

國立政治大學111學年度第二學期 期中R程式考題

Department: _____ ID: _____ Name: _____

Subject: **Regression Analysis (I)**

Date: 2023/04/20

Time: 11:00~12:00 (60 minutes)

60

注意事項:

1. 本次考題以R程式(Rgui或RStudio)方式作答，其他程式不允許。
2. 考試過程中可查詢書本、教學講義或上網，禁止利用messenger, IG, Line等等通訊軟體。
3. 禁止疑似作弊行為。
4. 本答案卷上請務必於 _____ 內複製「執行後的程式碼及結果(含圖形)」，於本答案卷貼上(Courier New, 10點字，白底黑字)，不能只有程式碼，不能只有報表。最後，將每小題之答案(不能只印出報表，要助教去找答案)，在小題最後以打字(英文)作答(Times New Roman, 12點字，白底黑字)。
5. 請依序註明題號: (1)a, (1)b, (2)a 等等。
6. 作答完請將此word檔存檔，檔名為「**學號-姓名-Regression-R-Midterm.docx**」(更改成自己「學號、姓名」)並上傳至<http://ftp.hmwu.idv.tw:8080/login.html?lang=tchinese>或點選教師網站首頁【作業考試上傳區】。
7. 帳號: reg111，密碼: 上課教室號碼，資料夾: 「**20230420-MidtermExam**」
8. 如果上傳網站出現「空白頁」，請將滑鼠移至「網址列」後，按「Enter」即可。若再不行，請換其它瀏覽器(IE/Edge/Firefox/Chrome)
9. 上傳檔案無法刪除，若要上傳更新檔，請於主檔名後加

「-2」，例如：「學號-姓名-Regression-R-Midterm-2.docx」。

Notes:

1. This is an Open Book exam; you are free to use any materials including laptop, tablet and internets.
2. Smart phone and the communication software/APP (e.g., Messenger, IG, LINE, WeChat,..) are prohibited.
3. Copy the R codes and the results from _____ and paste it to this answer sheet.
4. Change the file name of this answer sheet according to your ID and Full Name. Upload the answer sheet to <http://ftp.hmwu.idv.tw:8080/login.html?lang=tchinese>
5. Account: **rege111** , password: classroom number.

(1) Data file: Grade_Point_Average.csv

20% **Grade point average.** The director of admissions of a small college selected 120 students at random from the new freshman class in a study to determine whether a student's grade point average (GPA) at the end of the freshman year (Y) can be predicted from the ACT test score (X). The results of the study follow. Assume that first-order regression model (1.1) is appropriate.

i :	1	2	3	...	118	119	120
X_i :	21	14	28	...	28	16	28
Y_i :	3.897	3.885	3.778	...	3.914	1.860	2.948

- a. Obtain the least squares estimates of β_0 and β_1 , and state the estimated regression function.
- b. Plot the estimated regression function and the data. Does the estimated regression function appear to fit the data well?
- c. Obtain a point estimate of the mean freshman GPA for students with ACT test score $X = 30$.

Data input:

```
> df1 = read.table("Grade_Point_Average.csv")
```

a.

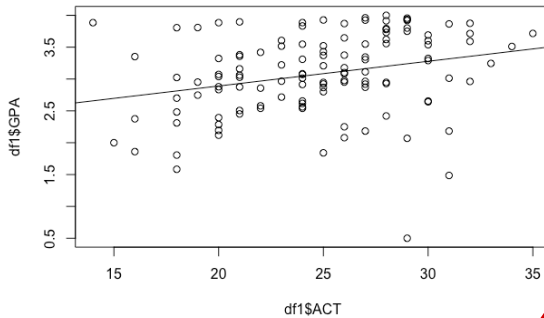
```
> model = lm(GPA~ACT,data = df1)
> model$coefficients #the least square estimates of beta0, beta1
(Intercept)          ACT
2.11404929  0.03882713
```

Hence, the least square estimates for beta0, beta1 is 1.037814e-14 and 1.000000e+00, respectively.

