



「-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 internet.
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.

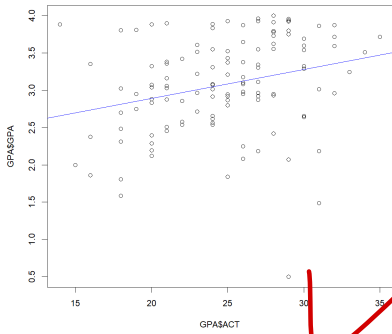
### (1) Data file: Grade\_Point\_Average.csv

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  $\beta_0$  and  $\beta_1$ , 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  $X = 30$ .

```
> GPA <- read.csv("Grade_Point_Average.csv")
> #a.
> GPA.lm <- lm(GPA~ACT, data=GPA)
> (b0 <- GPA.lm$coefficients[1])
(Intercept)
2.114049
> (b1 <- GPA.lm$coefficients[2])
ACT
0.03882713
> #b.
> plot(GPA$ACT, GPA$GPA)
> abline(GPA.lm, col="blue")
> #c.
> predict(GPA.lm, data.frame(ACT=30))
1
3.278863
```



(1)

a.

b.

The estimated regression function seems to fit the data well, since most of the points scatter around the regression line.

c.

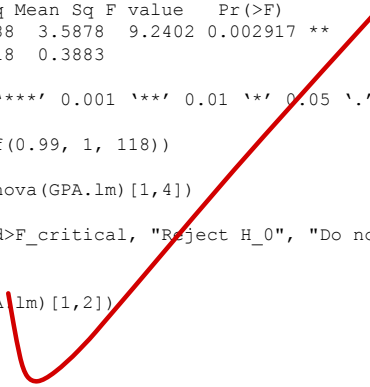
(2) Data file: Grade\_Point\_Average.csv

20% Refer to **Grade point average**

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

```
> #a.
> anova(GPA.lm)
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.
> (F_critical <- qf(0.99, 1, 118))
[1] 6.854641
> (F_observed <- anova(GPA.lm)[1,4])
[1] 9.240243
> ifelse(F_observed>F_critical, "Reject H_0", "Do not reject H_0")
[1] "Reject H_0"
> #c.
> (SSR <- anova(GPA.lm)[1,2])
[1] 3.587846
```



(2)

a.

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.

① versus

② Level of significance:

③ Test statistic:

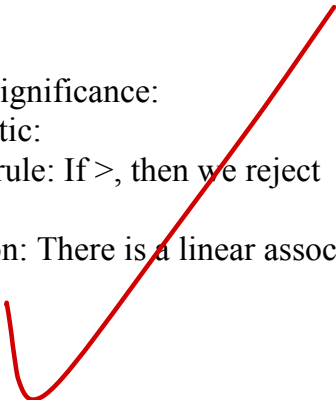
④ Decision rule: If  $>$ , then we reject

⑤ Decision:

⑥ Conclusion: There is a linear association between X and Y.

c.

SSR=3.5878



(3) Data file: [Grade\\_Point\\_Average\\_X.csv](#)

30% Refer to **Grade point average**

- a. Prepare a box plot for the ACT scores  $X_i$ . Are there any noteworthy features in this plot?
- b. Prepare a dot plot of the residuals. What information does this plot provide?
- c. 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?
- d. 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?
- e. 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)?
- f. 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

```
> #a.  
> boxplot(GPA$ACT, horizontal = TRUE)  
> #b.  
> dotchart(GPA.lm$residuals)  
> #c.  
> plot(GPA.lm$fitted.values, GPA.lm$residuals)  
> #d.  
> qqnorm(GPA.lm$residuals); qqline(GPA.lm$residuals)  
>  
> residuals.rank <- rank(GPA.lm$residuals)
```

