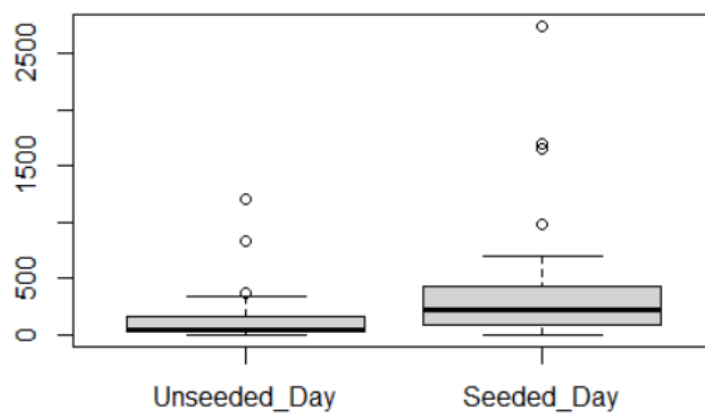


PS3

曾海翔 12032760

```
1 library(tidyr)
2 library(dplyr)
3 library(ggplot2)
4
5 #1
6 #1.1
7 #set the array
8 Unseeded_Day <- c(1202.6, 830.1, 372.4, 345.5, 321.2, 244.3, 163.0, 147.8, 95.0, 87.0, 81.2, 68.5, 47.3, 41.1, 36.1)
9 Seeded_Day <- c(2745.6, 1697.1, 1656.4, 978.0, 703.4, 489.1, 430.0, 334.1, 302.8, 274.7, 274.7, 255.0, 242.5, 200.1)
10 Rainfall_day <- cbind(Unseeded_Day, Seeded_Day)
11 boxplot(Rainfall_day)
12 #1.2
13 Rainfall_day_tbl <- as_tibble(Rainfall_day)
14 Rainfall_day_tbl %>%
15   summarise(
16     mean_Seeded = mean(Rainfall_day_tbl$Seeded_Day),
17     mean_Unseeded = mean(Rainfall_day_tbl$Unseeded_Day),
18     increse_rate = mean(((Rainfall_day_tbl$Seeded_Day)-mean(Rainfall_day_tbl$Unseeded_Day))/mean(Rainfall_day_tbl$Unseeded_Day))
19 )
```

#1



通过 cbind 把接种和未接种的天气降雨量组成一个矩阵，使用 boxplot 画图

```
# A tibble: 1 x 3
  mean_Seeded mean_Unseeded increse_rate
  <dbl>        <dbl>        <dbl>
1    442.        165.         1.69
> t.test(Unseeded_Day, Seeded_Day)

Welch Two Sample t-test

data: Unseeded_Day and Seeded_Day
t = -1.9983, df = 33.856, p-value = 0.05376
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -559.552533  4.752533
sample estimates:
mean of x mean of y
 164.5731  441.9731
```

分别的平均数，降雨增加率以及 t 检验结果

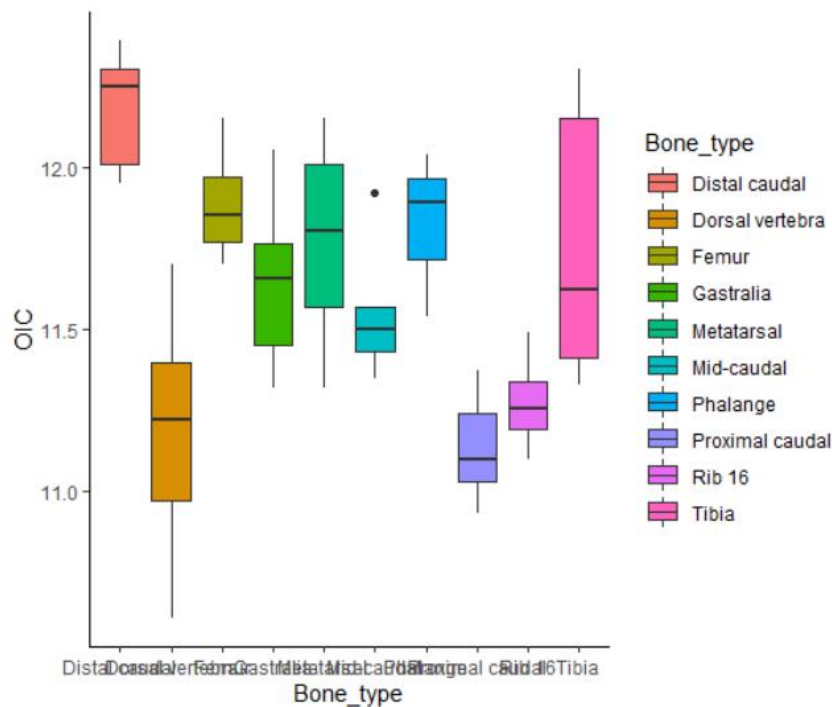
#2

```

22 #2
23 TR_bone <- read.csv("PS3_data.csv",header = T)
24 TR_bone_tbl <- as_tibble(TR_bone)
25 TR_bone_tbl %>%
26   group_by(Bone_type) %>%
27   summarise(
28     count = n(),
29     mean_OIC = mean(OIC),
30     sd_OIC = sd(OIC)
31   )
32 #use boxplot to check if there is a different between the distribution
33 ggplot(TR_bone_tbl, aes(x=Bone_type, y=OIC, fill=Bone_type))+
34   geom_boxplot()+
35   theme_classic()
36 #anova one way test
37 AOW_TR <- aov(OIC~Bone_type, data = TR_bone_tbl)
38 summary(AOW_TR)

