

生物统计学作业整理

计算所-刘栋梁

2019-06-20

生物统计学作业整理

计算所-刘栋梁

2019-06-20

Preface

homework-3

question 1

R language application (25')

R language application. (25')

question 4

homework-4

question 1

question 2

question 3

question 4

homework-5

question-1

question-2

question-3

question-4

homework-6

question-1

question-2

question-3

question 4

question 5

question 6

Preface

只对作业三到六进行了整理，作业一、二和七相对不重要，大家直接看答案即可。

肯定有不正确的地方，欢迎大家讨论，预祝大家考试顺利。

homework-3

question 1

一、实验室欲购进一批灯泡，打算在两个供货商之间选择一家购买。选购考虑的主要因素就是灯泡使用寿命的方差大小，为此需要对供货商提供的20个样品进行检测，得到的数据如下表所示。(20')

供货商1	供货商2
6802	5884
5730	5871
5823	5797
5915	5957
5774	5803
5880	5862
5870	5814
5773	5885
5830	5856
5841	5940
5763	5945
5851	5803
5789	5864
5796	5851
5818	5714
5685	5943
5602	5830
5841	5858
5723	5922
5757	5866

1、检验两家供货商的灯泡使用寿命的方差有无显著差异 ($\alpha=0.05$) (10')

解：检验两供应商的灯泡使用寿命的方差有无显著差异即为两个样本方差的同质性检验，可用F检验。

(1) .假设 $H_0 : \sigma_1^2 = \sigma_2^2$, 即两供应商的灯泡使用寿命的方差无显著差异; $H_A : \sigma_1^2 \neq \sigma_2^2$ 。

(2) .确定显著水平 $\alpha=0.05$ 。

(3) .检验计算：

```

1 sup_1<-
  c(6802,5730,5823,5915,5774,5880,5870,5773,5830,5841,5763,5851,5789,5796,5818,5685,5602,
    5841,5723,5757)
2 sup_2<-
  c(5884,5871,5797,5957,5803,5862,5814,5885,5856,5940,5945,5803,5864,5851,5714,5943,5830,
    5858,5922,5866)
3 n1<-length(sup_1)
4 n2<-length(sup_2)
5 var1<-var(sup_1)
6 var2<-var(sup_2)
7 F<-var1/var2
8 F

```

```

1 ## [1] 15.2795

```

```

1 p<-2*(1-pf(F,df1=(n1-1),df2=(n2-1)))
2 p

```

```

1 ## [1] 1.799681e-07

```

答：拒绝原假设，认为两家供货商灯泡使用寿命的方差存在显著差异。

2、选择最合适的检验方法检验两家供应商的灯泡使用寿命有无差别。（10'）

解：检验两家供应商的灯泡使用寿命有无差别即总体方差未知，且两样本都属于小样本，故应该用两组平均数差异显著性的双尾t检验。

先做方差齐性检验，第一问其实已经得出结论：方差有显著差异。当然我们也可以直接通过stats包中的var.test进行检测。

假设 $H_0 : \sigma_1^2 = \sigma_2^2$,即两供应商的灯泡使用寿命的方差无显著差异; $H_A : \sigma_1^2 \neq \sigma_2^2$ 。

```

1 var.test(sup_1,sup_2,conf.level = 0.95,alternative = "two.sided")

```

```

1 ##
2 ## F test to compare two variances
3 ##
4 ## data: sup_1 and sup_2
5 ## F = 15.279, num df = 19, denom df = 19, p-value = 1.8e-07
6 ## alternative hypothesis: true ratio of variances is not equal to 1
7 ## 95 percent confidence interval:
8 ## 6.047811 38.602899
9 ## sample estimates:
10 ## ratio of variances
11 ## 15.2795

```

因为两个供货商的方差存在显著差异，所以var.equal=false。以下做t检验。

原假设：两供应商的灯泡使用寿命无显著差别。

备择假设：两供应商的灯泡使用寿命有显著差别。

```
1 | t.test(sup_1,sup_2,var.equal=FALSE,alternative = "two.sided",conf.level = 0.95)
```

```
1 | ##
2 | ##  welch Two Sample t-test
3 | ##
4 | ## data:  sup_1 and sup_2
5 | ## t = -0.36748, df = 21.476, p-value = 0.7169
6 | ## alternative hypothesis: true difference in means is not equal to 0
7 | ## 95 percent confidence interval:
8 | ##  -133.69464   93.49464
9 | ## sample estimates:
10 | ## mean of x mean of y
11 | ##    5843.15    5863.25
```

答：因为 $p=0.7169>0.05$,因此不拒绝原假设，认为两供应商的灯泡使用寿命无显著差别。

R language application (25')

Please use R to resolve the following issues and display your R code and results.

1. For a normal random variable X with mean 4.0, and standard deviation 1.0,

a)find the probability that X is less than 2.0. (4')

b)find the value K so that $P(X>K) = 0.05$. (4')

```
1 | # find the probability that x is less than 2.0.
2 | pnorm(2,mean = 4,sd = 1)
```

```
1 | ## [1] 0.02275013
```

```
1 | # find the value K so that P(X>K) = 0.05.
2 | qnorm(0.05,mean = 4,sd =1,lower.tail = FALSE)
```

```
1 | ## [1] 5.644854
```

2. When tossing a fair coin 8 times,

a)find the probability of seeing no heads (Hint: this is a binomial distribution.) (3')

b)find the probability of seeing exactly 4 heads. (3')

c)find the probability of seeing more than 5 heads. (3')

It is a fair coin,so the prop = 0.5.

```
1 | # find the probability of seeing no heads
2 | dbinom(0,8,0.5)
```

```
1 ## [1] 0.00390625
```

```
1 # find the probability of seeing exactly 4 heads.  
2 dbinom(4,8,0.5)
```

```
1 ## [1] 0.2734375
```

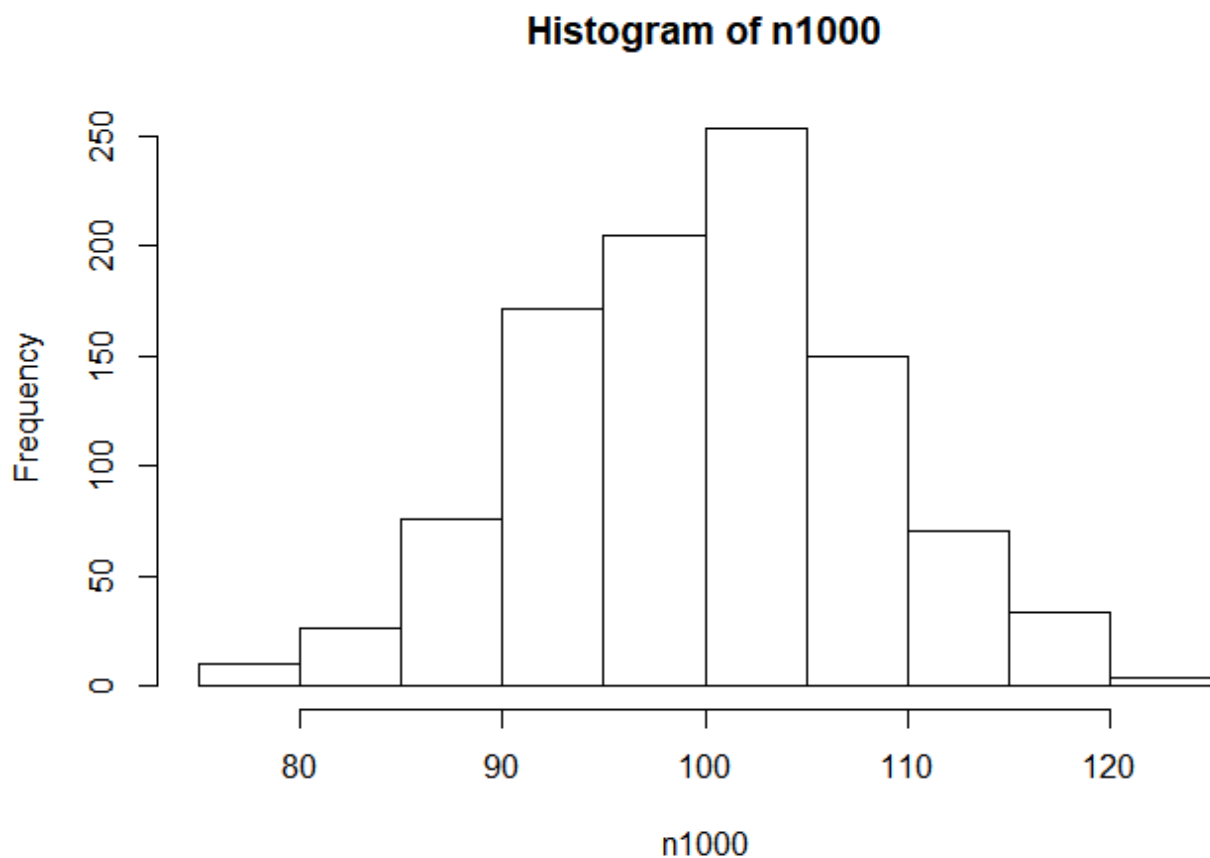
```
1 # find the probability of seeing more than 5 heads.  
2 pbinom(5,8,0.5,lower.tail = FALSE)
```

```
1 ## [1] 0.1445313
```

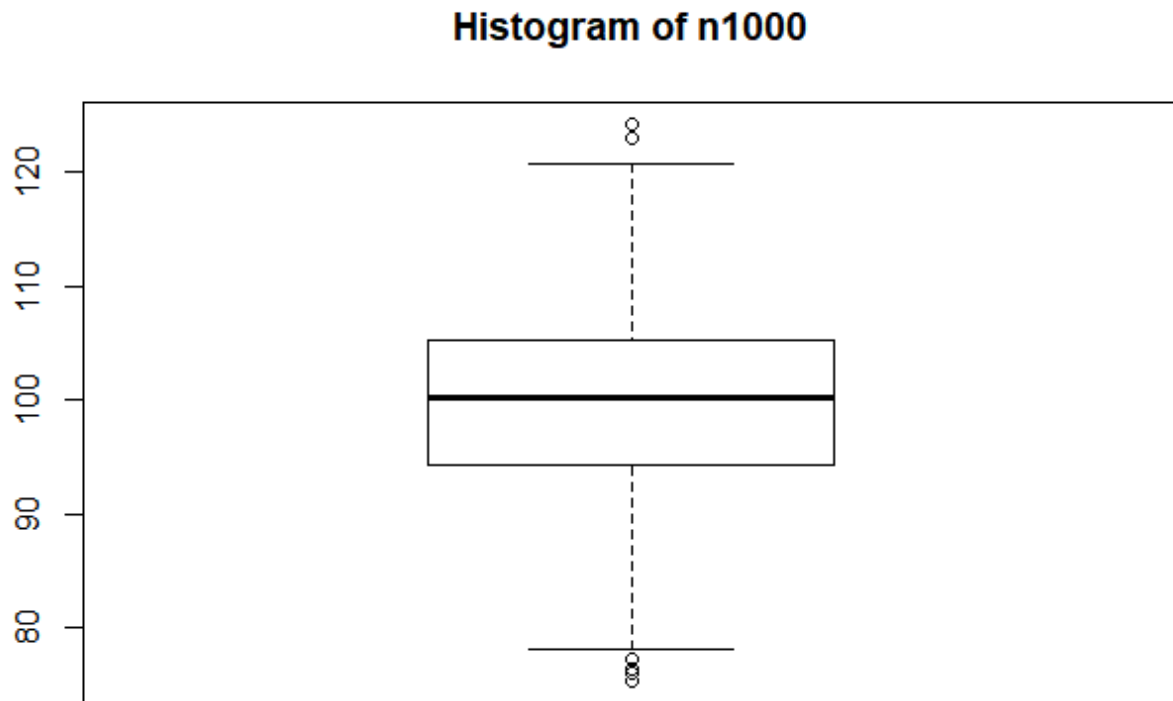
3. Simulate a sample of 1000 random data points from a normal distribution with mean 100 and standard deviation 8, and store the result in a vector.

a) plot a histogram and a boxplot of the vector you just created. (4')

```
1 set.seed(46)  
2 n1000 <- rnorm(1000, mean = 100, sd = 8)  
3 hist(n1000, main = "Histogram of n1000")
```



```
1 | boxplot(n1000,main = "Histogram of n1000")
```



b)using the data above, test the hypothesis that the mean equals 100 (using t.test). (4')

- (1) .假设 $H_0 : \mu_1 = \mu_2$,即样本与总体均值之间并无显著差异; $H_A : \mu_1 \neq \mu_2$ 。
- (2) .确定显著水平 $\alpha=0.05$ 。
- (3) .检验计算:

```
1 | t.test(n1000,mu=100)
```

```
1  ##
2  ##  One Sample t-test
3  ##
4  ## data:  n1000
5  ## t = -0.31056, df = 999, p-value = 0.7562
6  ## alternative hypothesis: true mean is not equal to 100
7  ## 95 percent confidence interval:
8  ##   99.40978 100.42893
9  ## sample estimates:
10 ## mean of x
11 ##   99.91936
```

因为p值大于0.05, 所以不拒绝原假设, 认为样本与总体之间并无显著差异。

R language application. (25')

(Please use R to read-in and manipulate data, code and results should be displayed.)

In order to detect air quality, a city's environmental protection department conducted a random test of PM2.5 in the air every few weeks. It is known that the average value of PM2.5 per cubic meter of air in the city is 82ug/m3. In the most recent test, the value of PM2.5 per cubic meter of air(ug/m3) is shown in the homework3_data.

