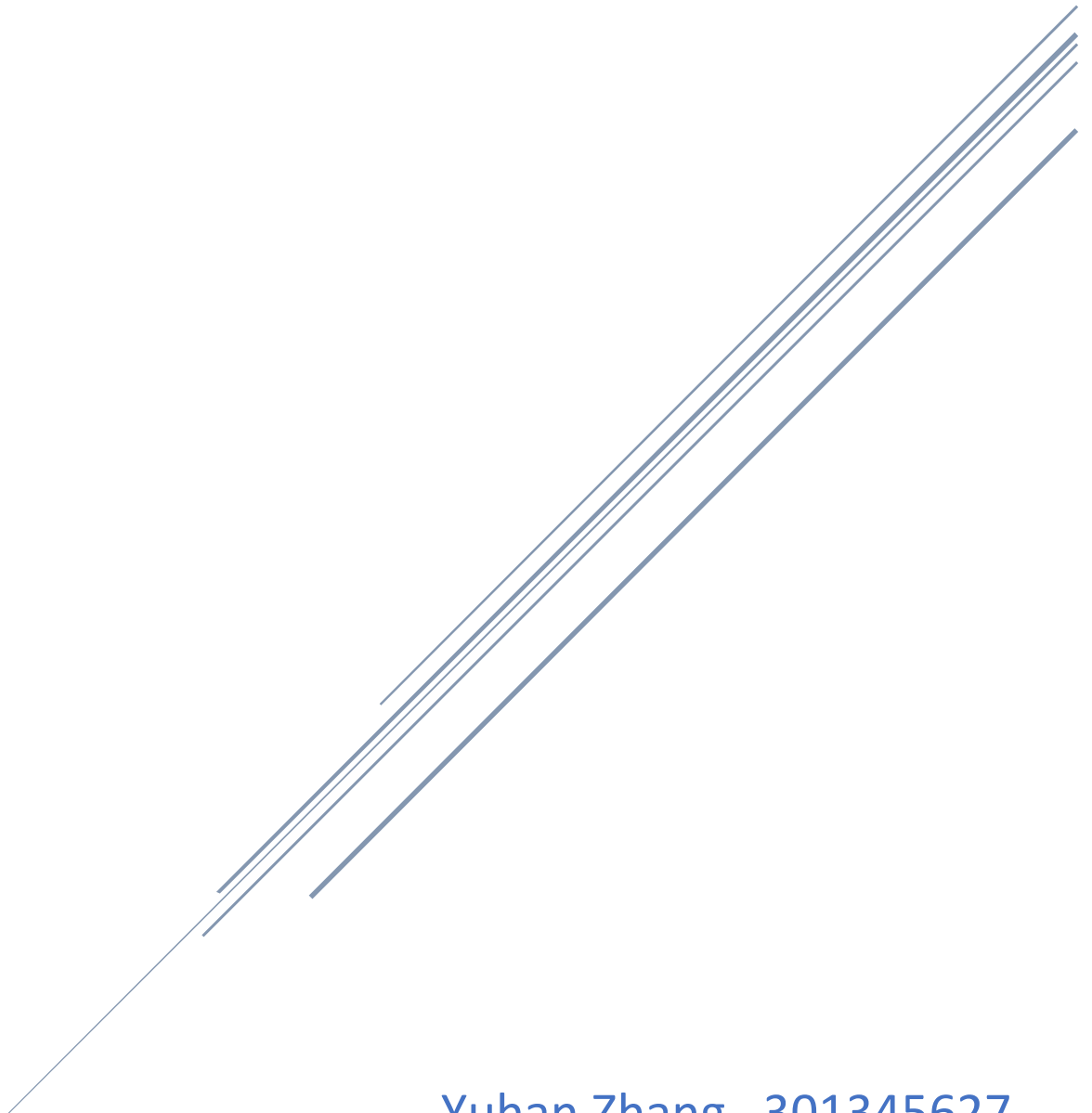


# STAT350 ASSIGNMENT1



Yuhan Zhang 301345627

## Q1 and Q2

1.  $n$  indept pairs of  $(x_i, y_i)$ . Show  $\sum_{i=1}^n x_i e_i = 0$ .

