

UNIVERSITY OF TORONTO

Faculty of Arts and Science

JUNE EXAMINATION 2004

STA 302 H1F / STA 1001 H1F

Duration - 3 hours

Aids Allowed: Calculator

NAME: _____

STUDENT NUMBER: _____

- There are 17 pages including this page.
- The last page is a table of formulae that may be useful.
- Tables of the t distribution can be found on page 15 and tables of the F distribution can be found on page 16.
- Total marks: 75

1abcde	1fg	2ab	2cd	3	4

5a	5b	5cd	5ef	6ab	6cde

1. Olympic gold medal performances in track and field improve over time. A regression was run with dependent variable `longjump`, the winning distance in the long jump (in inches), and independent variable `year`, the year the Olympics was held after 1900 (counting 1900 as year 0). Data from the Olympics held from 1900 through 1984 were used (some Olympics were missed during the World Wars). Here are the data:

<code>year</code>	0	4	8	12	20	24	28	32	36	48
<code>longjump</code>	282.9	289.0	294.5	299.3	281.5	293.1	304.8	300.8	317.3	308.0
<code>year</code>	52	56	60	64	68	72	76	80	84	
<code>longjump</code>	298.0	308.3	319.8	317.8	350.5	324.5	328.5	336.3	336.3	

Some output from SAS is given below.

The REG Procedure
Descriptive Statistics

Variable	Sum	Mean	Uncorrected SS	Variance	Standard Deviation
Intercept	19.00000	1.00000	19.00000	0	0
<code>year</code>	824.00000	43.36842	49184	747.13450	27.33376
<code>longjump</code>	5890.81250	310.04276	1833077	370.73440	19.25446

The REG Procedure
Model: MODEL1
Dependent Variable: `longjump`

Analysis of Variance						
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	
Model	1	5054.88650	5054.88650	53.10	<.0001	
Error	17	1618.33267	95.19604			
Corrected Total	18	6673.21916				

Root MSE	9.75685	R-Square	0.7575
Dependent Mean	310.04276	Adj R-Sq	0.7432
Coeff Var	3.14694		

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	
Intercept	1	283.45427	4.28064	66.22	<.0001	
<code>year</code>	1	0.61308	0.08413	7.29	<.0001	

Questions based on this output are on the next 2 pages.

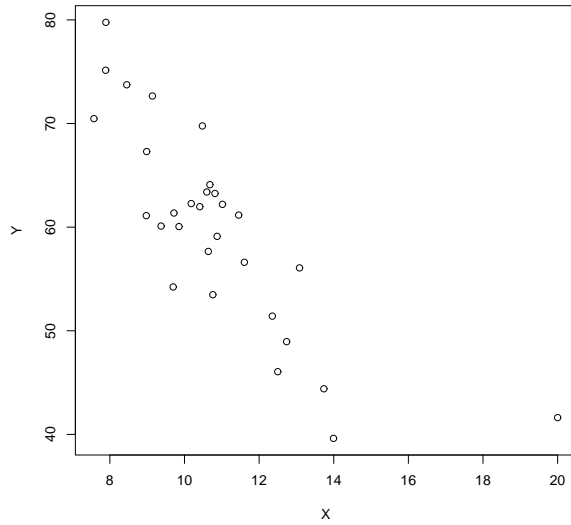
- (a) (2 marks) Estimate the mean change in the winning long jump distance from one Olympics to the next, assuming that the Olympics are held every 4 years.
- (b) (1 mark) Estimate the variance in winning long jump distance that is not explained by the year the Olympics were held.
- (c) (1 mark) Estimate the percent of total variability in winning long jump distance that is explained by the year the Olympics were held.
- (d) (1 mark) Estimate the correlation between winning long jump distance and the year the Olympics were held.
- (e) (2 marks) What is the observed value of the test statistic for the test $H_0 : \beta_1 = 0$ versus $H_a : \beta_1 > 0$? What is the p -value for this test?

(f) (3 marks) Construct simultaneous 99% confidence intervals for the slope and intercept of the regression line.

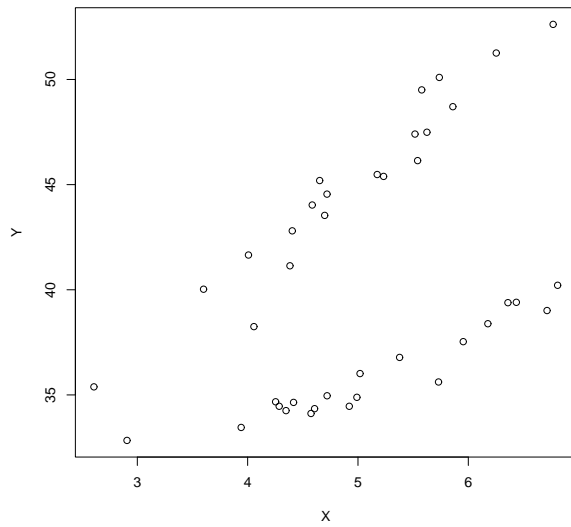
(g) (5 marks) It has been suggested that the Mexico City Olympics in 1968 saw unusually good track and field performances, possibly because of the high altitude. Construct an appropriate 95% interval for the predicted winning long jump distance in 1968. Do the data support these suggestions? Explain.