1)Show your work directory (2')

```
1 | getwd()
```

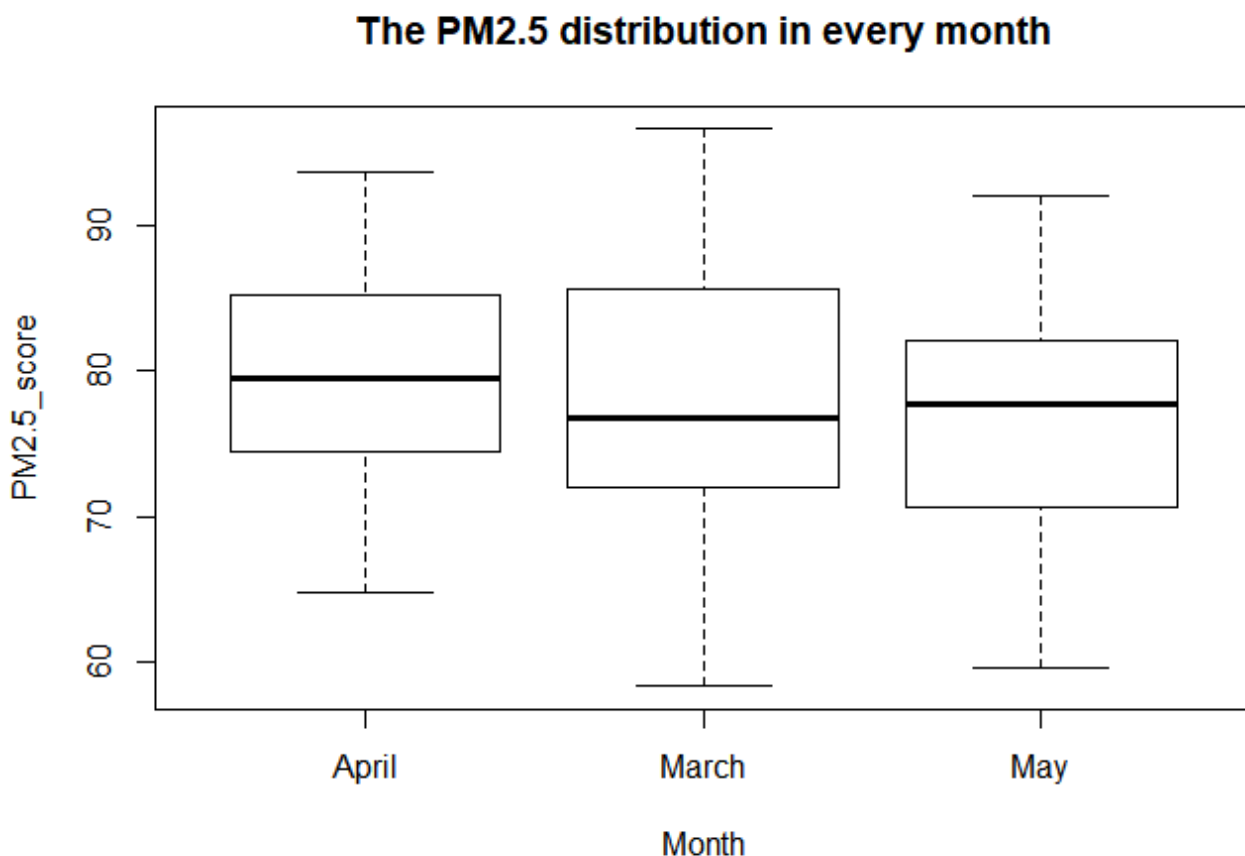
```
1 | ## [1] "C:/Users/liu/Documents/biostats_homework_summary"
```

2)use R to read in the data (3')

```
1 | # File is in your work directory, '~' means your work directory.  
2 | pm25 <- read.delim("./data/homework3_data.txt")
```

3)use boxplot to show the PM2.5 distribution in every month (10')

```
1 | boxplot(PM2.5_score~Month,pm25,main = "The PM2.5 distribution in every month")
```



4) get the data of month equal to March and store in data_march (10')

```
1 data_march <- subset(pm25, Month=="March")
2 head(data_march)
```

```
1 ## PM2.5_score Month
2 ## 1      81.6 March
3 ## 2      86.6 March
4 ## 3      80.8 March
5 ## 4      85.8 March
6 ## 5      78.6 March
7 ## 6      58.3 March
```

question 4

Suppose we draw a sample of size 20 of birthweights from a hospital, the details can be found in the homework data. The mean of national-wide birthweights is 118. (30')

```
1 library(readr)
2 birthweights <- read_csv("./data/homework4_data.csv")
```

```
1 ## Parsed with column specification:
2 ## cols(
3 ##   Individual_ID = col_double(),
4 ##   Birthweight = col_double()
5 ## )
```

```
1 head(birthweights)
```

```
1 ## # A tibble: 6 x 2
2 ##   Individual_ID Birthweight
3 ##         <dbl>         <dbl>
4 ## 1           1          123
5 ## 2           2           98
6 ## 3           3          115
7 ## 4           4          120
8 ## 5           5          105
9 ## 6           6          135
```

```
1 summary(birthweights$Birthweight)
```

```
1 ##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
2 ##   98.0   111.2   119.5   120.0   128.8   141.0
```

1) What is the probability that the mean birthweight of the sample falls between 100.0 and 126.0? Please list the formulas for this and also the R code for it. (5')

because σ^2 is unknown and n is less than 20, so we should use t -distribution to calculate probability that we need.

$$S_{\bar{X}} = \frac{S}{\sqrt{n}} \quad (1)$$

$$P(100.0 \leq \mu \leq 126.0) = pt_{n-1} \left(\frac{126 - \bar{X}}{S_{\bar{X}}} \right) - pt_{n-1} \left(\frac{100 - \bar{X}}{S_{\bar{X}}} \right) \quad (2)$$

```
1 | n <- 20
2 | mu <- 118
3 | df <- n-1
4 | mean <- mean(birthweights$Birthweight)
5 | sd <- sd(birthweights$Birthweight)
6 | Sx <- sd/sqrt(n)
7 | p_126 <- pt((126-mu)/Sx,df)
8 | p_100 <- pt((100-mu)/Sx,df)
9 | p <- p_126-p_100
10 | p
```

```
1 | ## [1] 0.9961359
```

So the probability that the mean birthweight of the sample falls between 100.0 and 126.0 is 0.9961359.

2)What is the 95% confidence interval of the sample mean?(5')

$$P \left(\bar{X} - t_{df,1-\alpha/2} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{df,1-\alpha/2} \frac{s}{\sqrt{n}} \right) = 1 - \alpha \quad (3)$$

```
1 | a <- 0.05
2 | P <- 1-a/2
3 | L1 <- mean - qt(P,n-1)*Sx
4 | L2 <- mean + qt(P,n-1)*Sx
5 | L1;L2
```

```
1 | ## [1] 114.3777
```

```
1 | ## [1] 125.6223
```

So the 95% confidence interval of the sample mean is [114.378,125.622].

3)Can we say the underlying mean birthweight from this hospital is higher than the national average?

Please list the formulas for this and also the R code for it.(5')

(1) .假设 $H_0 : \mu_1 \leq \mu_2$, 这家医院的新生儿体重小于或等于全国平均水平; $H_A : \mu_1 > \mu_2$ 。

(2) .确定显著水平 $\alpha=0.05$ 。

(3) .检验计算:

$$S_{\bar{X}} = \frac{S}{\sqrt{n}} \quad (4)$$

$$t = \frac{\bar{x} - \mu_0}{S_{\bar{X}}} \quad (5)$$

$$p = 1 - pt_{n-1}(t) \quad (6)$$

```
1 | t.test(birthweights$Birthweight, alternative = "greater", mu = 118)
```

```
1 ##
2 ## One Sample t-test
3 ##
4 ## data: birthweights$Birthweight
5 ## t = 0.74454, df = 19, p-value = 0.2328
6 ## alternative hypothesis: true mean is greater than 118
7 ## 95 percent confidence interval:
8 ## 115.3552      Inf
9 ## sample estimates:
10 ## mean of x
11 ##      120
```

$p > 0.05$, so maybe we can not say the underlying mean birthweight from this hospital is higher than the national average.

4) Test the hypothesis that the mean birthweight of sample size 20 is different from the national average (Significance level 0.05). Please list the formulas for this and also the R code for it. (5')

(1) .假设 $H_0: \mu_1 = \mu_2$, 即这家医院新生儿体重和全国新生儿体重平均值之间并无显著差异; $H_A: \mu_1 \neq \mu_2$ 。

(2) .确定显著水平 $\alpha = 0.05$ 。

(3) .检验计算:

$$S_{\bar{X}} = \frac{S}{\sqrt{n}} \quad (7)$$

$$t = \frac{\bar{x} - \mu_0}{S_{\bar{X}}} \quad (8)$$

$$p = 2 \times \min \left(pt_{n-1} \left(\frac{\bar{x} - \mu_0}{s/\sqrt{n}} \right), 1 - pt_{n-1} \left(\frac{\bar{x} - \mu_0}{s/\sqrt{n}} \right) \right) \quad (9)$$

```
1 | t.test(birthweights$Birthweight, alternative = "two.sided", mu = 118)
```

```

1 ##
2 ## One Sample t-test
3 ##
4 ## data: birthweights$Birthweight
5 ## t = 0.74454, df = 19, p-value = 0.4657
6 ## alternative hypothesis: true mean is not equal to 118
7 ## 95 percent confidence interval:
8 ## 114.3777 125.6223
9 ## sample estimates:
10 ## mean of x
11 ## 120

```

p-value = 0.4657 > 0.05, 不拒绝原假设, 认为这家医院新生儿体重和全国新生儿体重平均值之间并无显著差异。

5) Compute the power of the test performed in (4) with significance level 0.05.(5')

$$Power = F \left[t_{\alpha/2} + \frac{|\bar{X} - \mu_0|}{s/\sqrt{n}} \right] \quad (10)$$

```

1 power <- pt(qt(a/2,df)+abs(mean-mu)/Sx,n-1)
2 power

```

```

1 ## [1] 0.09667934

```

So, the power of the test performed in (4) with significance level is 0.09667934.

6) To see the significance difference between the sample mean and the national mean and ensure the type II error to be $\beta = 0.05$, what is the appropriate sample size with significance level is 0.01?(5')

$$n = \frac{(z_{1-\beta} + z_{1-\alpha/2})^2 S^2}{(\bar{x} - \mu_0)^2} \quad (11)$$

```

1 # 使用pwr包进行估计
2 # install.packages("pwr")
3 library("pwr")
4 d = abs(mean-mu)/sd
5 pwr.t.test(d=d, power=0.95, sig.level=0.01, type="one.sample", alternative="two.sided")

```

```

1 ##
2 ## One-sample t test power calculation
3 ##
4 ## n = 646.0381
5 ## d = 0.1664842
6 ## sig.level = 0.01
7 ## power = 0.95
8 ## alternative = two.sided

```

So, the appropriate sample size with significance level=0.01 is 647.

homework-4

question 1

一、为研究某种新药对抗凝血酶活力的影响，随机安排新药组病人12例，对照组病人10例，分别测定器抗凝血酶活力（单位：），其结果如下：

```
1  新药组: 126 125 138 128 123 138 142 116 110 108 113 140
2
3  对照组: 160 175 177 170 175 153 168 159 160 162
```

试分析新药组和对照组病人的抗凝血酶活力有无差别 ($\alpha = 0.05$)

解：

(1) 检验两组样本方差是否相同 (15')

假设检验：

H_0 : 两组样本方差不存在差异。

H_1 : 两组样本方差不存在差异。

先做正态检验，再做方差检验。

```
1  new_drug <- c(126,125,138,128,123,138,142,116,110,108,113,140)
2  control <- c(160,175,177,170,175,153,168,159,160,162)
3  # 原假设：样本符合正态分布；备择假设：样本不符合正态分布。
4  shapiro.test(new_drug)
```

```
1  ##
2  ##  shapiro-wilk normality test
3  ##
4  ## data:  new_drug
5  ## w = 0.92304, p-value = 0.3121
```

```
1  shapiro.test(control)
```

```
1  ##
2  ##  shapiro-wilk normality test
3  ##
4  ## data:  control
5  ## w = 0.92374, p-value = 0.3892
```

```
1  var.test(new_drug,control)
```

```

1  ##
2  ##  F test to compare two variances
3  ##
4  ## data:  new_drug and control
5  ## F = 2.1512, num df = 11, denom df = 9, p-value = 0.26
6  ## alternative hypothesis: true ratio of variances is not equal to 1
7  ## 95 percent confidence interval:
8  ##  0.5498769 7.7181417
9  ## sample estimates:
10 ## ratio of variances
11 ##          2.151159

```

结果显示两样本符合正态分布, $p=0.26>0.05$, 可认为两组样本方差一致 ($\alpha = 0.05$)

(2) 选择最合适的检验方法检验新药组和对照组病人的抗凝血酶活力有无差别。(15')

假设检验:

H_0 : 两组样本抗凝血酶活力无差别。

H_1 : 两组样本抗凝血酶活力有差别。

```

1 | t.test(new_drug, control, var.equal = TRUE, alternative = "two.sided", conf.level = 0.95)

```

```

1  ##
2  ##  Two Sample t-test
3  ##
4  ## data:  new_drug and control
5  ## t = -8.9578, df = 20, p-value = 1.947e-08
6  ## alternative hypothesis: true difference in means is not equal to 0
7  ## 95 percent confidence interval:
8  ##  -49.70504 -30.92829
9  ## sample estimates:
10 ## mean of x mean of y
11 ##  125.5833  165.9000

```

$p < 0.05$, 可认为新药组和对照组病人的抗凝血酶活力有显著差别。 ($\alpha = 0.05$)

question 2

二、对7位健康成年人的血液测量其中的尿酸浓度, 分别用手工 (X) 和仪器 (Y) 两种方法测量, 结果如下表所示, 请用wilcoxon signed-rank test来检测两种测量方法的精度是否存在差异? ($\alpha = 0.05$) (20')

