

# Report for HW2

12032924 李熹成

## Problem #1

```
1 # 1.Significant earthquakes since 2150 B.C.
2 ##1.1
3 library(tidyr)
4 library(dplyr)
5 library(ggplot2)
6 SED<-read.csv(file = 'signif.txt',header = T,sep = '\t')
7 class(SED)
8 Seq_Eqs<-as_tibble(SED)
9 ##1.2
10 Seq_Eqs %>%
11   group_by(COUNTRY) %>%
12   select(DEATHS, YEAR, COUNTRY) %>%
13   summarise(total_num_dth=sum(DEATHS)) %>%
14   arrange(desc(total_num_dth))>rank_death
15 rank_death[1:10,]

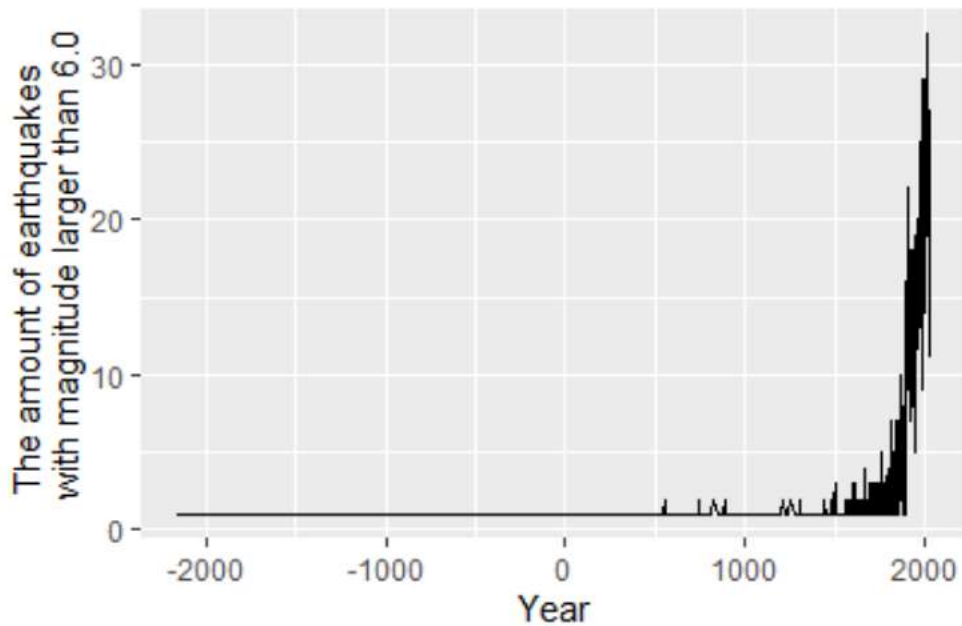
16 ##1.3
17 Seq_Eqs %>%
18   filter(EQ_PRIMARY>6) %>%
19   group_by(YEAR) %>%
20   summarise(ersq_amount=n()) %>%
21   ggplot(aes(x=YEAR,y=ersq_amount))+
22   geom_line()+
23   scale_x_continuous(name = 'Year')+
24   scale_y_continuous(name = 'The amount of earthquakes
25     with magnitude larger than 6.0')
26 ###From the plot above, the earthquakes recorded is more frequently.
27 ###However, it could also result from the detective techniques advanced.

29 ##1.4
30 Acountry<-readline(prompt="Please enter a country you want to observe:")
31 CountEq_LargestEq<-function(Acountry){
32   Seq_Eqs %>%
33     filter(COUNTRY==Acountry& EQ_PRIMARY!='NA') %>%
34     mutate(ThatDate=paste(YEAR,MONTH,DAY,sep = '-')) %>%
35     select(ThatDate,EQ_PRIMARY) %>%
36     summarise(ersq_amount2=n(),max_level_date=
37       ThatDate[which(EQ_PRIMARY==max(EQ_PRIMARY))])>C_D
38   C_D_NH<-unname(C_D)#remove dimname
39   return(C_D_NH)
40 }

42 ###Get rid of countries with earthquake magnitude equal to 'NA'
43 Seq_Eqs %>%
44   filter(EQ_PRIMARY!='NA')>Seq_Eqs_noNA
45
46 i=1
47 NewMat<-matrix(ncol = 3,nrow = length(unique(Seq_Eqs_noNA$COUNTRY)))
48 for(CountryName in unique(Seq_Eqs_noNA$COUNTRY)){
49   NewMat[i,<-c(as.character(CountryName),
50     as.numeric(CountEq_LargestEq(CountryName)[1,1]),
51     as.character(CountEq_LargestEq(CountryName)[1,2]))
52   i=i+1
53 }
54 #Sort in descending order by earthquake numbers.
55 NewMat_Order<-NewMat[order(as.numeric(NewMat[,2]),decreasing=T),]
56 NewMat_Order
```