The estimated regression function is : $\hat{Y} = 1.037814e-14 + 1.000000e+00 * X$

b.

```
> plot(x = df1$ACT, y = df1$GPA)
> abline(model)
```



I think the model does not fit well since there are a lot of points relatively far from the regression line.

C.
 > predict(model, data.frame(ACT = 30))
 1
 3.278863

The point estimate of the mean freshman GPA for students with ACT=30 is 3.278863.

(2) Data file: [Grade_Point_Average.csv](#)

20% Refer to **Grade point average**

- a. Set up the ANOVA table.
- b. Conduct an F test of whether or not $\beta_1 = 0$. Control the α risk at .01. State the alternatives, decision rule, and conclusion.
- c. What is the absolute magnitude of the reduction in the variation of Y when X is introduced into the regression model?

Data input:

```
> df1 = read.table("Grade_Point_Average.csv")
```

a.

```
> model = lm(GPA~ACT,data = df1)
> anova(model)
Analysis of Variance Table
```

Response: GPA

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
ACT	1	3.588	3.5878	9.2402	0.002917 **
Residuals	118	45.818	0.3883		

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

b.

H_0 : β_1 is 0 v.s H_1 : β_1 is not 0

The test statistic is $F_{\text{star}} = \text{MSR}/\text{MSE}$. If $F_{\text{star}} > F(0.95,1,118)$, then reject H_0 .

Since $F_{\text{star}} = 9.2402 > F(0.95,1,118)$ and $p\text{-value} = 0.002917$, we reject H_0 . The conclusion is there is a significant linear relation between freshmen GPA and ACT scores.

c.

The absolute magnitude of the reduction in variation of Y when X is introduced is $\text{SSR} = 3.588$.

30% Refer to Grade point average

- Prepare a box plot for the ACT scores X_i . Are there any noteworthy features in this plot?
- Prepare a dot plot of the residuals. What information does this plot provide?
- Plot the residual e_i against the fitted values \hat{Y}_i . What departures from regression model (2.1) can be studied from this plot? What are your findings?
- Prepare a normal probability plot of the residuals. Also obtain the coefficient of correlation between the ordered residuals and their expected values under normality. Test the reasonableness of the normality assumption here using Table B.6 and $\alpha = .05$. What do you conclude?
- Conduct the Brown-Forsythe test to determine whether or not the error variance varies with the level of X . Divide the data into the two groups, $X < 26$, $X \geq 26$, and use $\alpha = .01$. State the decision rule and conclusion. Does your conclusion support your preliminary findings in part (c)?
- Information is given below for each student on two variables not included in the model, namely, intelligence test score (X_2) and high school class rank percentile (X_3). (Note that larger class rank percentiles indicate higher standing in the class, e.g., 1% is near the bottom of the class and 99% is near the top of the class.) Plot the residuals against X_2 and X_3 on separate graphs to ascertain whether the model can be improved by including either of these variables. What do you conclude?

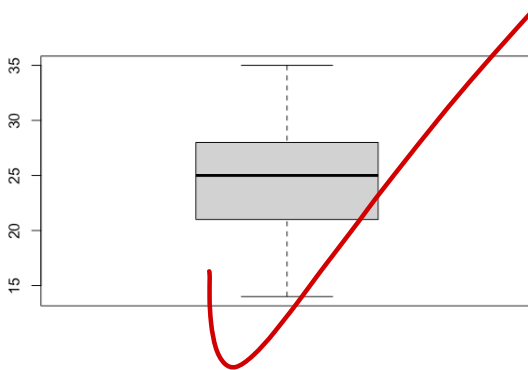
i :	1	2	3	...	118	119	120
X_2 :	122	132	119	...	140	111	110
X_3 :	99	71	75	...	97	65	85

3.

```
> df2 = read.table("Grade_Point_Average_X.csv")
```

a.

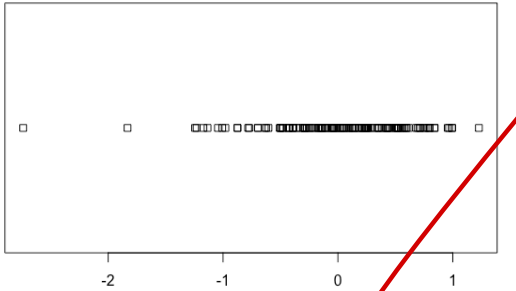
```
> boxplot(df2$ACT)
```



The box plot shows that the distribution of ACT scores are quite symmetric.

b.