2. (8 marks (2 marks each)) Each of the following plots is a scatterplot of a dependent variable versus an independent variable. We wish to study further the relationship between the two variables. Indicate an appropriate linear regression model based on the plot.

(a)

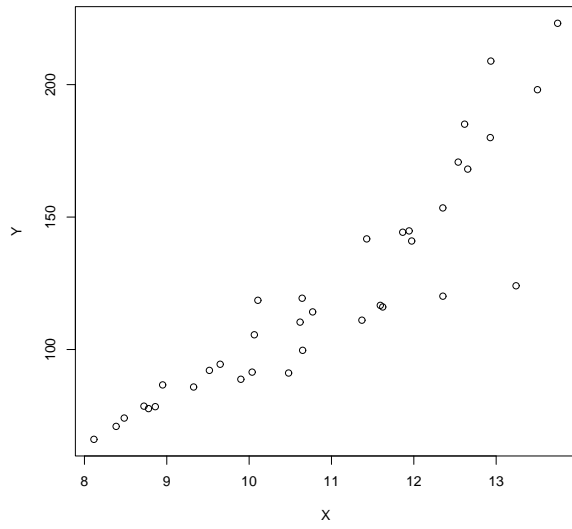


(b)

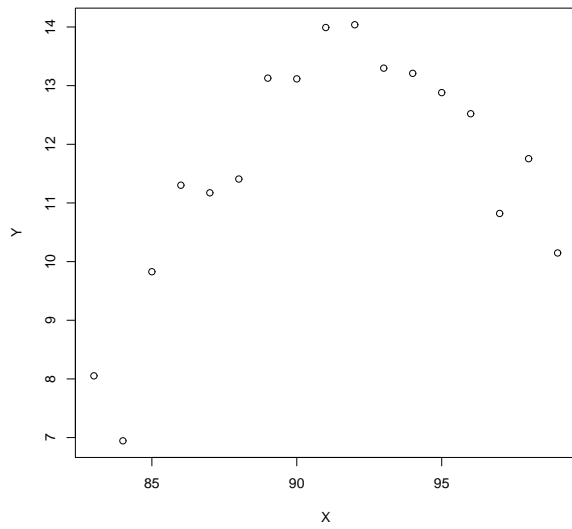


(This question continues on the next page.)

(c)



(d)



3. (3 marks) Show, in the case of simple linear regression, that the fitted line passes through the point (\bar{X}, \bar{Y}) .

4. (a) (5 marks) A multiple linear regression model is to be constructed to examine the relationship between a response variable Y and 3 predictor variables X_1 , X_2 , and X_3 . Suppose measurements of these 4 variables have been taken on n items. State the multiple linear regression model in matrix terms, defining all of your matrices, including the standard assumptions.

(b) (3 marks) Derive the expression for the covariance matrix of the least squares estimators of the model coefficients of your model in part (a). (I.e., Show $\text{Cov}(\mathbf{b}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.)

5. The data considered in the analysis below are observations on the acceleration (**accel**) of different vehicles along with their weight-to-horsepower ratio (**whp**), the speed at which they were travelling (**speed**), and the grade of the road (**grade**) which takes values 0, 2, and 6 (a value of 0 indicates the road was horizontal). There are 50 observations in the data set.

(a) (5 marks) The first model tried for these data was

$$\text{accel} = \beta_0 + \beta_1\text{whp} + \beta_2\text{speed} + \beta_3\text{grade} + \epsilon$$

Some output from SAS is given below.

The REG Procedure					
Model: MODEL1					
Dependent Variable: accel					
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	3	164.99430	54.99810	25.45	<.0001
Error	46	99.41390	2.16117		
Corrected Total	49	264.40820			
Root MSE		1.47009	R-Square	0.6240	
Dependent Mean		2.60600	Adj R-Sq	0.5995	
Coeff Var		56.41183			
Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	7.19950	0.60087	11.98	<.0001
whp	1	-0.01838	0.00269	-6.83	<.0001
speed	1	-0.09347	0.01307	-7.15	<.0001
grade	1	-0.15548	0.09040	-1.72	0.0922

What does this output tell you about the ability of the vehicles to accelerate under various conditions? Your answer should explain the affects of each of **whp**, **speed**, and **grade** on acceleration.

(b) (3 marks) The next model fit was

$$\text{accel} = \beta_0 + \beta_1\text{whp} + \beta_2\text{speed} + \beta_3\text{grade0} + \beta_4\text{grade2} + \epsilon$$

where $\text{grade0} = 1$ if grade is 0 and is 0 otherwise, and $\text{grade2} = 1$ if grade is 2 and is 0 otherwise. Some output from SAS for this model follows.