手工(X)	4.5	6.5	6	9.2	10	12	8.3
仪器(Y)	4	7.2	8	14	8.8	10	11.5

解: 假设检验:

H_0 : 两种测量方法的精度不存在差异。

H_1 : 两种测量方法的精度存在差异。

```

1 handmade <- c(4.5,6.5,6,9.2,10,12,8.3)
2 device <- c(4,7.2,8,14,8.8,10,11.5)
3 wilcox.test(handmade,device,alternative = "two.sided",paired = TRUE,exact=FALSE)

```

```

1 ##
2 ## wilcoxon signed rank test with continuity correction
3 ##
4 ## data: handmade and device
5 ## v = 8.5, p-value = 0.3972
6 ## alternative hypothesis: true location shift is not equal to 0

```

因为 $p\text{-value} = 0.3972 > 0.05$ ，所以不拒绝原假设，认为两种测量方法的精度不存在差异。

question 3

三、在某保险种类中，一次关于2018年的索赔数额（单位：元）的随机抽样为（按升幂排列）：

```

1 4152, 4579, 5053, 5112, 5745, 6250, 7081, 9048,
2
3 12095, 14430, 17220, 20610, 22836, 48950, 67200

```

已知2017年的索赔数额的中位数为7520元。问2018年索赔的中位数与前一年是否有所变化？（ $\alpha = 0.05$ ）(15')

Hint: You can use wilcox.test

解：假设检验问题： H_0 :2018年索赔的中位数与前一年无变化。

H_0 :2018年索赔的中位数与前一年有变化。

```

1 insurance <-
  c(4152,4579,5053,5112,5745,6250,7081,9048,12095,14430,17220,20610,22836,48950,67200)
2
3 wilcox.test(insurance,mu=7520)

```

```

1 ##
2 ## wilcoxon signed rank test
3 ##
4 ## data: insurance
5 ## v = 87, p-value = 0.1354
6 ## alternative hypothesis: true location is not equal to 7520

```

$p\text{-value} = 0.1354 > 0.05$,所以不拒绝原假设，认为2018年索赔的中位数与前一年无变化。（ $\alpha = 0.05$ ）

question 4

Type 1 diabetes is a multigenic disease caused by T-cell mediated destruction of the insulin producing β -cells. Although conventional (targeted) approaches of identifying causative genes have advanced our knowledge of this disease, many questions remain unanswered.

Here we have a gene data from NOD mouse after(case) and before(control) treatment. The data can be found in "Data.txt". Use the information mentioned above to answer the following questions:

a) use paired t-test to find genes which have significant expression ($p < 0.05$) between case and control sample. Give the number of differential expressed genes and give the names of top 10 significantly differential expression genes. hint: "apply(data,1,function(x){...})" can apply function to every row in data more quickly than "for{}", "names()" or "rownames()" can be used to extract names of differentially expressed genes. (20')

```
1 # 导入数据
2 expression_data <- read.table('./data/Data.txt',header = T,stringsAsFactors =F )
3 # View(expression_data) ##观察数据可发现前十列为control组, 后十列为case组
4
5
6 # 定义函数根据方差检验结果做成对t.test
7 t.test.p.value <- function(x){
8   p_value <- t.test(x[1:10],x[11:20],paired = TRUE)$p.value
9   p_value
10 }
11
12 p.t.test <- apply(expression_data,1,t.test.p.value)
13
14 # Give the number of differential expressed genes
15 sum(p.t.test<0.05)
```

```
1 ## [1] 2296
```

```
1 # give the names of top 10 significantly differential expression genes
2 names(p.t.test[order(p.t.test,decreasing=F)[1:10])])
```

```
1 ## [1] "21744" "4008" "4817" "28474" "6816" "2593" "17191" "6786"
2 ## [9] "17677" "8678"
```

b) Adjust the p-values in question a) with bonferroni and FDR method to find differentially expressed genes in stringent way(list the differentially expressed gene names and the adjusted p-value). (15')

Hint: you can do the adjustment according to the fomular, or use "p.adjust()" instead.

```
1 p.bonf <- p.adjust(p.t.test,'bonferroni')
2 p.fdr <- p.adjust(p.t.test,'fdr')
3 p.bonf[p.bonf<0.05]
```

```
1 ## named numeric(0)
```

```
1 p.fdr[p.fdr<0.05]
```

```
1 ## named numeric(0)
```

经过校正后发现,并无显著差异表达基因。

homework-5

question-1

一、在一个农业实验中，育种人员测试了3种不同的种子的粮食产量（单位：共计/亩），结果记录在数据yield.txt中。（20'）请问

解：

(1) 种子的品种是否影响粮食产量；（10'）

用方差分析方法检验种子的品种是否影响粮食产量

H_0 ：种子的品种不影响粮食产量；

H_1 ：至少一个种子的品种影响粮食产量。

```
1 yield <- read.delim("./data/yield.txt")
2
3 # H0:数据符合正态分布 H1: 数据不符合正态分布。
4 shapiro.test(yield$yield[yield$seed==1])
```

```
1 ##
2 ##  Shapiro-wilk normality test
3 ##
4 ## data:  yield$yield[yield$seed == 1]
5 ## W = 0.93955, p-value = 0.548
```

```
1 shapiro.test(yield$yield[yield$seed==2])
```

```
1 ##
2 ##  Shapiro-wilk normality test
3 ##
4 ## data:  yield$yield[yield$seed == 2]
5 ## W = 0.93528, p-value = 0.4667
```

```
1 shapiro.test(yield$yield[yield$seed==3])
```

```
1 ##
2 ##  Shapiro-wilk normality test
3 ##
4 ## data:  yield$yield[yield$seed == 3]
5 ## W = 0.97869, p-value = 0.9561
```

```
1 # 满足正态分布。
2 # H0: 数据方差齐性 H1: 数据方差非齐性
3 bartlett.test(yield~seed,data = yield)
```



```

1 ##
2 ## Bartlett test of homogeneity of variances
3 ##
4 ## data: yield by seed
5 ## Bartlett's K-squared = 4.0299, df = 2, p-value = 0.1333

```

```

1 # 满足方差齐次性。
2
3 seed_aov <- aov(yield~factor(seed),data = yield)
4 summary(seed_aov)

```

```

1 ##              Df Sum Sq Mean Sq F value    Pr(>F)
2 ## factor(seed)  2   4364   2182.2    10.04 0.000586 ***
3 ## Residuals    26   5649    217.3
4 ## ---
5 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

答：检验结果显示p值小于0.05，因此拒绝原假设，种子的品种会影响粮食产量。

(2) 如果受品种影响，那么哪一个品种和其他品种的产量有差异。(10')

```

1 | tukey_results <- TukeyHSD(seed_aov);tukey_results

```

```

1 ##      Tukey multiple comparisons of means
2 ##      95% family-wise confidence level
3 ##
4 ## Fit: aov(formula = yield ~ factor(seed), data = yield)
5 ##
6 ## $`factor(seed)`
7 ##           diff          lwr          upr      p adj
8 ## 2-1 -25.32727 -41.330375 -9.324171 0.0015630
9 ## 3-1  -0.10000 -17.473293 17.273293 0.9998872
10 ## 3-2  25.22727   8.208575 42.245971 0.0029482

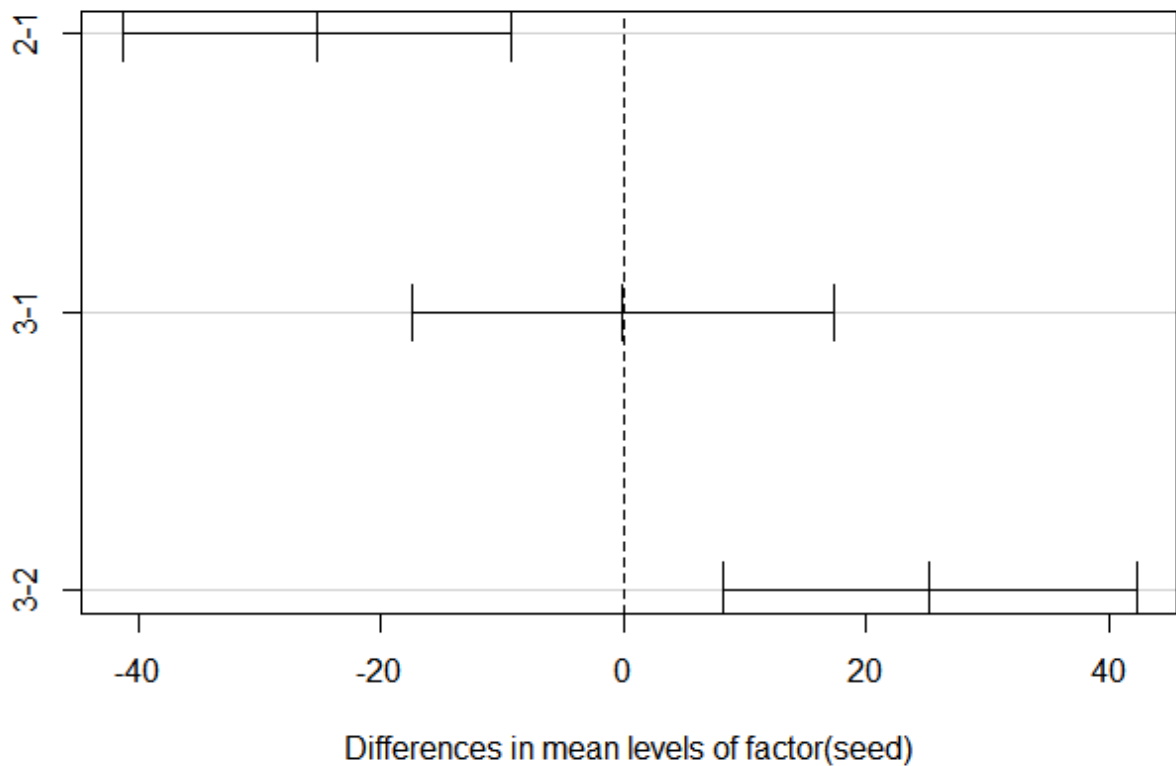
```

```

1 | plot(tukey_results)

```

95% family-wise confidence level



```
1 # 置信区间包含0说明差异并不显著。
```

答：品种二和其他品种的产量有差异。

或者用pairwise-t检验分析哪个品种的产量和其他品种不同。

```
1 pairwise.t.test(yield$yield,yield$seed,p.adjust.method = "none")
```

```
1 ##
2 ## Pairwise comparisons using t tests with pooled SD
3 ##
4 ## data: yield$yield and yield$seed
5 ##
6 ##      1      2
7 ## 2 0.00056 -
8 ## 3 0.98870 0.00106
9 ##
10 ## P value adjustment method: none
```

结果显示种子1和种子3的t检验的p值大于0.05，而种子2和其他两者的t检验p值均小于0.05，因此种子2和其他种子相比产量有差异。

question-2

二、为研究茶多酚保健饮料对急性缺氧的影响，某研究者将60只小白鼠随机分为低、中、高三个剂量组和一个对照组，每组15只老鼠。对照组给予蒸馏水0.25ml灌胃，低中高分别给予递增剂量的饮料，并将饮料溶于0.2~0.3ml蒸馏水后灌胃。每天一次，40天后，对小鼠进行耐缺氧存活时间试验，结果见数据文件。（30'）

解：

1.在本次试验中，为研究不同剂量的茶多酚保健饮料对延长小白鼠的平均耐缺氧存活时间有无差异。问因素或者处理是什么？，与之相对的有多少个分组或者水平？（5'）

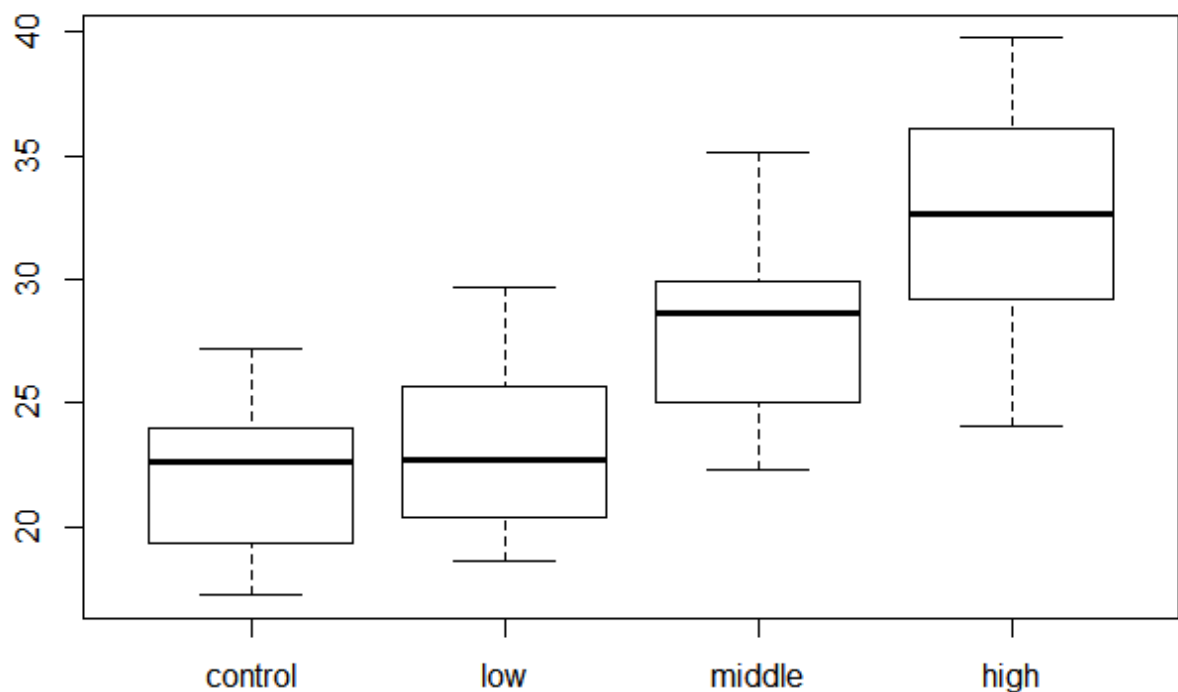
答：

因素：不同剂量的茶多酚保健饮料

四个处理：对照、低、中、高组

2.请将数据读入R中，并用Boxplot查看总体数据的情况，最后查看一下数据的最小值、中位数、平均数等信息。（5'）