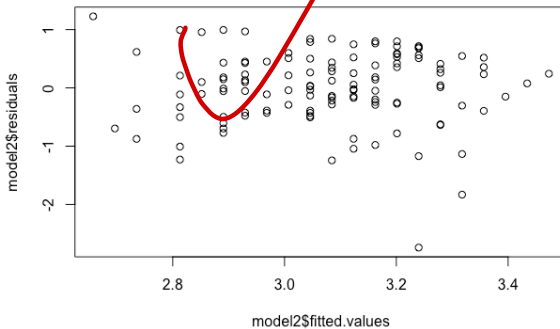
```
model2 = lm(GPA~ACT, data = df2)  
stripchart(model2$residuals, method = "stack")
```



The dot plot shows that residual are quite dense around 0.

c.

```
> plot(model2$fitted.values, model2$residuals)
```



The residual plot shows that the variance of residual is increasing when the fitted values increase. This matched the departure of constant variance from model assumption.

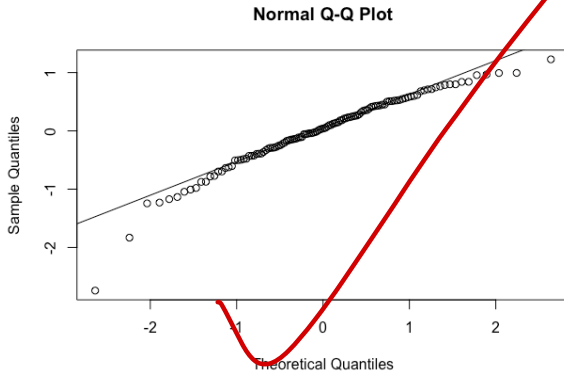
d.

```
> qqnorm(model2$residuals)
> qqline(model2$residuals)
>
> anova(model2)
Analysis of Variance Table
```

Response: GPA

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
ACT	1	3.588	3.5878	9.2402	0.002917 **
Residuals	118	45.818	0.3883		

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



Hence , $MSE = 0.3883$

```
> ranked_residual = rank(model2$residuals)
> MSE = 0.3883 #obtained by ANOVA table
> expected_e = sqrt(MSE)*qnorm(ranked_residual - 0.375/120+0.25)
警告訊息:
於 qnorm(ranked_residual - 0.375/120 + 0.25): 產生了 NaNs
>
> cor(model2$residuals,ranked_residual)
[1] 0.942459
```

The correlation between residuals and expected values are high, which indicates that the residuals might follow a normal distribution.

To conclude, by QQ-plot and correlation test, we can say that there is a significant proof that residuals follow a normal distribution.

(4) [Data file: Solution_concentration.csv](#)

30%

Solution concentration. A chemist studied the concentration of a solution (Y) over time (X). Fifteen identical solutions were prepared. The 15 solutions were randomly divided into five sets of three, and the five sets were measured, respectively, after 1, 3, 5, 7, and 9 hours. The results follow.

i :	1	2	3	...	13	14	15
X_i :	9	9	9	...	1	1	1
Y_i :	.07	.09	.08	...	2.84	2.57	3.10

- Prepare a scatter plot of the data. What transformation of Y might you try, using the prototype patterns in Figure 3.15 to achieve constant variance and linearity?
- Use the Box-Cox procedure and standardization (3.36) to find an appropriate power transformation. Evaluate SSE for $\lambda = -.2, -.1, 0, .1, .2$. What transformation of Y is suggested?
- Use the transformation $Y' = \log_{10} Y$ and obtain the estimated linear regression function for the transformed data.
- Plot the estimated regression line and the transformed data. Does the regression line appear to be a good fit to the transformed data?
- Obtain the residuals and plot them against the fitted values. Also prepare a normal probability plot. What do your plots show?
- Express the estimated regression function in the original units.