```

把数据复制进 excel 表格中，另存为 csv 格式（见附件）。使用 ggplot 中的 boxplot 画图用 aov() 进行单因素显著性差异分析，组间有显著性差异性，所以霸王龙可能为变温动物。



	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Bone_type	9	5.688	0.6320	7.922	1.01e-06 ***
Residuals	42	3.351	0.0798		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

#3

```
40 #3
41 V_Z <- read.csv("PS3_data.csv",header = T)
42 V_Z_tbl <- as_tibble(V_Z)
43 names(V_Z_tbl)
44 V_Z_tbl %>%
45   filter(people_type != "") %>%
46   select(people_type,zine_level) %>%
47   group_by(people_type) %>%
48   summarise(
49     count = n(),
50     mean_Zine = mean(zine_level),
51     sd_Zine = sd(zine_level)
52   )
53 AOW_Zine <- aov(zine_level~people_type, data = V_Z_tbl)
54 summary(AOW_Zine)
```

在管道里用 filter 筛选出有效的数据，按实验对象类型分类，分析 Zn 平均值。

```
`summarise()` ungrouping output (override with `.groups` argument)
# A tibble: 3 x 4
  people_type      count mean_Zine sd_Zine
  <chr>          <int>    <dbl>  <dbl>
1 Nonpregnant vegetarians     5    179.   27.3
2 Pregnant nonvegetarians     6    178   14.5
3 Pregnant vegetarians    12    177.   20.9
> AOW_Zine <- aov(zine_level~people_type, data = V_Z_tbl)
> summary(AOW_Zine)
              Df Sum Sq Mean Sq F value Pr(>F)
people_type    2     16    8.1   0.018  0.982
Residuals    20   8816   440.8
29 observations deleted due to missingness
> |
```