1.2 I used desc() to rank top 10 countries along with the total number of death in 1.2.

1.3 Plot the graph and added labels to the axis. The result is as below:



The earthquake magnitude more than 6 seems to become more frequently. However, it probably because the monitoring technique is getting more advanced over hundreds of years.

1.4 The results are as below:

```
> NewMat_Order
      [,1]      [,2]      [,3]
[1,] "CHINA"    "575"    "1668-7-25"
[2,] "JAPAN"    "343"    "2011-3-11"
[3,] "INDONESIA" "314"    "2004-12-26"
[4,] "IRAN"     "249"    "856-12-22"
[5,] "USA"      "215"    "1964-3-28"
[6,] "TURKEY"   "206"    "1912-8-9"
[7,] "GREECE"   "152"    "365-7-21"
[8,] "PERU"     "146"    "1716-2-6"
[9,] "CHILE"    "145"    "1960-5-22"
[10,] "RUSSIA"  "139"    "1952-11-4"
[11,] "PHILIPPINES" "132" "1897-9-21"
[12,] "MEXICO"  "119"    "1899-1-24"
[13,] "ITALY"   "96"     "1915-1-13"
[14,] "TAIWAN"  "93"     "1920-6-5"
[15,] "PAPUA NEW GUINEA" "89" "1919-5-6"
[16,] "INDIA"   "81"     "1950-8-15"
[17,] "NEW ZEALAND" "62" "1826-NA-NA"
```

Create a function named *CountEq\_LargestEq()* function to pick the times and date-of-maxEq according to country name. Then using a *for* loop to pick all countries names from the table and run *CountEq\_LargestEq()*. Noted that before the operations above, I remove the country name with earthquake magnitude of NA. Finally, I combined the element in each output dataframe and convert them into a matrix and output the result in decreasing order.

## Problem #2

```
58 #2. wind speed in Shenzhen during the past 10 years
59
60 shenzhenData<-read.csv(file = '2281305.csv',header = T)
61 class(shenzhenData)
62 SZD_tbl<-as_tibble(shenzhenData)
63
64 SZD_tbl %>%
65   mutate(WD_angle=as.numeric(substr(WND,1,3)),
66          WD_DQC=substr(WND,5,5),
67          WD_TC=substr(WND,7,7),
68          WD_Speed=as.numeric(substr(WND,9,12)),
69          WD_SQC=substr(WND,14,14),
70          Months=substr(DATE,1,7)) %>%
71   mutate(WD_angle_New=ifelse(WD_angle==999,'NA',WD_angle),
72          WD_DQC_New=ifelse(WD_DQC=='3'|WD_DQC=='7','NA',WD_DQC),
73          WD_TC_New=ifelse(WD_TC=='9','NA',WD_TC),
74          WD_Speed_New=ifelse(WD_TC==9999,'NA',WD_Speed),
75          WD_SQC_New=ifelse(WD_SQC=='3'|WD_SQC=='7','NA',WD_SQC)) %>%
76
77   # Filter(WD_angle_New!='NA',
78   #       WD_DQC_New!='NA',
79   #       WD_TC_New!='NA',
80   #       WD_Speed_New!='NA',
81   #       WD_SQC_New!='NA') %>%
82   select(WD_angle_New, WD_DQC_New,
83          WD_TC_New,WD_Speed_New,WD_SQC_New,Months) %>%
84   group_by(Months) %>%
85   summarise(MonAvgWS=mean(WD_Speed_New)) %>%
86   mutate(Months_day=paste(Months,'1',sep='-')) %>%
87   ggplot(aes(x=as.Date(Months_day),y=MonAvgWS,group = 1))+
88   geom_line()+
89   xlab('Year')+
90   ylab('wind speed')+
91   labs(title = 'Average monthly wind speed in Shenzhen in 2010-2020')
```

