicate(q2.round, expr=getInterval2(n, lambda))
66     accuracy2 = sum(interval2s) / q2.round
67     q2.b accuracies[i, j, 2] = sum(interval2s) / q2.round
68   }
69 }
70 print(q2.b accuracies)

```

```

> print(q2.b accuracies)
, , 1

      [,1] [,2] [,3] [,4]
[1,] 0.8066 0.8754 0.9152 0.9424
[2,] 0.8156 0.8778 0.9132 0.9386
[3,] 0.8096 0.8594 0.9182 0.9356
[4,] 0.8178 0.8642 0.9194 0.9414

, , 2

      [,1] [,2] [,3] [,4]
[1,] 0.8946 0.9228 0.9436 0.9514
[2,] 0.8972 0.9230 0.9372 0.9472
[3,] 0.8902 0.9174 0.9394 0.9454
[4,] 0.9082 0.9176 0.9416 0.9456

```

For  $\lambda \in \{0.01, 0.1, 1, 10\}$ , it is unclear to show the result in regular scale. So, I use ggplot2 library to make x coordinate log scale when talking about different  $\lambda$ .

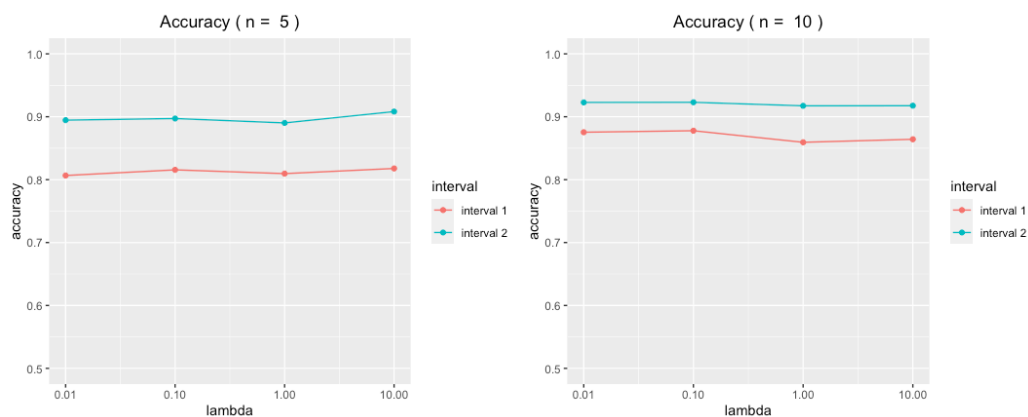
```