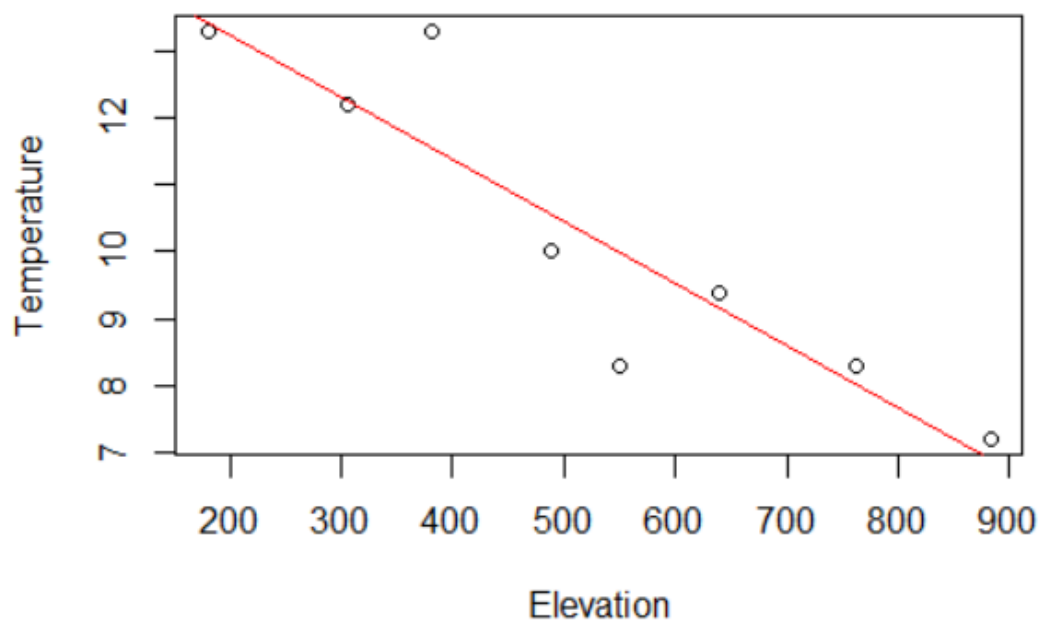
经分析发现组间无显著性差异

#4

```
56 #4
57 Elevation <- c(180,305,381,488,549,640,762,883)
58 Temperature <- c(13.3,12.2,13.3,10.0,8.3,9.4,8.3,7.2)
59 matrix1 <- data.frame(Elevation, Temperature)
60 plot(Temperature~Elevation, data = matrix1)
61 regression_line <- lm(Temperature~Elevation, data = matrix1)
62 abline(regression_line, col='red')
63 summary(regression_line)$coefficients[2,1]
64 #summary(regression_line)$coefficients[,1] or coef(regression_line) 第一个是截距，第二个是斜率
65 #learn from https://blog.csdn.net/dingchenxixi/article/details/50543822
66 lapse_rate <- summary(regression_line)$coefficients[2,1]*1000
67 lapse_rate
68
```

用 `lm()` 函数生成拟合直线，用 `abline` 把直线添加进已画好的图表里。

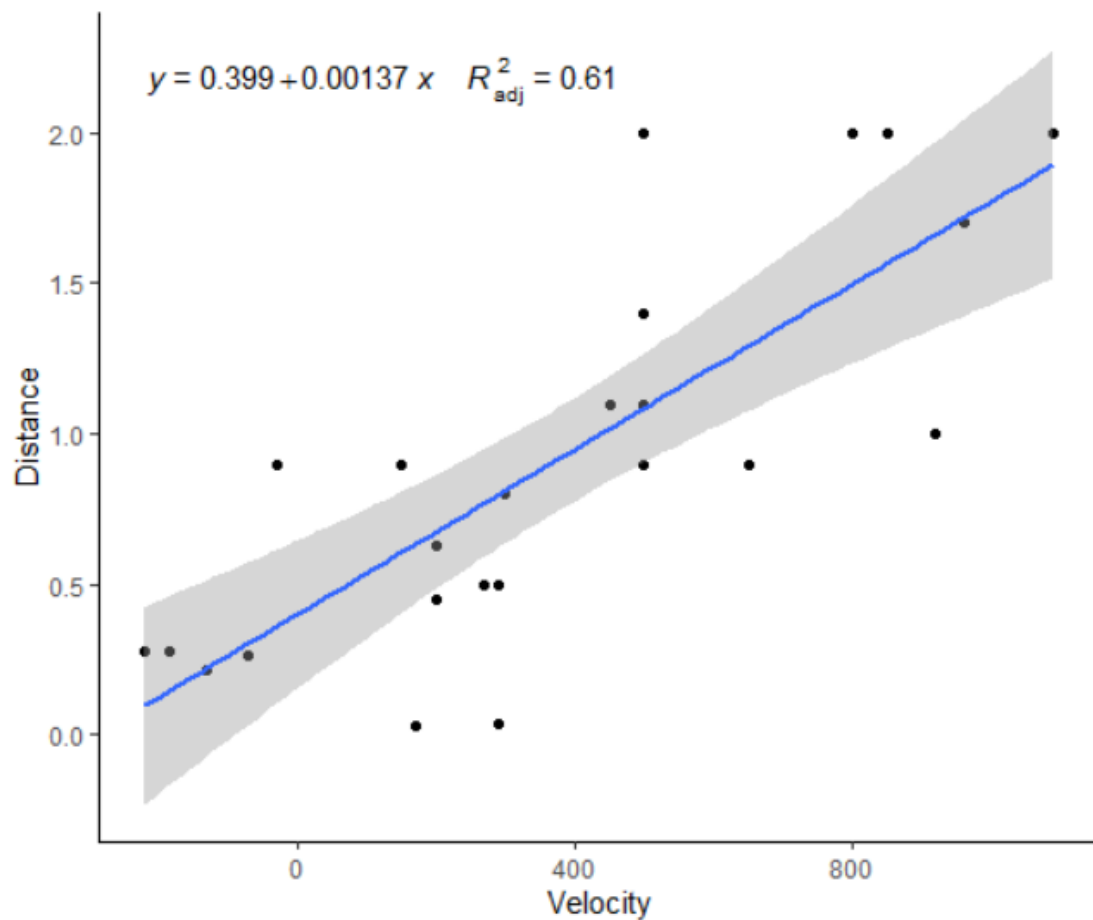
在 `summary` 中寻找斜率的对应值，经过调整单位运算后得出结果。计算的斜率为-9.312104



```
25 observations deleted due to missingness
> Elevation <- c(180,305,381,488,549,640,762,883)
> Temperature <- c(13.3,12.2,13.3,10.0,8.3,9.4,8.3,7.2)
> matrix1 <- data.frame(Elevation, Temperature)
> plot(Temperature~Elevation, data = matrix1)
> regression_line <- lm(Temperature~Elevation, data = matrix1)
> abline(regression_line, col='red')
> summary(regression_line)$coefficients[2,1]
[1] -0.009312104
> #summary(regression_line)$coefficients[,1] or coef(regression_line) 第一个是截距，第二个是斜率
> #learn from https://blog.csdn.net/dingchenxixi/article/details/50543822
> lapse_rate <- summary(regression_line)$coefficients[2,1]*1000
> lapse_rate
[1] -9.312104
>
```