I refer to the descriptive document and remove the data of complete error.

The results are as below:



## Problem#3

```
95 #3. Revisit a data set
96 Madrid<-read.csv(file = 'c82210-1.csv',header = T)
97 class(Madrid)
98 class(Madrid$Tm)
99 Madrid_tbl<-as_tibble(Madrid)
100
101 Madrid_tbl %>%
102   mutate(Date=as.Date(paste(Y,M,D,sep = '-')))) %>%
103   ggplot(aes(x=Date,y=AT))+
104   geom_line()+
105   xlab('Year')+
106   ylab('Temperature(°C)')+
107   labs(title = 'Average daily temperature
108         in Madrid from 1991-1995')
109
110
111 Madrid_tbl %>%
112   mutate(MonthM=as.Date(paste(Y,M,'1',sep = '-')))) %>%
113   group_by(MonthM) %>%
114   summarise(avgmonthAT=mean(AT)) %>%
115   ggplot(aes(x=MonthM,y=avgmonthAT))+
116   geom_line()+
117   xlab('Year')+
118   ylab('Temperature(°C)')+
119   labs(title = 'Average monthly temperature
120         in Madrid from 1991-1995')
121
122 Madrid_tbl %>%
123   mutate(Date2=as.Date(paste(Y,M,D,sep = '-')))) %>%
124   mutate(diff9195=TM-Tm) %>%
125   ggplot(aes(x=Date2,y=diff9195))+
126   geom_line()+
127   xlab('Year')+
128   ylab('Temperature(°C)')+
129   labs(title = 'Daily temperature difference in 1991-1995')
130
131 class(diff9195)
132 shapiro.test(diff9195)
133 ### p-value is 1.515e-10 satisfying normal distribution.
134
135 Madrid_tbl %>%
136   mutate(Date2=as.Date(paste(Y,M,D,sep = '-')))) %>%
137   mutate(diff9195=TM-Tm) %>%
138   ggplot(aes(x=as.numeric(H),y=diff9195))+
139   geom_point()+
140   xlab('Humidity')+
141   ylab('difference of temperature')+
142   labs(title = 'Scatter plot of Humidity in respect of Temperature Difference')
143
144
145 Madrid_tbl %>%
146   filter(H!='-') %>%
147   mutate(Date2=as.Date(paste(Y,M,D,sep = '-')))) %>%
148   mutate(diff9195=TM-Tm) %>%
149   pull(diff9195)->diff9195
150
151 Madrid_tbl %>%
152   mutate(Date2=as.Date(paste(Y,M,D,sep = '-')))) %>%
153   filter(H!='-') %>%
154   pull(as.numeric(H))->H
155
156 ##p-value = 5.204e-16,is normally distributed
157 class(diff9195)
158 shapiro.test(diff9195)
159 ##p-value < 2.2e-16,is normally distributed
160 class(H)
161 shapiro.test(as.numeric(H))
162
163 r3<-cor(as.numeric(H),diff9195)
```