```
1 data <- read.delim("../data/data2.txt")
2 boxplot(data)
```



```
1 summary(data)
```

	##	control	low	middle	high
2	##	Min. :17.22	Min. :18.64	Min. :22.29	Min. :24.07
3	##	1st Qu.:19.35	1st Qu.:20.39	1st Qu.:25.05	1st Qu.:29.21
4	##	Median :22.60	Median :22.69	Median :28.67	Median :32.63
5	##	Mean :21.98	Mean :23.23	Mean :28.13	Mean :32.84
6	##	3rd Qu.:23.96	3rd Qu.:25.69	3rd Qu.:29.95	3rd Qu.:36.14
7	##	Max. :27.21	Max. :29.67	Max. :35.12	Max. :39.76

3.有研究员对上述资料采用了两样本均数t检验进行了两两比较。问这样处理是否合理，为什么？应采用何种处理方法。（5'）

答：不合理，t检验适合两组数据的检验，用于多组数据会增大犯一型错误的概率。对于多组数据的差异研究，应该采用方差分析。

4.为了用更好的方法来处理数据，请说明数据应满足的哪三个基本条件？试检验这批数据是否满足这些条件。（提示：一般需要满足三个条件；shapiro.test()函数可检验正态性，bartlett.test()可检验多个正态总体的方差齐次性）（5'）

答：

方差分析的基本假设：

(1) 各总体的方差必须相等。

```
1 # H0: 数据方差齐性 H1: 数据方差非齐性
2 bartlett.test(data)
```

```
1 ##
2 ## Bartlett test of homogeneity of variances
3 ##
4 ## data: data
5 ## Bartlett's K-squared = 2.1206, df = 3, p-value = 0.5478
```

满足方差齐性。

(2) 各总体必须服从正态分布。

```
1 # H0: 数据符合正态分布 H1: 数据不符合正态分布。
2 shapiro.test(data$control)
```

```
1 ##
2 ## Shapiro-Wilk normality test
3 ##
4 ## data: data$control
5 ## W = 0.95001, p-value = 0.5245
```

```
1 shapiro.test(data$low)
```

```

1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: data$low
5 ## W = 0.93877, p-value = 0.3671

```

```

1 shapiro.test(data$middle)

```

```

1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: data$middle
5 ## W = 0.94679, p-value = 0.4754

```

```

1 shapiro.test(data$high)

```

```

1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: data$high
5 ## W = 0.95715, p-value = 0.6431

```

经检验，各总体服从正态分布。

(3) 各观测值相互独立

可以通过控制抽样过程来控制独立性，无具体的检验方法。

5.请根据上述验证情况，对数据进行分析，并给出差异的配对组别（10'）

H_0 四个总体均值相等。 H_A 至少有一个不等， $\alpha=0.05$

```

1 time <- c(data$control,data$low,data$middle,data$high)
2
3 levels <- factor(c(rep("control",15),rep("low",15),rep("middle",15),rep("high",15)))
4 time_aov <- aov(lm(time~levels))
5 summary(time_aov)

```

```

1 ##              Df Sum Sq Mean Sq F value    Pr(>F)
2 ## levels        3 1109.1    369.7    24.46 3.03e-10 ***
3 ## Residuals    56  846.5      15.1
4 ## ---
5 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

1 TukeyHSD(time_aov)

```

```

1 ## Tukey multiple comparisons of means
2 ## 95% family-wise confidence level
3 ##
4 ## Fit: aov(formula = lm(time ~ levels))
5 ##
6 ## $levels
7 ##          diff          lwr          upr      p adj
8 ## high-control 10.860000    7.100912 14.6190884 0.0000000
9 ## low-control   1.254667   -2.504422  5.0137551 0.8132480
10 ## middle-control 6.151333    2.392245  9.9104218 0.0003513
11 ## low-high     -9.605333  -13.364422 -5.8462449 0.0000001
12 ## middle-high  -4.708667   -8.467755 -0.9495782 0.0084746
13 ## middle-low    4.896667    1.137578  8.6557551 0.0057640

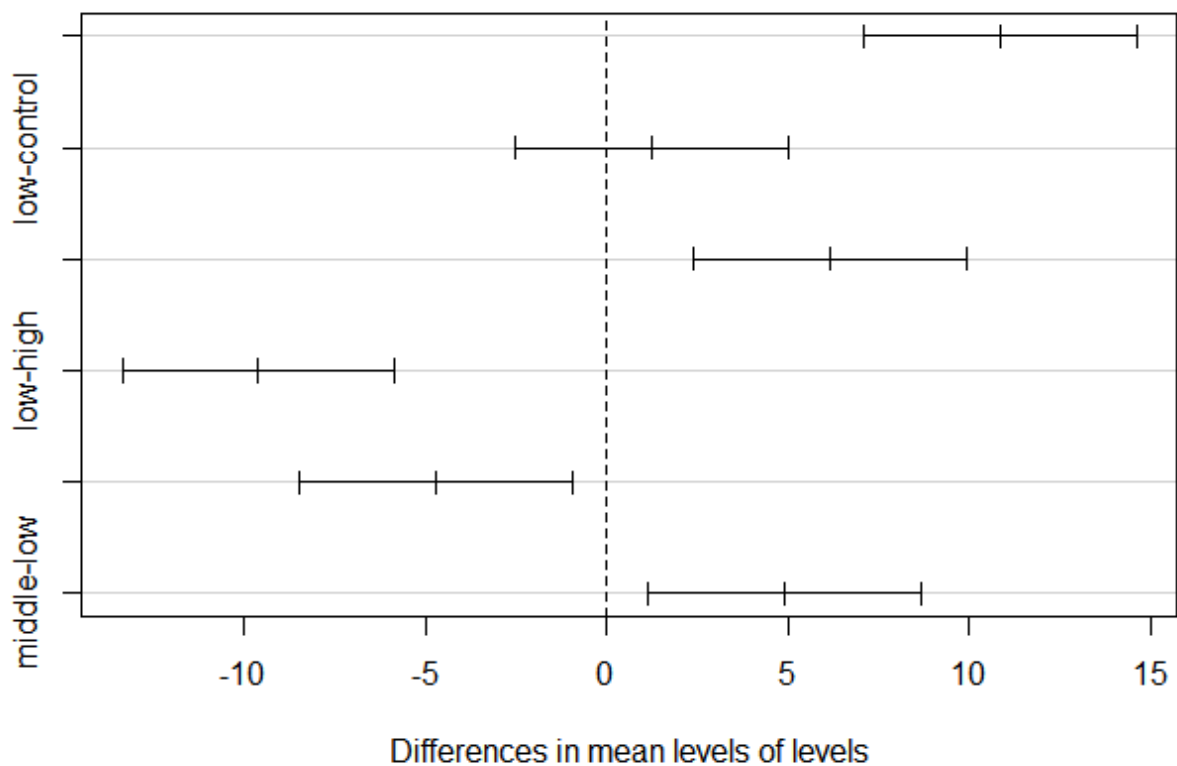
```

```

1 plot(TukeyHSD(time_aov))

```

95% family-wise confidence level



答:除low组和control组经检验不具有显著性差异，其他组都具有显著差异。

question-3

三、已知有三种药物都能促进小鼠肠道对营养的吸收，现将初始状态相近的两批成年小鼠分别进行给药，经过一段时间后测量其体重，然后得到这段时间内小鼠体重增加的值。问：这三种药导致的平均体重增加值有无统计学差异？

(20')

药物1 (增加的体重值g)	药物2 (增加的体重值g)	药物3 (增加的体重值g)
40	50	60
10	20	30
35	45	100
25	55	85
20	20	20
15	15	55
35	80	45
15	-10	30
-5	105	77
30	75	105
25	10	
70	60	
65	45	
45	60	
50	30	

解：求均值的统计学差异，用ANOVA来解决。

H0：三种药物对小鼠体重均值的影响都相同；

H1：三种药物对小鼠体重均值的影响至少有一个不同。

```

1 # 输入数据
2 weight1 <- c(40,10,35,25,20,15,35,15,-5,30,25,70,65,45,50)
3 weight2 <- c(50,20,45,55,20,15,80,-10,105,75,10,60,45,60,30)
4 weight3 <- c(60,30,100,85,20,55,45,30,77,105)
5 # 正态性检验
6 shapiro.test(weight1)

```

```

1 ##
2 ##  Shapiro-wilk normality test
3 ##
4 ## data:  weight1
5 ## W = 0.9778, p-value = 0.9523

```

```

1 shapiro.test(weight2)

```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: weight2
5 ## W = 0.98349, p-value = 0.9878
```

```
1 shapiro.test(weight3)
```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: weight3
5 ## W = 0.94005, p-value = 0.5536
```

```
1 weights <- c(weight1,weight2,weight3)
2 drugs <- factor(c(rep("drug1",15),rep("drug2",15),rep("drug3",10)))
3
4 # 方差齐性检验
5 bartlett.test(weights,g = drugs)
```

```
1 ##
2 ## Bartlett test of homogeneity of variances
3 ##
4 ## data: weights and drugs
5 ## Bartlett's K-squared = 2.4482, df = 2, p-value = 0.294
```

```
1 #方差分析
2 weight_aov <- aov(lm(weights~drugs))
3 summary(weight_aov)
```

```
1 ##              Df Sum Sq Mean Sq F value Pr(>F)
2 ## drugs          2    5062   2531.2    3.485 0.0411 *
3 ## Residuals     37   26877    726.4
4 ## ---
5 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
1 TukeyHSD(weight_aov)
```

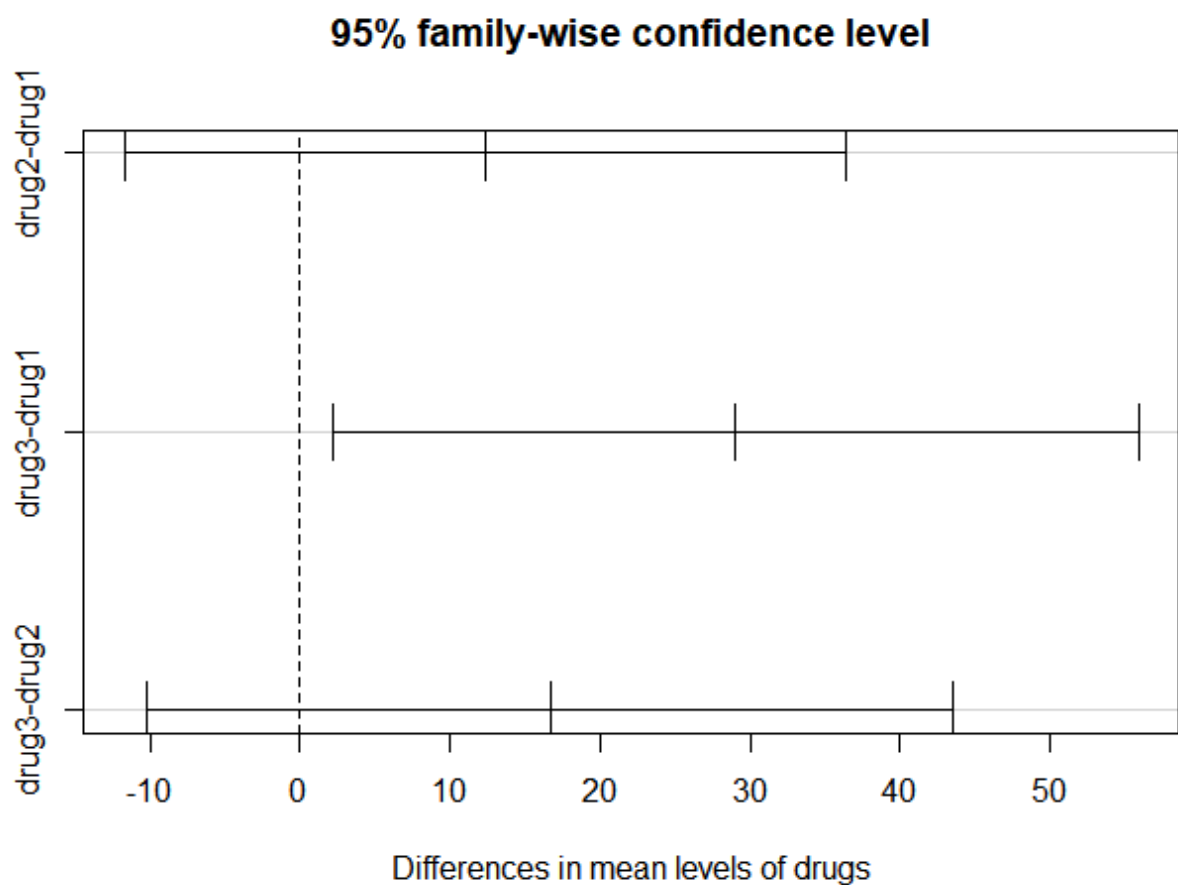


```

1 ## Tukey multiple comparisons of means
2 ## 95% family-wise confidence level
3 ##
4 ## Fit: aov(formula = lm(weights ~ drugs))
5 ##
6 ## $drugs
7 ##          diff          lwr          upr          p adj
8 ## drug2-drug1 12.33333 -11.694604 36.36127 0.4302024
9 ## drug3-drug1 29.03333   2.169283 55.89738 0.0317315
10 ## drug3-drug2 16.70000 -10.164050 43.56405 0.2944100

```