$$S(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$$\frac{dS}{d\hat{\beta}_1} = \sum_{i=1}^n 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-x_i)$$

$$= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i$$

$$= \sum_{i=1}^n e_i \cdot x_i = 0$$

2. Show  $\sum_{i=1}^n \hat{y}_i e_i = 0$ .

$$\frac{dS}{d\hat{\beta}_0} = \sum_{i=1}^n 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-1)$$

$$= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)$$

$$= \sum_{i=1}^n y_i - \hat{y}_i = \boxed{\sum e_i = 0}$$

$$\sum_{i=1}^n \hat{y}_i e_i = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i) e_i$$

$$= \hat{\beta}_0 \sum_{i=1}^n e_i + \hat{\beta}_1 \sum_{i=1}^n x_i e_i$$

$$\therefore \sum_{i=1}^n e_i = 0 \quad \text{and} \quad \sum_{i=1}^n e_i x_i = 0 \quad (\text{from Q1})$$

$$\therefore \sum_{i=1}^n \hat{y}_i e_i = 0$$

### Q3

3.  $y = \beta x + \epsilon$ .  $\epsilon \sim N(0, \sigma^2)$

$\therefore \epsilon \sim N(0, \sigma^2)$   $\therefore \epsilon$  is normal distribution.

$$\therefore f(\epsilon) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\beta x)^2}{2\sigma^2}}$$

$$L = f(\epsilon_1) \cdot f(\epsilon_2) \cdots f(\epsilon_i)$$

$$= \prod_{i=1}^n f(\epsilon_i)$$

$$= (2\pi\sigma^2)^{-n/2} e^{-\frac{\sum (y_i - \beta x_i)^2}{2\sigma^2}}$$

$$l = \ln L = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{\sum (y_i - \beta x_i)^2}{2\sigma^2}$$

$$0 = \frac{dl}{d\beta} = -\frac{1}{\sigma^2} \sum (y_i - \beta x_i)^2$$

$$\therefore \frac{d(\sum (y_i - \beta x_i)^2)}{d\beta} = 0$$

$\Rightarrow$  Same as OLS

$$0 = \frac{dl}{d\sigma^2} = -\frac{n}{2} \left( \frac{2\pi}{2\pi\sigma^2} \right) + \frac{\sum (y_i - \beta x_i)^2}{2\sigma^4}$$

$$\therefore \frac{n}{2} \frac{1}{\sigma^2} = \frac{\sum (y_i - \beta x_i)^2}{2\sigma^4} = \frac{\sum e_i^2}{2\sigma^4}$$

$$\frac{6n}{2} = \frac{\sum e_i^2}{2}$$

$\hat{\sigma}^2 = \sum e_i^2 / n \sim \text{MLE}$  which is biased as the unbiased estimator is  $\hat{\sigma}^2 = \sum e_i^2 / (n-2)$ .

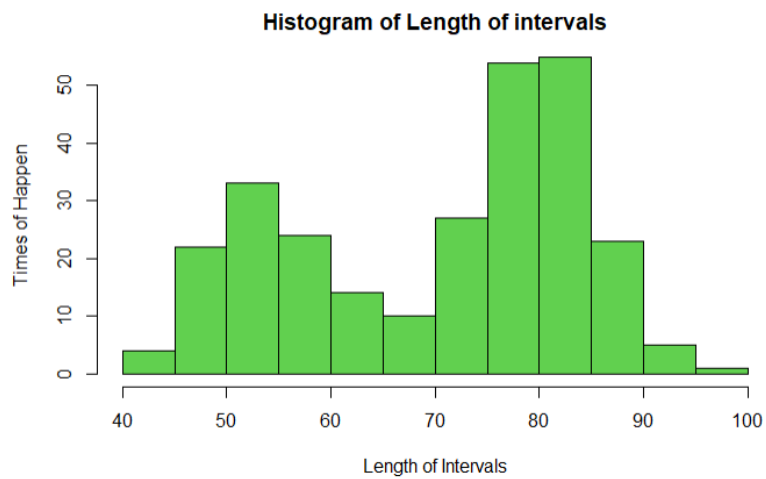
## Q4

```
2 #####Stat350 Homework1
3 {r}
4 mydata <- read.csv("geyser.csv",header=TRUE)
5 mydata
6
```