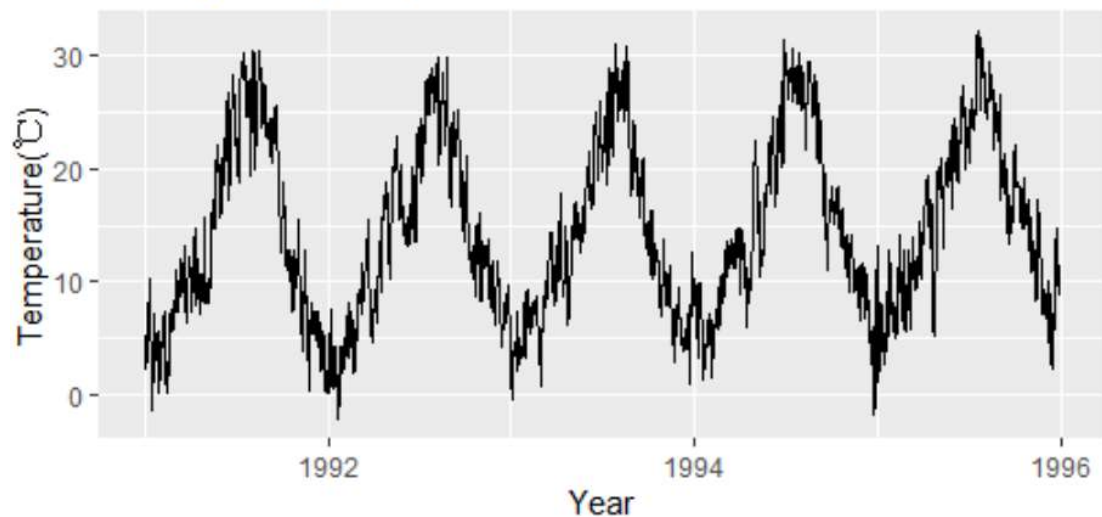
```

164 Madrid_tbl %>%
165   filter(H!='-') %>%
166   mutate(diff9195=TM-Tm) %>%
167   group_by(H) %>%
168   summarise(diffavg=mean(diff9195)) %>%
169   mutate(nH=as.numeric(H)) %>%
170   ggplot(aes(x=nH, y=diffavg)) +
171     geom_point()+
172     xlab('Humidity')+
173     ylab('difference of temperature')+
174     labs(title = 'Scatter plot of Humidity in respect of Temperature Difference
175              (doing mean to Humidity in each Humidity)')
176
177 hdregrression<-lm(HD$diffavg~HD$nH,data=HD)
178 summary(hdregrression)
179 abline(hdregrression)

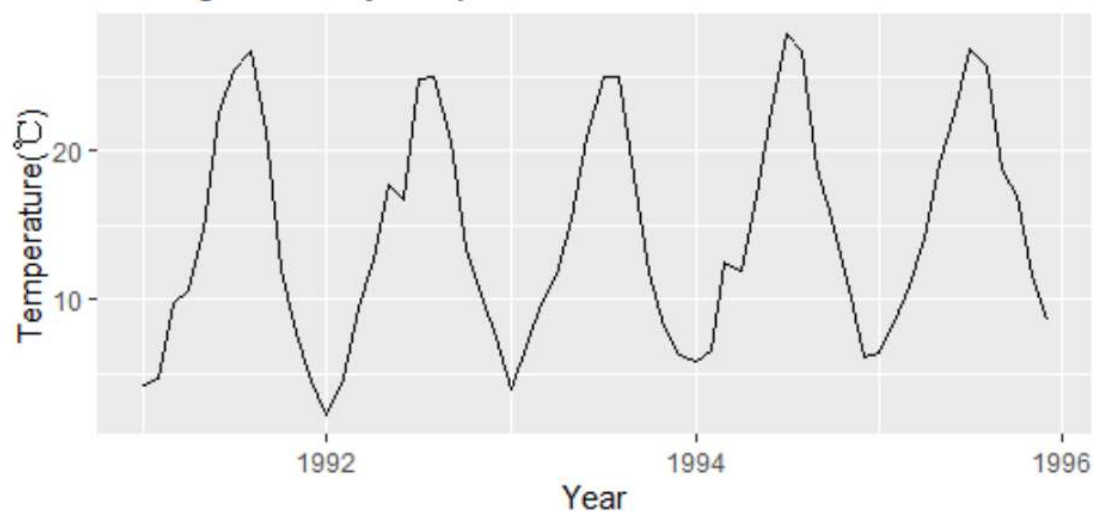
```