```
1 plot(TukeyHSD(weight_aov))
```



答:经检验，药物3与药物1导致的平均体重增加值之间存在显著差异，其他药物导致的平均体重增加值之间不存在显著差异。

question-4

四、在一生物实验中，为了研究不同饲养条件对大鼠体重的影响，现将60只8周龄体重相等的大鼠随机分为六组，分别放入以下饲养条件中培养两周：

A: 饲养温度（℃） 4（A1）， 25（A2）， 30（A3）

B: 饲料：普通饲料（B1），高脂饲料（B2）

	A1	A2	A3
B1	282.1	296.7	300.1
	264.2	318	307.5
	274.2	295.3	294.2
	276.4	292.8	312
	283.7	304.5	300.2
	288	305.9	292.6
	274.3	312.3	302
	278.4	311.4	306.9
	293.5	307.6	313.3
	271.5	292.7	312.4
B2	284.5	296.6	304
	263.7	323.9	312.4
	292.3	296.5	297.8
	270.3	298.1	318.4
	281.3	310	302.4
	286.9	312.5	295.3
	271.2	317.6	305.5
	275.6	305.5	309.9
	289.4	305.8	319.2
	289.8	295.4	320.2

解：

(1) 检验体重数据对于因素A和因素B是否是正态的？是否满足方差齐性的要求？（10分）

将数据整理成csv并导入。

```
1 library(readr)
2 rat_weight <- read_csv("./data/biostats_homework_weight.csv")
```

```
1 ## Parsed with column specification:
2 ## cols(
3 ##   weight = col_double(),
4 ##   A = col_character(),
5 ##   B = col_character()
6 ## )
```

```
1 shapiro.test(rat_weight$weight[rat_weight$A=="A1"])
```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data:  rat_weight$weight[rat_weight$A == "A1"]
5 ## W = 0.96118, p-value = 0.5676
```

```
1 shapiro.test(rat_weight$weight[rat_weight$A=="A2"])
```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data:  rat_weight$weight[rat_weight$A == "A2"]
5 ## W = 0.93573, p-value = 0.1989
```

```
1 shapiro.test(rat_weight$weight[rat_weight$A=="A3"])
```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data:  rat_weight$weight[rat_weight$A == "A3"]
5 ## W = 0.96408, p-value = 0.6281
```

```
1 shapiro.test(rat_weight$weight[rat_weight$B=="B1"])
```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data:  rat_weight$weight[rat_weight$B == "B1"]
5 ## W = 0.95026, p-value = 0.1718
```

```
1 shapiro.test(rat_weight$weight[rat_weight$B=="B2"])
```

```

1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: rat_weight$weight[rat_weight$B == "B2"]
5 ## W = 0.96797, p-value = 0.4853

```

```

1 bartlett.test(weight~A,data = rat_weight)

```

```

1 ##
2 ## Bartlett test of homogeneity of variances
3 ##
4 ## data: weight by A
5 ## Bartlett's K-squared = 0.19888, df = 2, p-value = 0.9053

```

```

1 bartlett.test(weight~B,data = rat_weight)

```

```

1 ##
2 ## Bartlett test of homogeneity of variances
3 ##
4 ## data: weight by B
5 ## Bartlett's K-squared = 0.20953, df = 1, p-value = 0.6471

```

答：经检验，体重数据对于因素A和因素B是正态的，并满足方差齐性的要求。

(2) 试分析因素A、因素B以及两因素的相互作用对大鼠体重有无显著影响？（10分）

```

1 rat_fit<-aov(weight~A+B+A:B,data = rat_weight)
2 summary(rat_fit)

```

```

1 ##           Df Sum Sq Mean Sq F value    Pr(>F)
2 ## A           2   9080    4540  56.809 5.22e-14 ***
3 ## B           1    127     127   1.589   0.213
4 ## A:B         2     17      9    0.108   0.897
5 ## Residuals  54   4316      80
6 ## ---
7 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

答:因素A对大鼠体重有显著影响

(3) 对 $A_i B_j$ 条件下平均产量作多重比较，并指出与 $A_2 B_1$ 组有显著差异的组。（10分）

```

1 TukeyHSD(rat_fit)

```

```

1 ## Tukey multiple comparisons of means
2 ## 95% family-wise confidence level
3 ##
4 ## Fit: aov(formula = weight ~ A + B + A:B, data = rat_weight)
5 ##

```

```

6  ## $A
7  ##          diff          lwr          upr          p adj
8  ## A2-A1 25.39 18.576905 32.203095 0.0000000
9  ## A3-A1 26.75 19.936905 33.563095 0.0000000
10 ## A3-A2  1.36 -5.453095  8.173095 0.8805321
11 ##
12 ## $B
13 ##          diff          lwr          upr          p adj
14 ## B2-B1 2.91 -1.717782 7.537782 0.2128406
15 ##
16 ## $`A:B`
17 ##          diff          lwr          upr          p adj
18 ## A2:B1-A1:B1 25.09 13.277922 36.90208 0.0000009
19 ## A3:B1-A1:B1 25.49 13.677922 37.30208 0.0000006
20 ## A1:B2-A1:B1  1.87 -9.942078 13.68208 0.9970585
21 ## A2:B2-A1:B1 27.56 15.747922 39.37208 0.0000001
22 ## A3:B2-A1:B1 29.88 18.067922 41.69208 0.0000000
23 ## A3:B1-A2:B1  0.40 -11.412078 12.21208 0.9999985
24 ## A1:B2-A2:B1 -23.22 -35.032078 -11.40792 0.0000050
25 ## A2:B2-A2:B1  2.47 -9.342078 14.28208 0.9892562
26 ## A3:B2-A2:B1  4.79 -7.022078 16.60208 0.8358408
27 ## A1:B2-A3:B1 -23.62 -35.432078 -11.80792 0.0000035
28 ## A2:B2-A3:B1  2.07 -9.742078 13.88208 0.9952518
29 ## A3:B2-A3:B1  4.39 -7.422078 16.20208 0.8800277
30 ## A2:B2-A1:B2 25.69 13.877922 37.50208 0.0000005
31 ## A3:B2-A1:B2 28.01 16.197922 39.82208 0.0000001
32 ## A3:B2-A2:B2  2.32 -9.492078 14.13208 0.9919360

```

答: A_1B_1 , A_1B_2 组和 A_2B_1 组有显著差异 (3分)

homework-6

question-1

一、某地29名13岁儿童身高 (cm) , 体重 (kg) 和肺活量 (L) 数据见data, 求: (1) 由身高, 体重推算肺活量的回归方程; (2) 求出的方程是否有意义; (3) 剩余标准差

```

1 library(readxl)
2 homework_6_1_data <- read_excel("./data/homework-6.1-data.xlsx")
3 y <- homework_6_1_data$y
4 x1 <- homework_6_1_data$x1
5 x2 <- homework_6_1_data$x2
6 ff <- lm(y~x1+x2)
7 ff

```

```

1 ##
2 ## Call:
3 ## lm(formula = y ~ x1 + x2)
4 ##
5 ## Coefficients:
6 ## (Intercept)          x1          x2
7 ##   -0.565664    0.005017    0.054061

```

```

1 summary(ff)

```

```

1 ##
2 ## Call:
3 ## lm(formula = y ~ x1 + x2)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -0.54117 -0.25524 -0.00266  0.22039  0.55425
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept) -0.565664   1.240127  -0.456  0.65208
12 ## x1           0.005017   0.010575   0.474  0.63920
13 ## x2           0.054061   0.015984   3.382  0.00228 **
14 ## ---
15 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16 ##
17 ## Residual standard error: 0.3137 on 26 degrees of freedom
18 ## Multiple R-squared:  0.546, Adjusted R-squared:  0.511
19 ## F-statistic: 15.63 on 2 and 26 DF, p-value: 3.485e-05

```

答：

(1)回归方程： $y = -0.565664 + 0.005017x_1 + 0.054061x_2$ 。

(2)因为 $F=15.63$ ， $p=3.485e-05 < 0.01$ ，所以方程有意义。(这个具体为什么用F检验，看书吧)

(3)剩余标准差(Residual standard error): 0.3137 on 26 degrees of freedom

question-2

二、某农场通过试验取得早稻收获量与春季降雨量和春季温度的数据如下：

收获量y(kg/mm2)	降雨量x1(mm)	温度x2(°C)
2250	25	6
3450	33	8
4500	45	10
6750	105	13
7200	110	14
7500	115	16
8250	120	17

建立早稻收获量对春季降雨量和春季温度的二元线性回归方程，计算各回归系数的置信区间，并对回归模型的线性关系和回归系数进行检验 ($\alpha = 0.05$)。

```
1 library(readr)
2 homework_6_2_data <- read_delim("./data/homework-6.2-data.txt",
3   "\t", escape_double = FALSE, trim_ws = TRUE)
```

```
1 ## Parsed with column specification:
2 ## cols(
3 ##   y = col_double(),
4 ##   x1 = col_double(),
5 ##   x2 = col_double()
6 ## )
```

```
1 y <- homework_6_2_data$y
2 x1 <- homework_6_2_data$x1
3 x2 <- homework_6_2_data$x2
4 ff <- lm(y~x1+x2)
5 ff
```

```
1 ##
2 ## Call:
3 ## lm(formula = y ~ x1 + x2)
4 ##
5 ## Coefficients:
6 ## (Intercept)          x1          x2
7 ##      -0.591       22.386      327.672
```

```
1 # 置信区间
2 confint(ff)
```

```

1 ##                2.5 %    97.5 %
2 ## (Intercept) -1402.707516 1401.52552
3 ## x1          -4.268921   49.04184
4 ## x2          53.364699   601.97873

```

降雨量x1的置信区间为 (-4.268921, 49.04184) , 含义是在温度不变的条件下, 降雨量每变动1mm, 收获量的平均变动在-4.268921到49.04184 kg/mm2之间。

温度x2的置信区间为 (53.364699, 601.97873) , 含义是在降雨量不变的条件下, 温度每变动1°C, 收获量的平均变动在53.364699到601.97873 kg/mm2之间。

```

1 summary(ff)

```

```

1 ##
2 ## Call:
3 ## lm(formula = y ~ x1 + x2)
4 ##
5 ## Residuals:
6 ##      1      2      3      4      5      6      7
7 ## -275.101  90.464 216.483 140.280 150.676 -316.599  -6.203
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept)   -0.591    505.004  -0.001   0.9991
12 ## x1             22.387     9.601   2.332   0.0801 .
13 ## x2            327.672    98.798   3.317   0.0295 *
14 ## ---
15 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16 ##
17 ## Residual standard error: 261.4 on 4 degrees of freedom
18 ## Multiple R-squared:  0.9913, Adjusted R-squared:  0.987
19 ## F-statistic: 228.4 on 2 and 4 DF,  p-value: 7.532e-05

```

线性关系检验是检验因变量y与k个自变量之间的关系是否显著, 也称总体显著性检验。具体步骤如下:

第1步 提出假设:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0 \quad (12)$$

$$H_1 : \beta_1, \beta_2, \dots, \beta_k \text{ 至少有一个不等于 } 0 \quad (13)$$

第2步 计算检验统计量F:

$$F = \frac{SSR/k}{SSE/(n-k-1)} \sim F(k, n-k-1) \quad (14)$$

第3步 做出决策。给定显著水平 α , 根据分子的自由度= k , 分母的自由度= $n-k-1$ 计算出统计量的P值。若 $P < \alpha$, 拒绝原假设, 表明y与k个自变量之间的线性关系显著。根据以上R输出结果, 检验统计量 $F = 228.4$, 显著水平 $P = 7.532e-05 < 0.05$, 拒绝 H_0 , 即收获量y与降雨量x1和温度x2之间的线性关系显著。

4、要判断每个自变量对因变量的影响是否都显著, 需要对各回归系数 β_i 分别进行t检验, 具体步骤如下: 第1步 提出假设。对于任意参数 $\beta_i (i = 1, 2, \dots, k)$, 有

$$H_0 : \beta = 0, H_1 : \beta \neq 0 \quad (15)$$

第2步 计算检验的统计量t:

$$t_i = \frac{\hat{\beta}_i}{s_{\hat{\beta}_i}} \sim t(n - k - 1) \quad (16)$$

其中： $s_{\hat{\beta}_i}$ 是回归系数的抽样分布的标准差。第3步 做出决策。给定显著性水平 α ，根据自由度= n-k-1计算出统计量的P值。若 $P < \alpha$ ，则拒绝原假设，表明回归系数 β_i 显著。根据R输出结果，降雨量x1和温度x2的回归系数相应的显著水平分别为0.0801和0.0295，只有温度对应的显著性水平小于0.05通过检验，这表明影响收获量的自变量中，只有温度对收获量的影响显著，而降雨量对收获量的影响不显著。

question-3

三、某葡萄酒爱好者想探索葡萄酒的品质与哪些因素相关。他有一个数据集包含了（1 -固定酸度，2 -挥发性酸度，3 -柠檬酸，4 -残余糖，5 -氯化物，6 -自由二氧化硫量，7 -二氧化硫总量，8 -密度，9 - pH值，10 -硫酸盐，11 -酒精浓度，和12 -品质(0 - 10分)。