X <int>	eruptions <dbl>	waiting <int>
1	3.600	79
2	1.800	54
3	3.333	74
4	2.283	62
5	4.533	85
6	2.883	55
7	4.700	88
8	3.600	85
9	1.950	51
10	4.350	85

```
##a)
{r}
hist(mydata$waiting, main="Histogram of Length of intervals", xlab="Length of Intervals", ylab= "Times of Happen",
col=3)

```

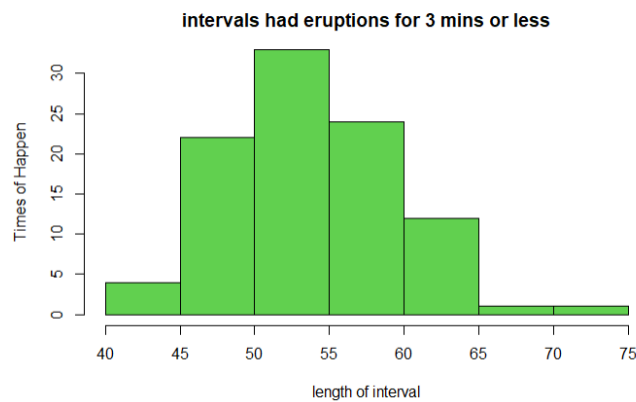
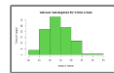


As there are two peaks, the data is not a normal distribution. Therefore, it is not faithful and is not reliable.

```

22- ##b)
23- {r}
24- mean(mydata$waiting)
25- sd(mydata$waiting)
26-
27- No, I would not use mean and standard deviation to describe the data as the data is not a normal distribution.
28-
29- ##c)
30- {r}
31- less_than_3 = df[mydata$eruptions <= 3,]
32- more_than_3 = df[mydata$eruptions > 3,]
33- hist(less_than_3$waiting, xlab="length of interval", ylab= "Times of Happen", main="intervals had eruptions for 3 mins
34- or less", col=3)
35- hist(more_than_3$waiting, xlab="length of interval", ylab= "Times of Happen", main="intervals had eruption for more
than 3 mins", col=7)

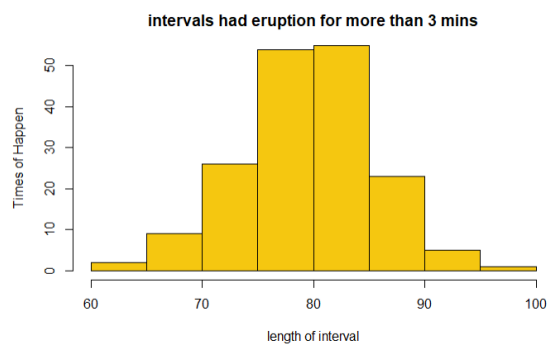
```



```

22- ##b)
23- {r}
24- mean(mydata$waiting)
25- sd(mydata$waiting)
26-
27- No, I would not use mean and standard deviation to describe the data as the data is not a normal distribution.
28-
29- ##c)
30- {r}
31- less_than_3 = df[mydata$eruptions <= 3,]
32- more_than_3 = df[mydata$eruptions > 3,]
33- hist(less_than_3$waiting, xlab="length of interval", ylab= "Times of Happen", main="intervals had eruptions for 3 mins
34- or less", col=3)
35- hist(more_than_3$waiting, xlab="length of interval", ylab= "Times of Happen", main="intervals had eruption for more
than 3 mins", col=7)

```





```

37-
38- {r}
39- mean(less_than_3$waiting, )
40- sd(less_than_3$waiting, )
41- mean(less_than_3$waiting) +c(-1,1)*sd(less_than_3$waiting) ##68%
42- mean(less_than_3$waiting) +2*c(-1,1)*sd(less_than_3$waiting) ##95%
43- mean(less_than_3$waiting) +3*c(-1,1)*sd(less_than_3$waiting) ##99.7%
44-