#5

```
69 #5
70 #5.1
71 install.packages("ggpmisc")
72 library(ggpmisc)
73 BBT <- read.csv("PS3_data.csv",header = T)
74 BBT_tbl <- as_tibble(BBT)
75 names(BBT_tbl)
76 BBT_tbl %>%
77   filter(Nebula != "") %>%
78   select(Nebula,Velocity,Distance) %>%
79   ggplot(aes(y = Distance,x = Velocity))+
80   geom_point()+
81   #5.2
82   geom_smooth(method="lm", formula = y ~ x)+ # learn from https://blog.csdn.net/weixin\_42933967/article/details/96
83   stat_poly_eq(aes(label = paste(..eq.label.., ..adj.rr.label.., sep = '~~~~')), formula = y ~ x, parse = T) +
84   theme_classic()
85 #5.3
86 #the first assumption sounds reasonable because the universe
87 #was come from a single point( big bang theory)
88 #5.4
89 |
```



调用了新的包：ggpmisc 用于计算回归方程。经计算得宇宙的年龄为 990 亿年。

#5.4

如果测量的点能更准确，得到的回归曲线也当然更精确。

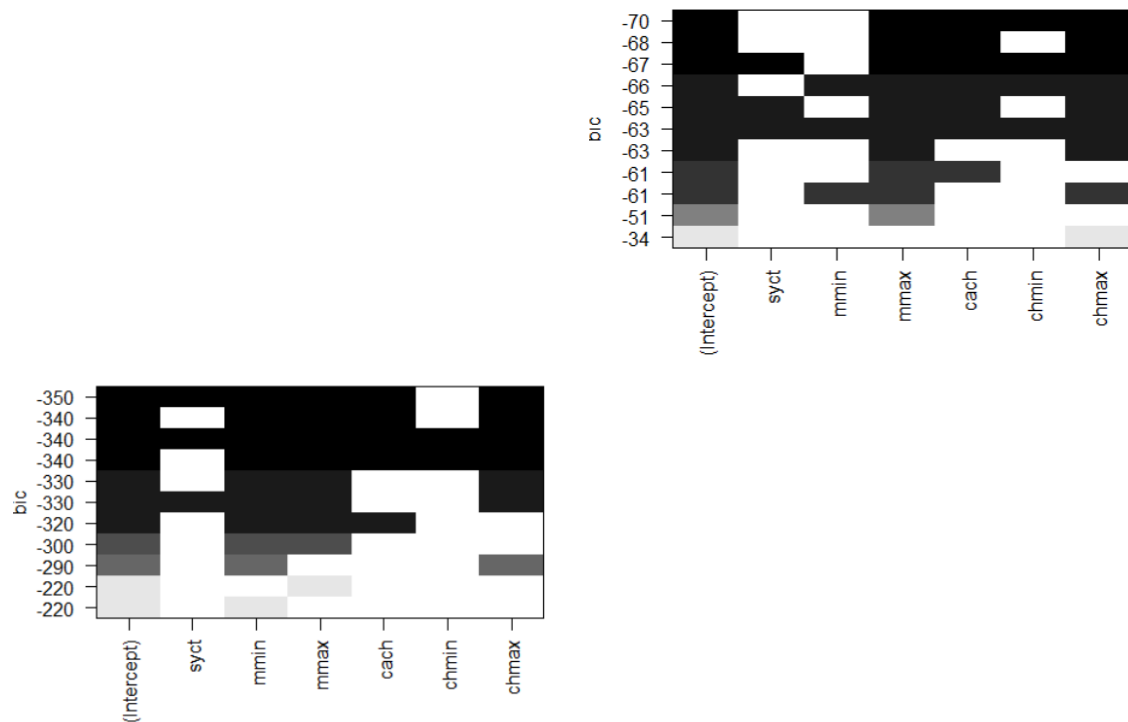
#6

```

90 #6
91 install.packages('leaps')
92 library(tidyr)
93 library(dplyr)
94 library(ggplot2)
95 library(leaps)
96 library(MASS)
97 data(cpus)
98 head(cpus)
99 nrow(cpus)
100 #make the order random
101 random_selcet <- runif(209,0,100)
102 cpu1 <- cbind(cpus,random_selcet)
103 cpu_random <- cpu1 %>%
104   arrange(desc(random_selcet))
105 #seperate it into two group
106 cpu_train_set <- cpu_random[1:167,]
107 cpu_test_set <- cpu_random[168:209,]
108 #6.1
109 train_subset_result <- regsubsets(perf ~ syct+mmin+mmax+cach+chmin+chmax, data=cpu_train_set, nbest=2, nvmax = 6)
110 plot(train_subset_result)
111
112 #6.2
113 test_subset_result <- regsubsets(perf ~ syct+mmin+mmax+cach+chmin+chmax, data=cpu_test_set, nbest=2, nvmax = 6)
114 plot(test_subset_result)

```

在 `cpu` 组后加上一列从 1 到 100 的随机数，按大小排列。用 `nrow` 读出矩阵长度，按 80%和 20%的数据量划分为两组。用 `regsubset` 进行回归拟合（虽然还看不懂）



左图为学习（train）组的，右图为测试（test）组的

#7

```
116 #7
117 data_Lab <- read.csv("CMMNDHgfinal.csv",header = T)
118 data_tbl <- as_tibble(data_Lab)
119 names(data_tbl)
120 #which kind of bird have the highest Hg concentration
121 data_tbl %>%
122   filter(Species != "") %>%
123   group_by(Species) %>%
124   #summarise(Hg_mean = mean(Hgppm)) %>%
125   #arrange(desc(Hg_mean)) %>%
126   ggplot(aes(x=Species,y=Hgppm),head = T, fill=Species)+
127   theme_classic()+
128   geom_boxplot()
129
130 #t test
131 MYWA <- data_tbl %>%
132   filter(Species == "MYWA")%>%
133   select(Hgppm)
134 NOWA <- data_tbl %>%
135   filter(Species == "NOWA")%>%
136   select(Hgppm)
137 t.test(MYWA,NOWA)
138
139 #ANOVA test
140 anova_bird <- aov(Hgppm~Species,data=data_tbl)
141 summary(anova_bird)
142
143 #linear regression model
144 NOWA_LG <- data_tbl %>%
145   filter(Species == "NOWA")
146 plot(Deterium~Hgppm, data = NOWA_LG)
147 regression_line <- lm(Deterium~Hgppm, data = NOWA_LG)
148 abline(regression_line, col='red')
149
```

用了组里马艳菊师姐的鸟类——重金属（Hg）的数据作为分析对象。对汞浓度最高的两种鸟类进行 t 检验和显著性差异分析。

回归曲线用了汞浓度和实验分析所需材料？进行分析（好像不是很好看）

```
> t.test(MYWA,NOWA)
```

Welch Two Sample t-test

```
data: MYWA and NOWA
t = 0.20376, df = 28.834, p-value = 0.84
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -2.656957  3.244789
sample estimates:
mean of x mean of y
 4.572010  4.278094
```

```
> summary(anova_bird)
              Df Sum Sq Mean Sq F value Pr(>F)
Species         7   430.8    61.55   21.66 <2e-16 ***
Residuals      477 1355.3     2.84
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
9 observations deleted due to missingness
```