1.查看数据集的前五行和数据集的总结。

```
1 library(readr)
2 winequality <- read_csv("../data/homework-6.3-winequality-red.csv")
```

```
1 ## Parsed with column specification:
2 ## cols(
3 ##   `fixed acidity` = col_double(),
4 ##   `volatile acidity` = col_double(),
5 ##   `citric acid` = col_double(),
6 ##   `residual sugar` = col_double(),
7 ##   chlorides = col_double(),
8 ##   `free sulfur dioxide` = col_double(),
9 ##   `total sulfur dioxide` = col_double(),
10 ##   density = col_double(),
11 ##   pH = col_double(),
12 ##   sulphates = col_double(),
13 ##   alcohol = col_double(),
14 ##   quality = col_double()
15 ## )
```

```
1 head(winequality,n=5)
```

```

1 ## # A tibble: 5 x 12
2 ##   `fixed acidity` `volatile acidity` `citric acid` `residual sugar` chlorides
3 ##           <dbl>           <dbl>           <dbl>           <dbl>           <dbl>
4 ## 1             7.4             0.7             0             1.9           0.076
5 ## 2             7.8             0.88            0             2.6           0.098
6 ## 3             7.8             0.76            0.04           2.3           0.092
7 ## 4            11.2             0.28            0.56           1.9           0.075
8 ## 5             7.4             0.7             0             1.9           0.076
9 ## # ... with 7 more variables: `free sulfur dioxide` <dbl>, `total sulfur
10 ## #   dioxide` <dbl>, density <dbl>, pH <dbl>, sulphates <dbl>,
11 ## #   alcohol <dbl>, quality <dbl>

```

```
1 summary(winequality)
```

```

1 ## fixed acidity volatile acidity citric acid residual sugar
2 ## Min. : 4.60 Min. :0.1200 Min. :0.000 Min. : 0.900
3 ## 1st Qu.: 7.10 1st Qu.:0.3900 1st Qu.:0.090 1st Qu.: 1.900
4 ## Median : 7.90 Median :0.5200 Median :0.260 Median : 2.200
5 ## Mean : 8.32 Mean :0.5278 Mean :0.271 Mean : 2.539
6 ## 3rd Qu.: 9.20 3rd Qu.:0.6400 3rd Qu.:0.420 3rd Qu.: 2.600
7 ## Max. :15.90 Max. :1.5800 Max. :1.000 Max. :15.500
8 ## chlorides free sulfur dioxide total sulfur dioxide
9 ## Min. :0.01200 Min. : 1.00 Min. : 6.00
10 ## 1st Qu.:0.07000 1st Qu.: 7.00 1st Qu.: 22.00
11 ## Median :0.07900 Median :14.00 Median : 38.00
12 ## Mean :0.08747 Mean :15.87 Mean : 46.47
13 ## 3rd Qu.:0.09000 3rd Qu.:21.00 3rd Qu.: 62.00
14 ## Max. :0.61100 Max. :72.00 Max. :289.00
15 ## density pH sulphates alcohol
16 ## Min. :0.9901 Min. :2.740 Min. :0.3300 Min. : 8.40
17 ## 1st Qu.:0.9956 1st Qu.:3.210 1st Qu.:0.5500 1st Qu.: 9.50
18 ## Median :0.9968 Median :3.310 Median :0.6200 Median :10.20
19 ## Mean :0.9967 Mean :3.311 Mean :0.6581 Mean :10.42
20 ## 3rd Qu.:0.9978 3rd Qu.:3.400 3rd Qu.:0.7300 3rd Qu.:11.10
21 ## Max. :1.0037 Max. :4.010 Max. :2.0000 Max. :14.90
22 ## quality
23 ## Min. :3.000
24 ## 1st Qu.:5.000
25 ## Median :6.000
26 ## Mean :5.636
27 ## 3rd Qu.:6.000
28 ## Max. :8.000

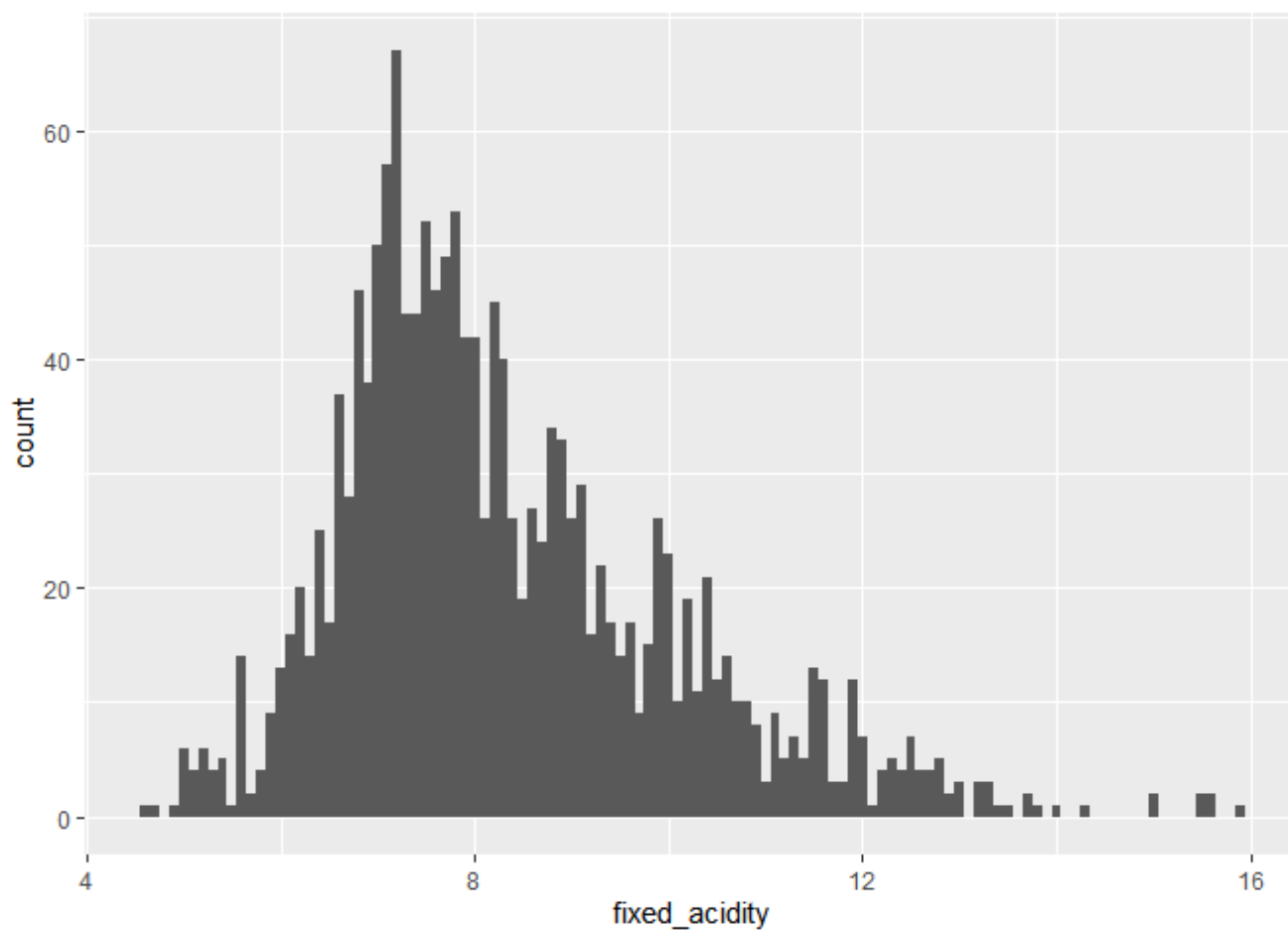
```

2.通过直方图展示固定酸度的分布和展示挥发性酸度与品质的散点图。

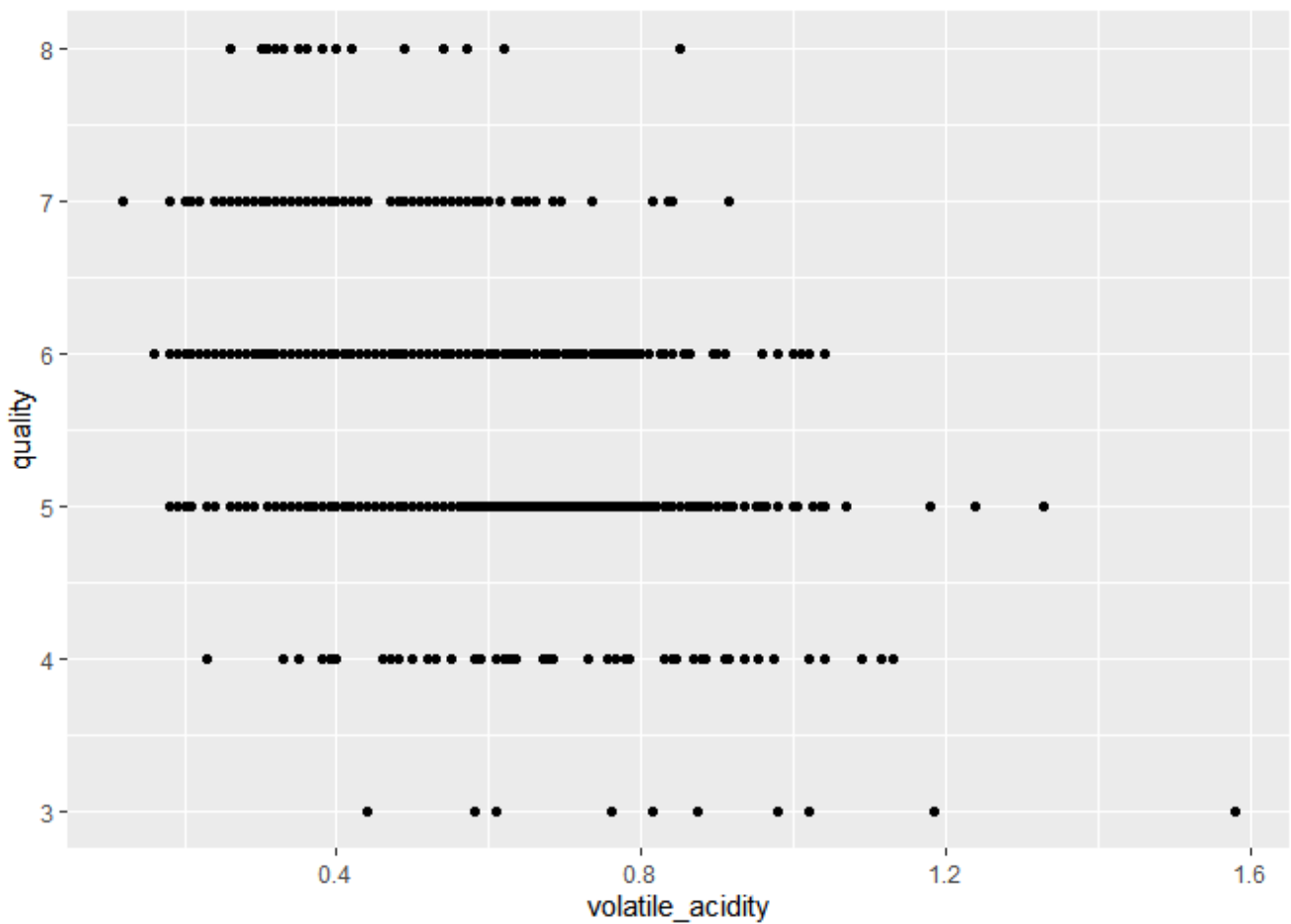
```

1 library(ggplot2)
2 fixed_acidity <- winequality$`fixed acidity`
3 hist=ggplot(winequality,aes(x=fixed_acidity))+geom_histogram(stat="bin",binwidth = 0.1)
4 hist

```



```
1 volatile_acidity <- winequality$`volatile acidity`  
2 quality <- winequality$quality  
3 point <- ggplot(winequality, mapping = aes(x=volatile_acidity, y=quality)) + geom_point()  
4 point
```



3.计算这些变量与品质的相关性。

```
1 | quality_cor <- cor(winequality[1:11],winequality$quality)
2 | quality_cor
```

```
1 | ##                                [,1]
2 | ## fixed acidity                   0.12405165
3 | ## volatile acidity                -0.39055778
4 | ## citric acid                     0.22637251
5 | ## residual sugar                  0.01373164
6 | ## chlorides                      -0.12890656
7 | ## free sulfur dioxide             -0.05065606
8 | ## total sulfur dioxide            -0.18510029
9 | ## density                        -0.17491923
10 | ## pH                             -0.05773139
11 | ## sulphates                      0.25139708
12 | ## alcohol                        0.47616632
```

4.通过方差分析不同品质的葡萄酒的酒精浓度是否有差异。

```
1 | alcohol <- winequality$alcohol
2 | alcohol_anova <- data.frame(alcohol,quality)
3 | # 方差齐性检验
4 | bartlett.test(alcohol_anova)
```

```
1 ##
2 ## Bartlett test of homogeneity of variances
3 ##
4 ## data: alcohol_anova
5 ## Bartlett's K-squared = 121.32, df = 1, p-value < 2.2e-16
```

```
1 # 正态检验
2 for (i in c(3:8)) {
3   print(shapiro.test(alcohol_anova$alcohol[alcohol_anova$quality==i]))
4 }
```

```
1 ##
2 ## Shapiro-wilk normality test
3 ##
4 ## data: alcohol_anova$alcohol[alcohol_anova$quality == i]
5 ## w = 0.9423, p-value = 0.5788
6 ##
7 ##
8 ## Shapiro-wilk normality test
9 ##
10 ## data: alcohol_anova$alcohol[alcohol_anova$quality == i]
11 ## w = 0.93444, p-value = 0.00607
12 ##
13 ##
14 ## Shapiro-wilk normality test
15 ##
16 ## data: alcohol_anova$alcohol[alcohol_anova$quality == i]
17 ## w = 0.84302, p-value < 2.2e-16
18 ##
19 ##
20 ## Shapiro-wilk normality test
21 ##
22 ## data: alcohol_anova$alcohol[alcohol_anova$quality == i]
23 ## w = 0.96945, p-value = 2.885e-10
24 ##
25 ##
26 ## Shapiro-wilk normality test
27 ##
28 ## data: alcohol_anova$alcohol[alcohol_anova$quality == i]
29 ## w = 0.99166, p-value = 0.3108
30 ##
31 ##
32 ## Shapiro-wilk normality test
33 ##
34 ## data: alcohol_anova$alcohol[alcohol_anova$quality == i]
35 ## w = 0.96336, p-value = 0.6676
```

```

1 #这部分有的没有通过正态性检验，答案中没有考虑这一点，大家见仁见智吧。
2 #H0:不同品质的葡萄酒的酒精浓度没有差异；HA:不同品质的葡萄酒的酒精浓度有显著差异
3 fit <- aov(winequality$alcohol~winequality$quality)
4 summary(fit)