```

```

[1] 54.49485
[1] 5.840098
[1] 48.65475 60.33494
[1] 42.81465 66.17504
[1] 36.97455 72.01514

```

```

45- {r}
46- mean(more_than_3$waiting, )
47- sd(more_than_3$waiting, )
48- mean(more_than_3$waiting) +c(-1,1)*sd(more_than_3$waiting)
49- mean(more_than_3$waiting) +2*c(-1,1)*sd(more_than_3$waiting)
50- mean(more_than_3$waiting) +3*c(-1,1)*sd(more_than_3$waiting)
51-

```

```

[1] 79.98857
[1] 5.994239
[1] 73.99433 85.98281
[1] 68.00009 91.97705
[1] 62.00585 97.97129

```

```

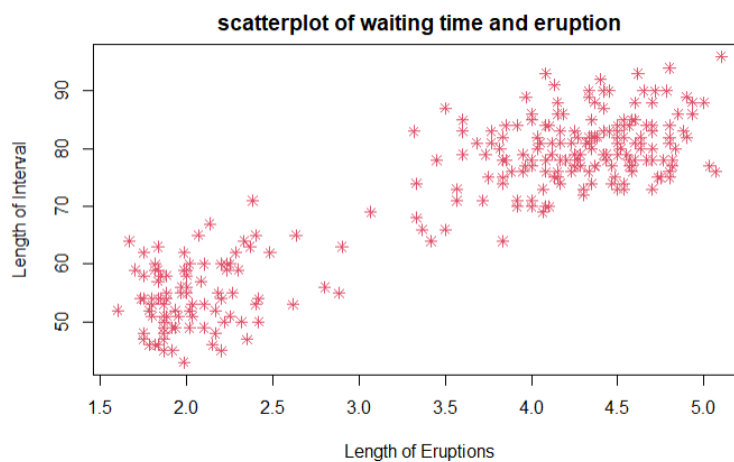
52
53
54

```

```

56
57- ###d)
58
59- {r}
60- plot(mydata$waiting-df$eruptions, main="scatterplot of waiting time and eruption", xlab="Length of
61- Eruptions",ylab="Length of Interval",col=2, pch=8)

```



```

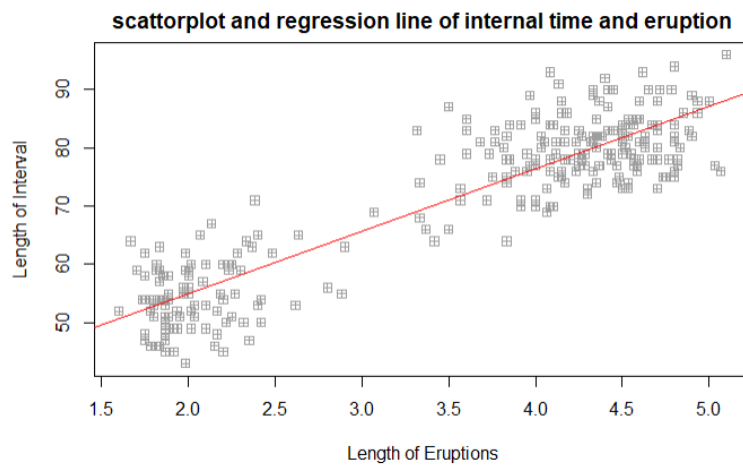
62- The length of internal between two eruption increases with the length of the precious eruption.
63- Yes, there is a positive linear relation between length of eruption and length of waiting interval.
64
65
66

```

```

75
76 ##e)
77 {r}
78 plot(mydata$waiting-mydata$eruptions, main="scatterplot and regression line of internal time and eruption",
      xlab="Length of Eruptions", ylab="Length of Interval", col=8, pch=12)
79 abline(regre_plot,col="red")
80

```



```

88
89 {r}
90 data.2 = data.frame(eruptions=2)
91 predict(regre_plot,data.2)
92

```

```

1
54.93368

```

```

93
94
95

```