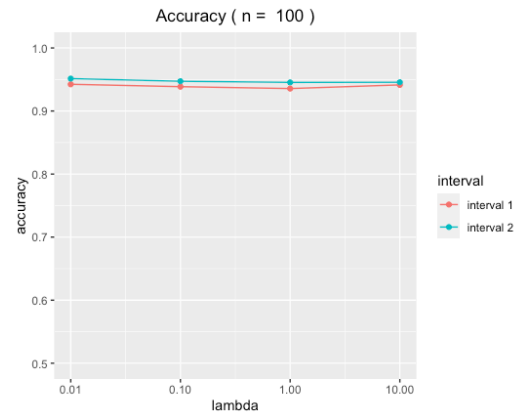
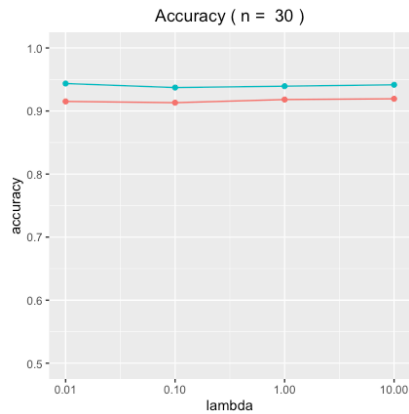
72 library(ggplot2)
73 # Different lambdas and one n.
74 for (j in 1:length(q2.b.ns)) {
75   df = data.frame(
76     lambda = rep(q2.b.lambdas, times=length(q2.intervals)),
77     interval = rep(q2.intervals, each=length(q2.b.lambdas)),
78     accuracy = c(q2.b accuracies[, j, ]))
79   img = ggplot(data=df, mapping=aes(x=lambda, y=accuracy, colour=interval)) +
80     geom_line() + geom_point() + # line plus scatter
81     scale_x_continuous(trans='log10') + # log scale in x coordinate
82     ylim(0.5, 1) + # y range: [0, 1]
83     labs(title=paste("Accuracy ( n = ", q2.b.ns[j], ")")) + # add title
84     theme(plot.title = element_text(hjust = 0.5)) # center the title
85   print(img)
86 }
87 # One lambda and different ns.
88 for (i in 1:length(q2.b.lambdas)) {
89   df = data.frame(
90     n = rep(q2.b.ns, times=length(q2.intervals)),
91     interval = rep(q2.intervals, each=length(q2.b.ns)),
92     accuracy = c(q2.b accuracies[i, , ]))
93   img = ggplot(data=df, mapping=aes(x=n, y=accuracy, colour=interval)) +
94     geom_line() + geom_point() + # line plus scatter
95     ylim(0.5, 1) + # y range: [0, 1]
96     labs(title=paste("Accuracy ( lambda = ", q2.b.lambdas[i], ")")) + # add title
97     theme(plot.title = element_text(hjust = 0.5)) # center the title
98   print(img)
99 }

```

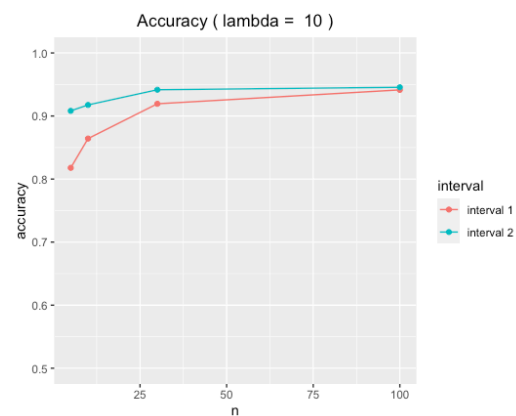
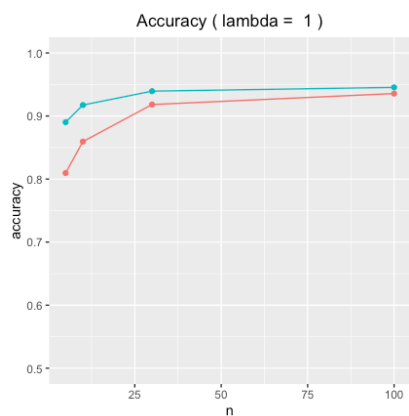
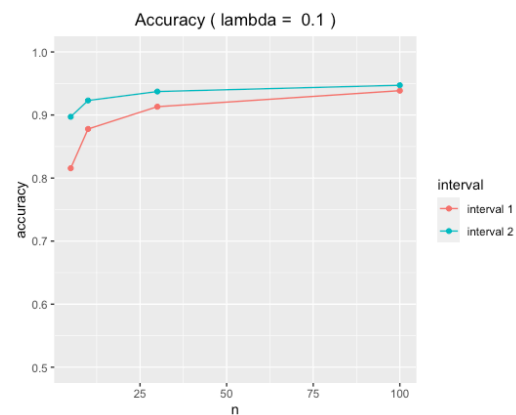
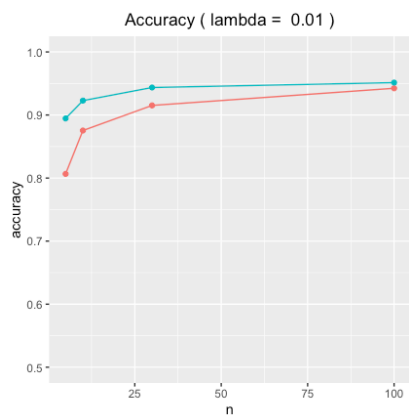
Then, we got the diagrams below:

Discuss the influence of  $\lambda$  with each n.





Discuss the influence of  $n$  with each  $\lambda$ .



From the eight diagrams above, we can make a brief summary now. The accuracy is independent from  $\lambda$ . While  $n$  increases, the accuracy increases and approaches convergence.

(c) Answer the following questions:

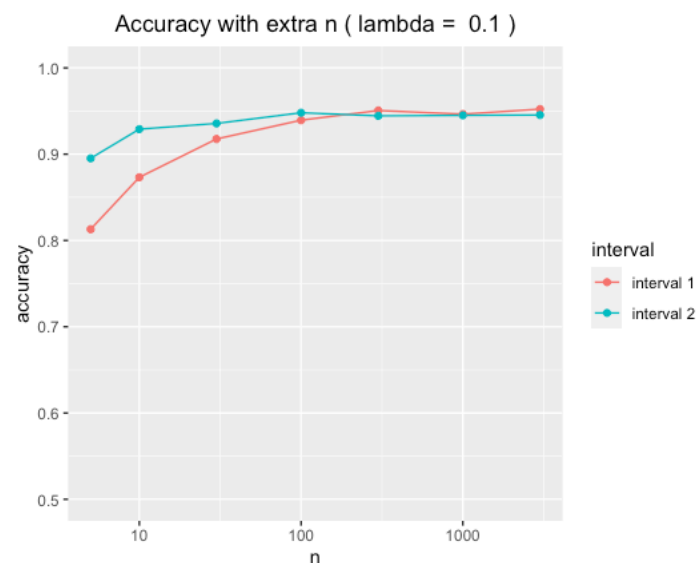
(1) How large  $n$  is needed for the large-sample z-interval?

$n = 300$  is large enough for the large-sample z-interval.

From the first four diagrams in (b), we can see that  $\lambda$  has no effect on its accuracy. However,

we cannot make sure whether the function converges at  $n = 100$  in interval 1. I tried some bigger  $n$  below:

```
101 # (c) Analyse the result we got in (b)
102 q2.c.lambda = 0.1
103 q2.c.ns = c(5, 10, 30, 100, 300, 1000, 3000)
104 q2.c accuracies = array(0,
105   dim=c(length(q2.c.ns), length(q2.intervals)))
106 for (i in 1:length(q2.c.ns)) {
107   n = q2.c.ns[i]
108   interval1s = replicate(q2.round, expr=getInterval1(n, q2.c.lambda))
109   accuracy1 = sum(interval1s) / q2.round
110   q2.c accuracies[i, 1] = sum(interval1s) / q2.round
111   interval2s = replicate(q2.round, expr=getInterval2(n, q2.c.lambda))
112   accuracy2 = sum(interval2s) / q2.round
113   q2.c accuracies[i, 2] = sum(interval2s) / q2.round
114 }
115 print(q2.c accuracies)
116
117 q2.c.df = data.frame(
118   n = rep(q2.c.ns, times=length(q2.intervals)),
119   interval = rep(q2.intervals, each=length(q2.c.ns)),
120   accuracy = c(q2.c accuracies[, ]))
121 ggplot(data=q2.c.df, mapping=aes(x=n, y=accuracy, colour=interval)) +
122   geom_line() + geom_point() + # line plus scatter
123   scale_x_continuous(trans='log10') + # log scale in x coordinate
124   ylim(0.5, 1) + # y range: [0, 1]
125   labs(title=paste(
126     "Accuracy with extra n ( lambda = ", q2.c.lambda, ")") + # add title
127   theme(plot.title = element_text(hjust = 0.5)) # center the title
> print(q2.c accuracies)
      [,1] [,2]
[1,] 0.8130 0.8952
[2,] 0.8732 0.9290
[3,] 0.9176 0.9356
[4,] 0.9394 0.9480
[5,] 0.9506 0.9444
[6,] 0.9464 0.9450
[7,] 0.9522 0.9454
```



From the diagram above, we think  $n = 300$  is large enough in interval 1.

(2) How large  $n$  is needed for bootstrap percentile method interval?

$n = 30$  is large enough for bootstrap percentile method interval.

From the diagrams in (b) and (c), we can see that the function converges at  $n = 30$  in interval 2.

(3) Do these answers depend on  $\lambda$ ?

No. From the first four diagrams we got in (b), The change on  $\lambda$  does not affect the result. Also, thinking from mathematics way, the change on  $\lambda$  does not change the overall shape of exponential distribution. Hence, the relationship between the mean and the confidence interval remains. The accuracy would not change by  $\lambda$ .

(4) Can we say that one method is more accurate than the other?

Yes, interval 2 is better than interval 1. From all eight diagrams, the accuracy of interval 2 is always larger than interval 1.

(5) Which interval would you recommend?

I would recommend interval 2 if we are talking about the accuracy of the confidence interval, especially for small  $n$ . However, if  $n$  is large enough, for example  $n = 100$ , calculation for interval 1 is quicker than interval 2 while the accuracy is close.

Overall, the bootstrap is good for estimating when the sample size is not large enough. However, it takes more resource to calculate while extra sampling in bootstrap. This can be seen from question (1) and (2).

(d) Do your conclusion in (c) depend on  $\lambda$  that were fixed in advance?

Actually, question (d) is quite similar the question (c) (3).

For different  $\lambda$ , the accuracy does not change. Hence:

Answer for (1) remains: The  $n = 300$  is large enough for large-sample z-interval.

Answer for (2) remains: The  $n = 30$  is large enough for bootstrap percentile method interval.

Answer for (3) remains: The change on  $\lambda$  does not affect the relationship between the mean and the confidence interval. All answers do not depend on  $\lambda$ .

Answer for (4) remains: The interval 2 is always more accurate. For its accuracy is always larger than interval 1's.

Answer for (5) remains: Still recommend interval 2. Only exceptional for large  $n$  on the consideration of calculation efficiency.