The results are as below:

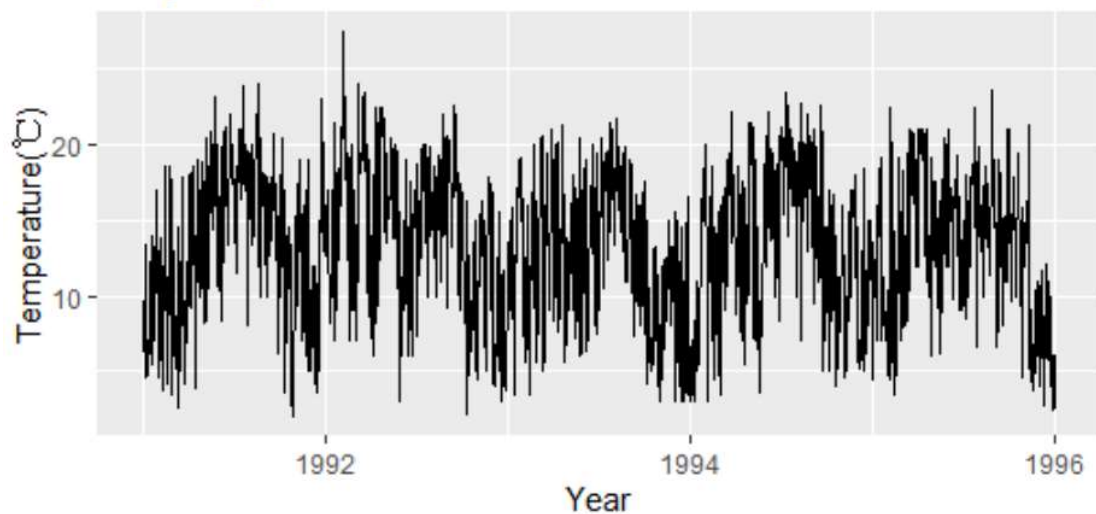
**Average daily temperature in Madrid from 1991-1995**