```

```

1 ##                Df Sum Sq Mean Sq F value Pr(>F)
2 ## winequality$quality    1  411.5    411.5    468.3 <2e-16 ***
3 ## Residuals              1597 1403.3      0.9
4 ## ---
5 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

答：不同品质的葡萄酒的酒精浓度有显著差异

5.通过多元线性回归建立一个品质预测模型，并说明哪些变量与品质显著相关。

```

1 ff<-lm(quality~.,data = winequality)
2 summary(ff)

```

```

1 ##
2 ## Call:
3 ## lm(formula = quality ~ ., data = winequality)
4 ##
5 ## Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -2.68911 -0.36652 -0.04699  0.45202  2.02498
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error t value Pr(>|t|)
11 ## (Intercept)      2.197e+01  2.119e+01   1.036   0.3002
12 ## `fixed acidity`    2.499e-02  2.595e-02   0.963   0.3357
13 ## `volatile acidity` -1.084e+00  1.211e-01  -8.948 < 2e-16 ***
14 ## `citric acid`      -1.826e-01  1.472e-01  -1.240   0.2150
15 ## `residual sugar`    1.633e-02  1.500e-02   1.089   0.2765
16 ## chlorides          -1.874e+00  4.193e-01  -4.470  8.37e-06 ***
17 ## `free sulfur dioxide` 4.361e-03  2.171e-03   2.009   0.0447 *
18 ## `total sulfur dioxide` -3.265e-03  7.287e-04  -4.480  8.00e-06 ***
19 ## density            -1.788e+01  2.163e+01  -0.827   0.4086
20 ## pH                 -4.137e-01  1.916e-01  -2.159   0.0310 *
21 ## sulphates          9.163e-01  1.143e-01   8.014  2.13e-15 ***
22 ## alcohol            2.762e-01  2.648e-02  10.429 < 2e-16 ***
23 ## ---
24 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
25 ##
26 ## Residual standard error: 0.648 on 1587 degrees of freedom
27 ## Multiple R-squared:  0.3606, Adjusted R-squared:  0.3561
28 ## F-statistic: 81.35 on 11 and 1587 DF,  p-value: < 2.2e-16

```

挥发性酸度，氯化物，自由二氧化硫量，二氧化硫总量，pH值，硫酸盐和酒精浓度与品质显著相关。

question 4

四、A study was performed to investigate new automobile purchases. A sample of 20 families was selected. Each family was surveyed to determine the age of their oldest vehicle and their total family income. A follow-up survey was conducted six months later to determine if they had actually purchased a new vehicle during that time period ($y = 1$ indicates yes and $y = 0$ indicates no). The data from this study are shown in the Table1.

(a) Fit a logistic regression model to the data.(40')

(b) Interpret the model coefficients β_1 and β_2 and write the logistic regression model formula.(20')

(c) What is the estimated probability that a family with an income of \$45,000 and a car that is five years old will purchase a new vehicle in the next six months?(40')

```
1 income_age <- read.csv("./data/homework-6.4-data.txt")
2 fit<-glm(y~.,data = income_age,family = binomial())
3 summary(fit)
```

```
1 ##
2 ## Call:
3 ## glm(formula = y ~ ., family = binomial(), data = income_age)
4 ##
5 ## Deviance Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -1.5635  -0.8045  -0.1397   0.9535   1.7915
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error z value Pr(>|z|)
11 ## (Intercept) -7.047e+00  4.674e+00  -1.508   0.132
12 ## Income       7.382e-05  6.371e-05   1.159   0.247
13 ## Age          9.879e-01  5.274e-01   1.873   0.061 .
14 ## ---
15 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16 ##
17 ## (Dispersion parameter for binomial family taken to be 1)
18 ##
19 ##      Null deviance: 27.726  on 19  degrees of freedom
20 ## Residual deviance: 21.082  on 17  degrees of freedom
21 ## AIC: 27.082
22 ##
23 ## Number of Fisher Scoring iterations: 5
```

```
1 coef(fit)
```

```
1 ##      (Intercept)      Income      Age
2 ## -7.047061e+00  7.381679e-05  9.878861e-01
```

```
1 test_data<-data.frame(Income=45000, Age=5)
2 test_data$probe<-predict(fit,newdata = test_data,type = "response")
3 test_data
```

```
1 ## Income Age probe
2 ## 1 45000 5 0.7710279
```

答:

公式: $odds = \exp(-7.047061 + 7.381679e - 05x_1 + 9.878861e - 01x_2)$

question 5

五、数据文件“Drivers.csv”为对45名司机的调查结果，其中四个变量的含义为：

- 1) x_1 : 表示视力状况，它是一个分类变量，1表示好，0表示有问题；
- 2) x_2 : 年龄，数值型；
- 3) x_3 : 驾车教育，它也是一个分类变量，1表示参加过驾车教育，0表示没有；
- 4) y : 一个分类型输出变量，表示去年是否出过事故，1表示出过事故，0表示没有；

问题:

(1) 请在R语言中调用logistic回归函数，计算视力状况、年龄、驾车教育与是否发生事故的logistic回归模型，并以“odds=.....”的形式写出回归公式。（10分）

(2) 指出（1）得到的模型中哪些因素对是否发生事故有显著性影响。如果存在对是否发生事故没有显著性影响的因素，请去除这些因素后重新计算logistic回归模型，并以“p=.....”的形式写出回归公式。（20分）

(3) A是一名参加过驾车教育，但视力有问题的50岁老司机；B是一名没有参加过驾车教育，但视力良好的20岁新手。现在A、B都想在某保险公司投保，但按公司规定，被保险人必须满足“明年出事故的概率不高于40%”的条件才能予以承保。请预测A、B两者明年出事故的概率，并告诉保险公司谁可以投保。（20分）

```
1 library(readr)
2 drivers <- read_csv("./data/homework-6.5-Drivers.csv")
```

```
1 ## Parsed with column specification:
2 ## cols(
3 ##   x1 = col_double(),
4 ##   x2 = col_double(),
5 ##   x3 = col_double(),
6 ##   y = col_double()
7 ## )
```

```
1 fit.full<-glm(y~., data=drivers, family=binomial)
```

回归公式为: $odds = \exp(0.597610 - 1.496084x_1 - 0.001595x_2 + 0.315865x_3)$

```
1 fit_x1<-glm(y~x1, data=drivers, family=binomial)
2 summary(fit_x1)
```

```
1 ##
2 ## Call:
```



```

3 ## glm(formula = y ~ x1, family = binomial, data = drivers)
4 ##
5 ## Deviance Residuals:
6 ##      Min       1Q   Median       3Q      Max
7 ## -1.4490  -0.8782  -0.8782   0.9282   1.5096
8 ##
9 ## Coefficients:
10 ##              Estimate Std. Error z value Pr(>|z|)
11 ## (Intercept)   0.6190     0.4688   1.320   0.1867
12 ## x1           -1.3728     0.6353  -2.161   0.0307 *
13 ## ---
14 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15 ##
16 ## (Dispersion parameter for binomial family taken to be 1)
17 ##
18 ##      Null deviance: 62.183  on 44  degrees of freedom
19 ## Residual deviance: 57.241  on 43  degrees of freedom
20 ## AIC: 61.241
21 ##
22 ## Number of Fisher Scoring iterations: 4

```

回归公式为 $p = \exp(0.6190 - 1.3728x_1) / (1 + \exp(0.6190 - 1.3728x_1))$

```

1 test_data<- data.frame(x1=c(0,1))
2 test_data$probe<-predict(fit_x1,test_data ,type='response')
3 test_data

```

```

1 ##   x1 probe
2 ## 1  0  0.65
3 ## 2  1  0.32

```

答：所以A、B两者明年出事故的概率分别为0.65和0.32,因只有B明年出事故的概率不高于40%，故只有B可以投保。

question 6

六、Many digitized image of a fine needle aspirate (FNA) of a breast mass are collected and computed to predict the diagnosis of breast cancer(data.csv).

Attribute information

1) ID number

2) Diagnosis (M = malignant, B = benign)

Ten real-valued features are computed for each cell nucleus:

- 1 a) radius (mean of distances from center to points on the perimeter)
- 2
- 3 b) texture (standard deviation of gray-scale values)
- 4
- 5 c) perimeter
- 6

```

7 d) area
8
9 e) smoothness (local variation in radius lengths)
10
11 f) compactness (perimeter^2 / area - 1.0)
12
13 g) concavity (severity of concave portions of the contour)
14
15 h) concave points (number of concave portions of the contour)
16
17 i) symmetry
18
19 j) fractal dimension ("coastline approximation" - 1)

```

The mean, standard error, and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.

All feature values are recorded with four significant digits. In total, there are 357 benign and 212 malignant samples.

You may need to use proper regression algorithm to train your data, and make predictions.

Instructions:

1) Use all mean features (such as: radius_mean, texture_mean...) to construct a logistic regression model

```

1 library(readr)
2 breast_cancer <- read_csv("./data/homework-6.6-data.csv")

```

```

1 ## Parsed with column specification:
2 ## cols(
3 ##   .default = col_double(),
4 ##   diagnosis = col_character()
5 ## )

```

```

1 ## See spec(...) for full column specifications.

```

```

1 labels <- breast_cancer$diagnosis
2 labels[labels=="M"] <- 1
3 labels[labels=="B"] <- 0
4 labels <- as.integer(labels)
5 breast_cancer$labels <- labels
6 fit.full <-
  glm(labels~radius_mean+texture_mean+perimeter_mean+area_mean+smoothness_mean+compactness_mean+concavity_mean+`concave points_mean`+symmetry_mean+fractal_dimension_mean, data = breast_cancer, family = binomial(), control=list(maxit=100))

```

```

1 ## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```

```
1 summary(fit.full)
```

```
1 ##
2 ## Call:
3 ## glm(formula = lables ~ radius_mean + texture_mean + perimeter_mean +
4 ##       area_mean + smoothness_mean + compactness_mean + concavity_mean +
5 ##       `concave points_mean` + symmetry_mean + fractal_dimension_mean,
6 ##       family = binomial(), data = breast_cancer, control = list(maxit = 100))
7 ##
8 ## Deviance Residuals:
9 ##      Min       1Q   Median       3Q      Max
10 ## -1.95590 -0.14839 -0.03943  0.00429  2.91690
11 ##
12 ## Coefficients:
13 ##              Estimate Std. Error z value Pr(>|z|)
14 ## (Intercept)      -7.35952    12.85259  -0.573   0.5669
15 ## radius_mean      -2.04930     3.71588  -0.551   0.5813
16 ## texture_mean       0.38473     0.06454   5.961 2.5e-09 ***
17 ## perimeter_mean    -0.07151     0.50516  -0.142   0.8874
18 ## area_mean         0.03980     0.01674   2.377   0.0174 *
19 ## smoothness_mean   76.43227    31.95492   2.392   0.0168 *
20 ## compactness_mean  -1.46242    20.34249  -0.072   0.9427
21 ## concavity_mean     8.46870     8.12003   1.043   0.2970
22 ## `concave points_mean` 66.82176    28.52910   2.342   0.0192 *
23 ## symmetry_mean     16.27824    10.63059   1.531   0.1257
24 ## fractal_dimension_mean -68.33703    85.55666  -0.799   0.4244
25 ## ---
26 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
27 ##
28 ## (Dispersion parameter for binomial family taken to be 1)
29 ##
30 ##    Null deviance: 751.44  on 568  degrees of freedom
31 ## Residual deviance: 146.13  on 558  degrees of freedom
32 ## AIC: 168.13
33 ##
34 ## Number of Fisher Scoring iterations: 9
```

Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred:这个问题实际是因为数据集问题，这个数据集本身接近线性可分了，所以导致模型过拟合。具体原因大家自己检索一下吧。可以看看这篇博客：<https://www.cnblogs.com/runner-ljt/p/4574275.html>

2)Then try to reduce the number of features from your last model, construct another regression model, and you will need to write down the equation of your logistic regression model(Tips:

$$\text{Logit}P = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p)$$

除去不显著相关的变量。

```
1 fit.reduced<-glm(lables~texture_mean+area_mean+smoothness_mean+`concave
points_mean`,data = breast_cancer,family = binomial(),control=list(maxit=100))
```

```
1 ## warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
1 summary(fit.reduced)
```

```
1 ##
2 ## Call:
3 ## glm(formula = lables ~ texture_mean + area_mean + smoothness_mean +
4 ##       `concave points_mean`, family = binomial(), data = breast_cancer,
5 ##       control = list(maxit = 100))
6 ##
7 ## Deviance Residuals:
8 ##      Min       1Q   Median       3Q      Max
9 ## -2.31798  -0.15623  -0.04212   0.01662   2.84201
10 ##
11 ## Coefficients:
12 ##              Estimate Std. Error z value Pr(>|z|)
13 ## (Intercept)    -23.677816    3.882774  -6.098 1.07e-09 ***
14 ## texture_mean     0.362687    0.060544   5.990 2.09e-09 ***
15 ## area_mean        0.010342    0.002002   5.165 2.40e-07 ***
16 ## smoothness_mean  59.471304   25.965153   2.290  0.022 *
17 ## `concave points_mean` 76.571210  16.427864   4.661 3.15e-06 ***
18 ## ---
19 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
20 ##
21 ## (Dispersion parameter for binomial family taken to be 1)
22 ##
23 ##    Null deviance: 751.44  on 568  degrees of freedom
24 ## Residual deviance: 156.44  on 564  degrees of freedom
25 ## AIC: 166.44
26 ##
27 ## Number of Fisher Scoring iterations: 8
```