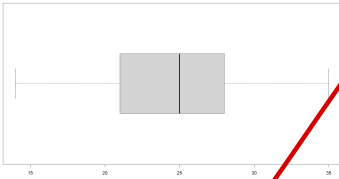
```

> MSE <- anova(GPA.lm)[2,3]
> residuals.expected <- sqrt(MSE)*qnorm((residuals.rank-0.375)/(nrow(GPA)+0.25))
> (corr <- cor(GPA.lm$residuals, residuals.expected))
[1] 0.9737275
> (B.6.critical <- 0.987)
[1] 0.987
> ifelse(corr>B.6.critical, "Reject H_0", "Do not reject H_0")
[1] "Do not reject H_0"
> #e.
> group <- ifelse(GPA$ACT<26, "A", "B")
> library(ALSM)
> bftest(GPA.lm, group)
      t.value  P.Value alpha  df
[1,] 0.8967448 0.371681 0.05 118

```

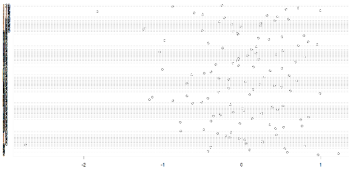
(3)

a.



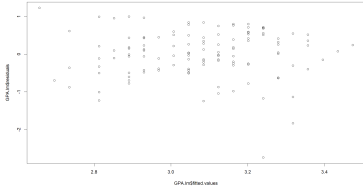
The data is skewed to the left.

b.



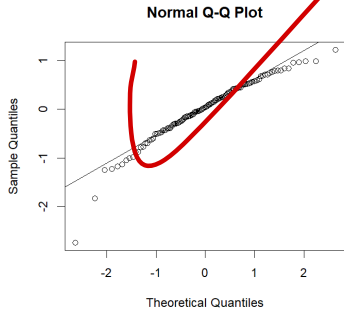
There are outliers.

c.



The error variance is fairly constant.

d.



The coefficient of correlation is 0.9737275.

$0.9737275 < 0.987$ , do not reject. Not normally distributed.

e.

①

② Level of significance:

③ Test statistic:

④ Decision rule: If  $>$ , then we reject

⑤ Decision:

⑥ Conclusion: The error variance is constant.



(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.

```
# (a)
plot(concentration$X, concentration$Y, main="3.16(a)")
abline(concentration.lm, col="blue", lwd=2)
## We need a transformation on Y, since the error variance is not constant and
## the linear regression seems not to fit the data.
## I might try Y'=ln(Y) or Y'=sqrt(-Y)
```

```
# (b)
# library(EnvStats)
# Y_bc <- boxcox(concentration$Y, optimize = TRUE)
# (lambda <- Y_bc$lambda)
boxcox(concentration.lm, lambda = seq(-2, 2, by = 0.1))
library(ALSM)
(SSE <- boxcox.sse(concentration$X, concentration$Y, l=seq(-.2, .2, by=.1)))
## Y'=ln(Y) is suggested.
```

```
# (c)
concentration.lm.log10Y <- lm(log10(Y)~X, data=concentration)
concentration.lm.log10Y$coefficients
```

```
## ^Y' = 0.6549-0.1954X
```

```
# (d)
```

```
plot(concentration$X, log10(concentration$Y), main="3.16(d)")
```

```
abline(concentration.lm.log10Y, col="blue", lwd=2)
```

```
## The regression line appears to be a good fit to the transformed data since
```

```
## the regression line pass through the center of data and the error variance
```

```
## is approximately constant.
```

```
# (e)
```

```
residuals <- concentration.lm.log10Y$residuals
```

```
fitted <- concentration.lm.log10Y$fitted.values
```

```
par(mfrow = c(1, 2))
```

```
plot(fitted, residuals, main="Residuals against the fitted values")
```

```
qqnorm(residuals); qqline(residuals)
```

```
## The error variance is approximately constant and normally distributed.
```

```
# (f)
```

```
# log10(^Y)=0.6549-0.1954X
```

```
# => ^Y = 10^(0.6549-0.1954X)
```

```
4
```

*output?*