**Average monthly temperature in Madrid from 1991-1995**

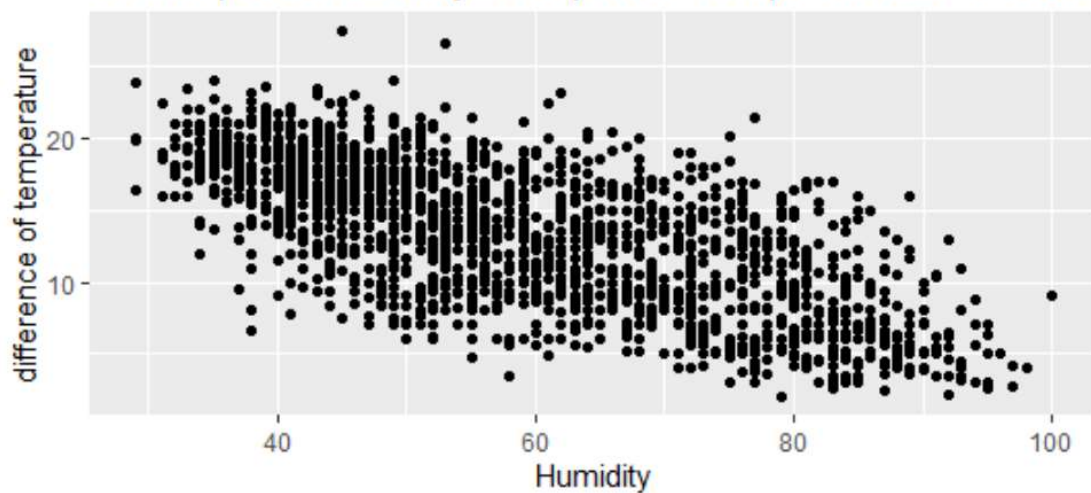


## Daily temperature difference in 1991-1995



After reiterating the steps done in homework 1. I tried to find the relationship between Difference of Temperature and Humidity.  
I firstly draw a scatter plot of all element.

## Scatter plot of Humidity in respect of Temperature Difference



It seems there is a linearly relationship between the two variables.

Then I test whether the two variable is normally distributed.

The results show they did:

```
> shapiro.test(diff9195)

      shapiro-wilk normality test

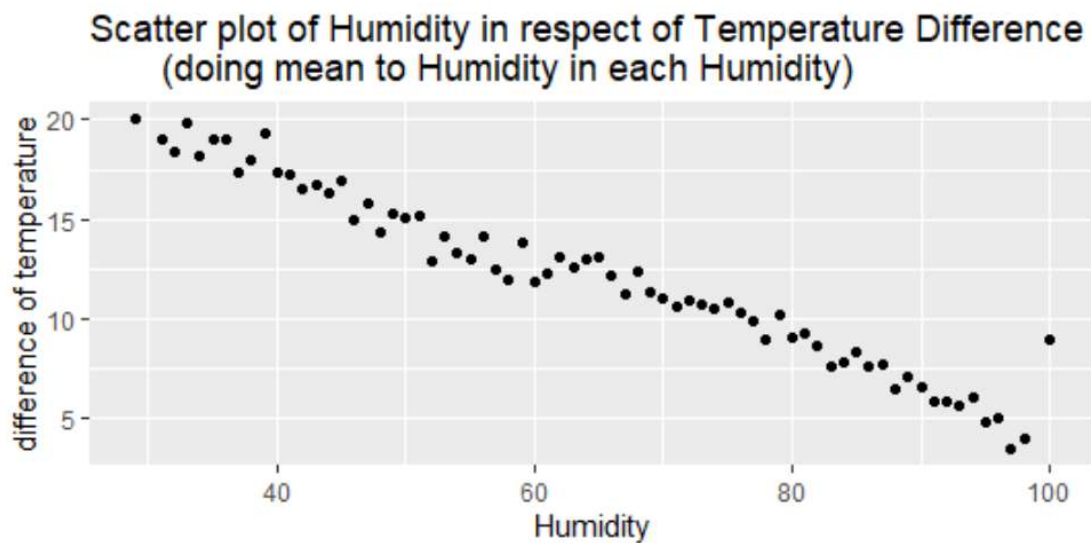
data:  diff9195
W = 0.97827, p-value = 5.204e-16

> ##p-value < 2.2e-16,is normally distributed
> class(H)
[1] "character"
> shapiro.test(as.numeric(H))

      shapiro-wilk normality test

data:  as.numeric(H)
W = 0.96192, p-value < 2.2e-16
```

The two p-values are both in a tiny amount. The two sets both obey the normal distribution.



```
> summary(hdregression)
```

Call:  
lm(formula = HD\$diffavg ~ HD\$nH, data = HD)

Residuals:

Min	1Q	Median	3Q	Max
-1.8678	-0.5157	0.0702	0.4196	4.2858

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	25.617606	0.360754	71.01	<2e-16 ***
HD\$nH	-0.209034	0.005336	-39.17	<2e-16 ***

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

Residual standard error: 0.9043 on 68 degrees of freedom  
Multiple R-squared: 0.9576, Adjusted R-squared: 0.9569  
F-statistic: 1535 on 1 and 68 DF, p-value: < 2.2e-16

The humidity and the temperature difference can well fit linear regression.