公式: $LogitP = -23.677816 + 0.362687x_2 + 0.010342x_4 + 59.47130x_5 + 76.571210x_8$ 3) Use proper test to test the difference between two models

对两个模型进行ANOVA分析。

```
1 #H0:两模型无显著差别; HA:两模型有显著差别
2 anova(fit.full, fit.reduced, test = "Chisq")
```

```
1 ## Analysis of Deviance Table
2 ##
3 ## Model 1: lables ~ radius_mean + texture_mean + perimeter_mean + area_mean +
4 ##       smoothness_mean + compactness_mean + concavity_mean + `concave points_mean` +
5 ##       symmetry_mean + fractal_dimension_mean
6 ## Model 2: lables ~ texture_mean + area_mean + smoothness_mean + `concave points_mean`
7 ##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
8 ## 1      558      146.13
9 ## 2      564      156.44 -6    -10.31   0.1122
```

$p=0.1122>0.05$, 两模型无显著差别。

4) You may split the data properly, use part of them to train your regression model and use another part to make predictions. Lastly, you may try to calculate the accuracy of your model. (Tips: To split the data, you can use the first 398 rows as training data, use the last 171 rows as prediction data. The predict function returns a value between 0 and 1, 0~0.5 belong to the first class, and 0.5~1 belong to second class in binary classification problems)

```
1 train_data <- breast_cancer[1:398,]
2 test_data <- breast_cancer[399:569,]
3 train_fit <- glm(lables ~ texture_mean + area_mean + smoothness_mean + `concave points_mean`, data
  = train_data, family = binomial())
4 summary(train_fit)
```

```
1 ##
2 ## Call:
3 ## glm(formula = lables ~ texture_mean + area_mean + smoothness_mean +
4 ##       `concave points_mean`, family = binomial(), data = train_data)
5 ##
6 ## Deviance Residuals:
7 ##      Min       1Q   Median       3Q      Max
8 ## -2.39278  -0.14454  -0.02447   0.03635   2.60665
9 ##
10 ## Coefficients:
11 ##              Estimate Std. Error z value Pr(>|z|)
12 ## (Intercept)    -27.47397     4.74798  -5.786 7.19e-09 ***
13 ## texture_mean      0.46244     0.08434   5.483 4.19e-08 ***
14 ## area_mean        0.01082     0.00235   4.606 4.11e-06 ***
15 ## smoothness_mean  90.11221    30.96961   2.910 0.003618 **
16 ## `concave points_mean` 59.01212    17.51779   3.369 0.000755 ***
17 ## ---
18 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
19 ##
20 ## (Dispersion parameter for binomial family taken to be 1)
21 ##
22 ##    Null deviance: 544.93  on 397  degrees of freedom
23 ## Residual deviance: 108.30  on 393  degrees of freedom
24 ## AIC: 118.3
25 ##
26 ## Number of Fisher Scoring iterations: 8
```

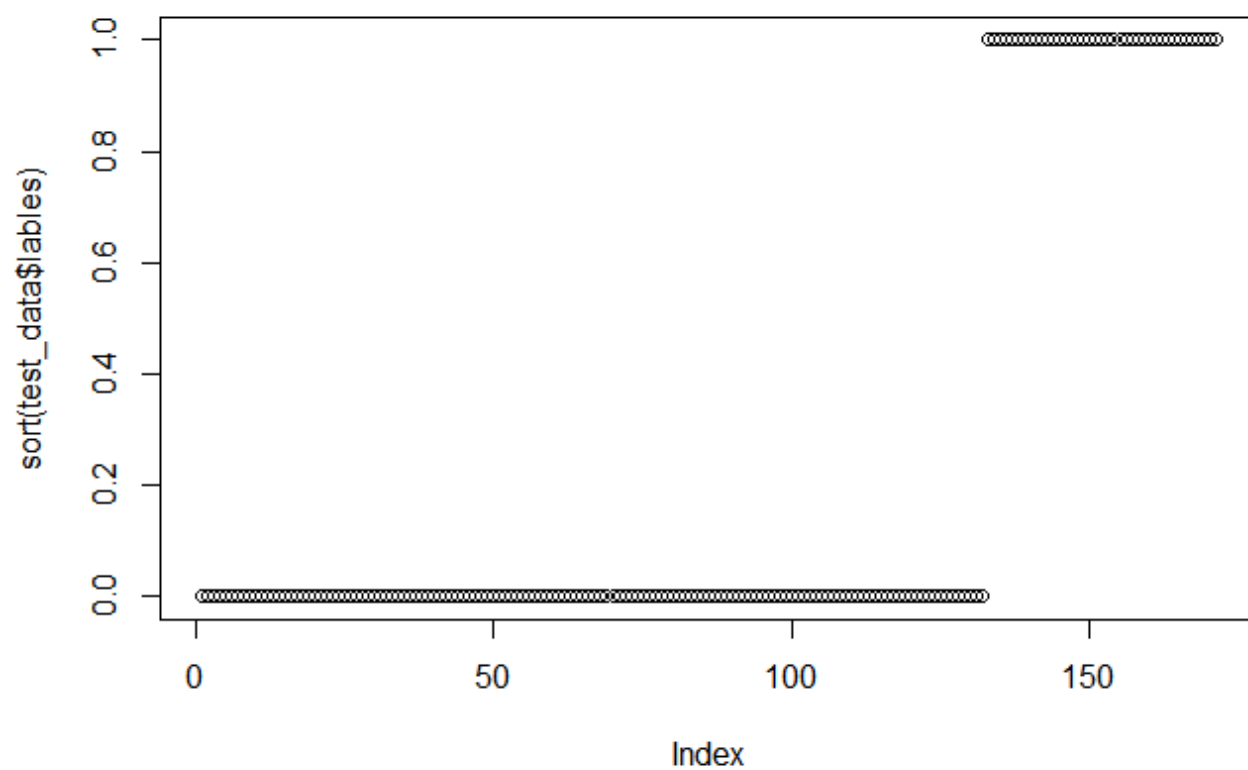
```
1 test_data$probe <- predict(train_fit, newdata = test_data, type = "response")
2 lables <- test_data$lables
3 pred_lables <- ifelse(test_data$probe > 0.5, 1, 0)
4 mean(pred_lables == lables)
```

```
1 ## [1] 0.9064327
```

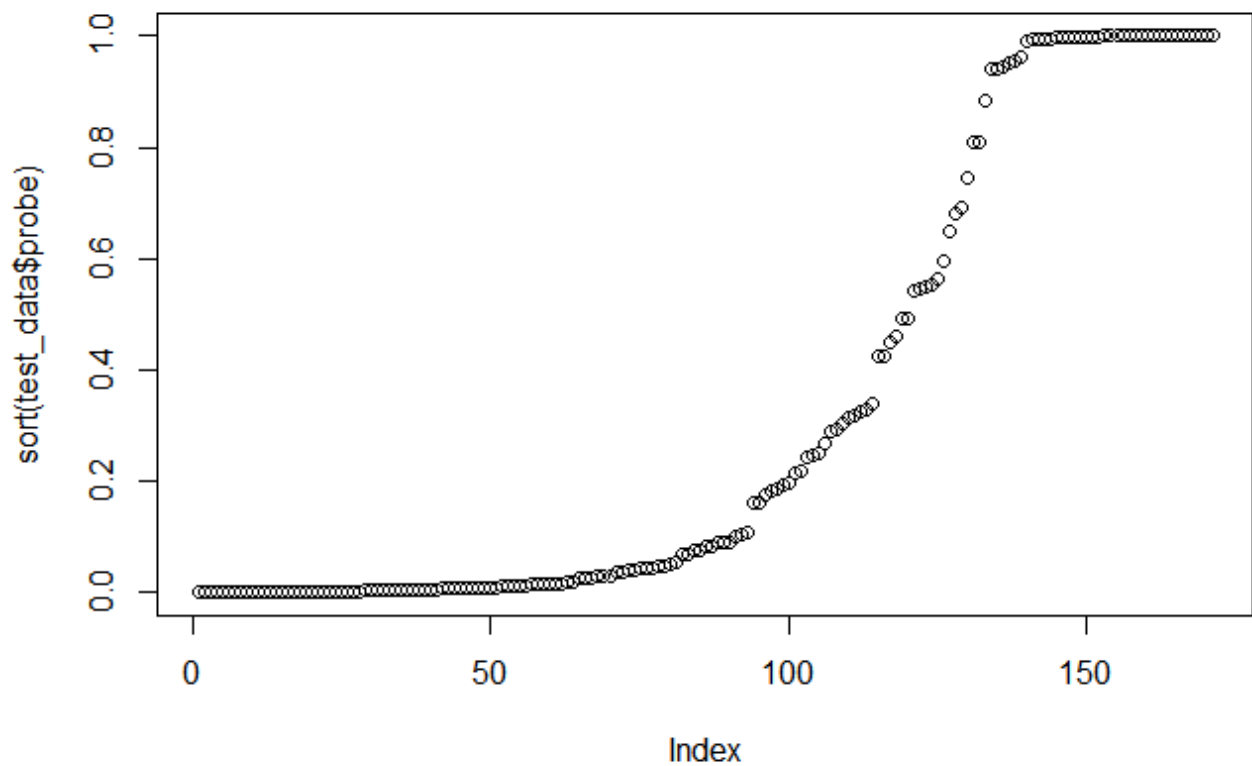
The accuracy is 0.9064

进一步地，可以看看它们的预测情况和ROC曲线以及AUC值。

```
1 | plot(sort(test_data$labels))
```



```
1 | plot(sort(test_data$probe))
```



```
1 # install.packages("pROC")
2 library('pROC')
```

```
1 ## Type 'citation("pROC")' for a citation.
```

```
1 ##
2 ## Attaching package: 'pROC'
```

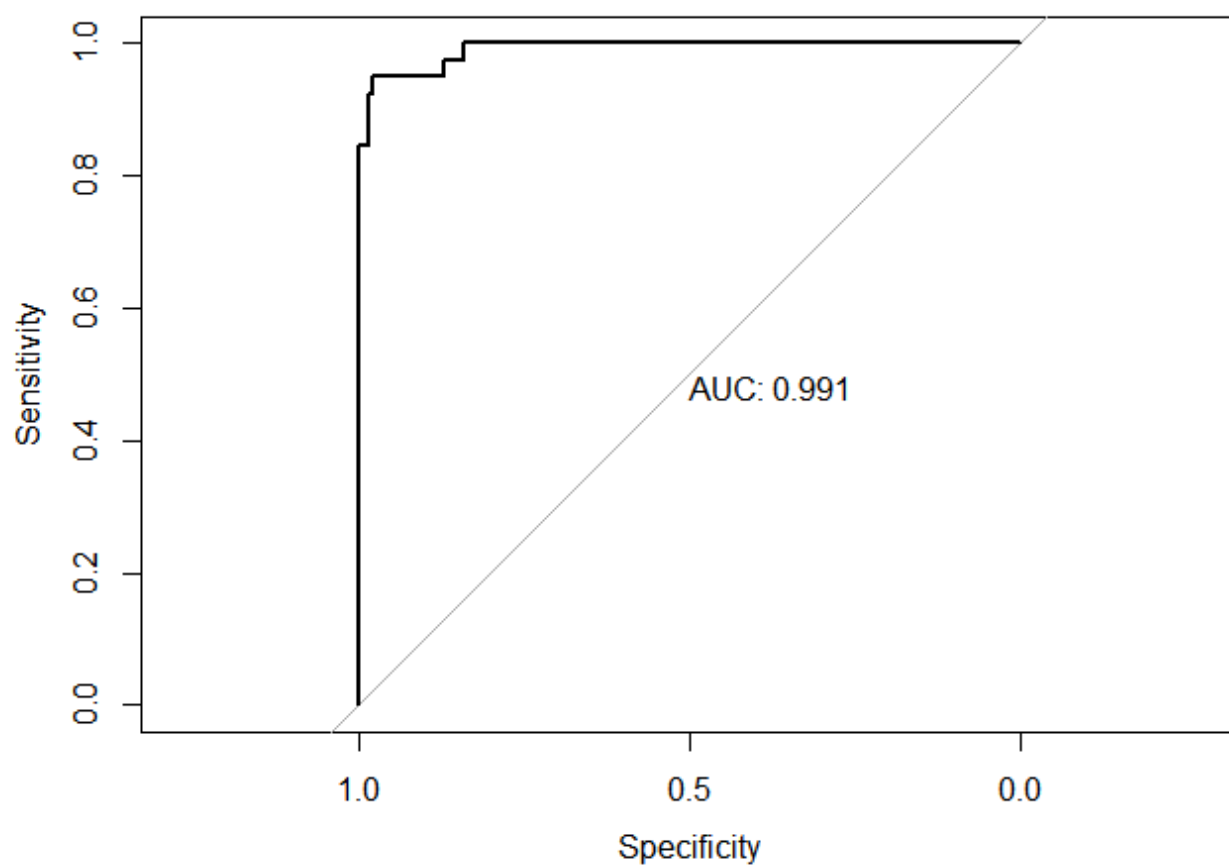
```
1 ## The following objects are masked from 'package:stats':
2 ##
3 ##     cov, smooth, var
```

```
1 model_roc <- roc(test_data$labes,test_data$probe)
```

```
1 ## Setting levels: control = 0, case = 1
```

```
1 ## Setting direction: controls < cases
```

```
1 plot(model_roc,print.auc = T)
```



AUC = 0.991, 可以说明分类效果非常好。