The REG Procedure
Dependent Variable: accel

Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	4	165.00484	41.25121	18.67	<.0001
Error	45	99.40336	2.20896		
Corrected Total	49	264.40820			

Root MSE	1.48626	R-Square	0.6241
Dependent Mean	2.60600	Adj R-Sq	0.5906
Coeff Var	57.03217		

Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Type I SS
Intercept	1	6.25693	0.61730	10.14	<.0001	339.56180
whp	1	-0.01839	0.00272	-6.76	<.0001	52.93850
speed	1	-0.09345	0.01322	-7.07	<.0001	105.66329
grade0	1	0.93165	0.54864	1.70	0.0964	3.48062
grade2	1	0.65181	0.56668	1.15	0.2561	2.92243

What is the purpose of using the two variables grade0 and grade2 ?

(c) (6 marks) Is the evidence for a grade effect stronger or weaker from the model in part (b) than the model in part (a)? Support your answer with appropriate statistical tests.

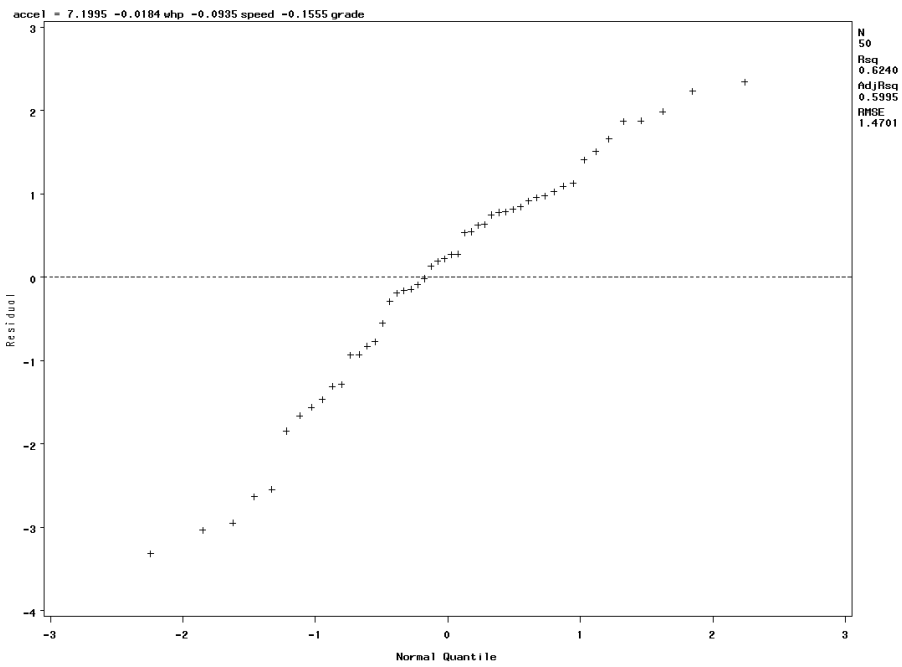
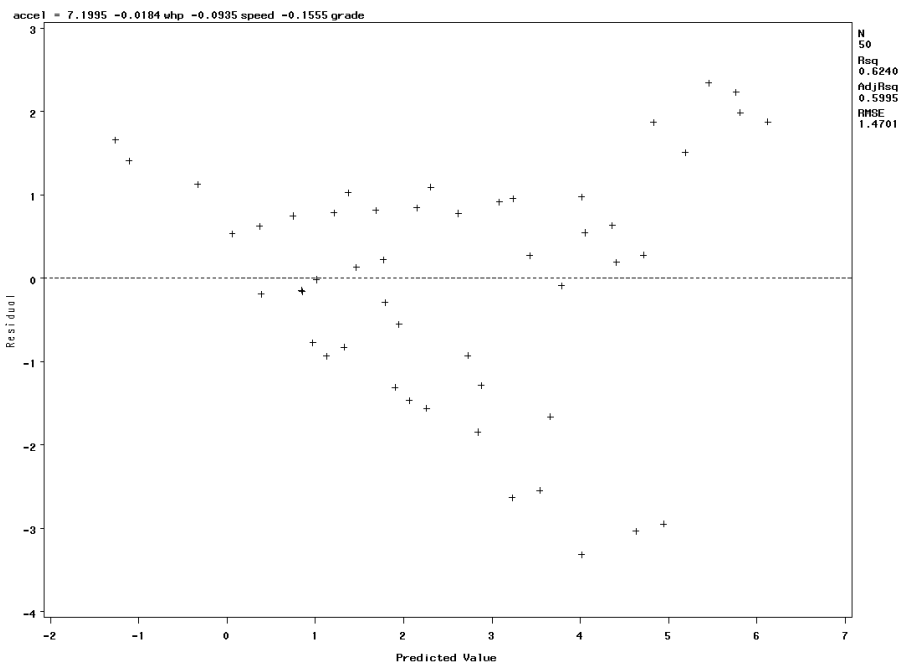
(d) (2 marks) Note that R^2 is almost identical for the models fit in parts (a) and (b). Choose another statistic that is useful for choosing between the two models and indicate which model is preferred.

(e) (3 marks) Another model that could be fit is

$$\text{accel} = \beta_0 + \beta_1 \text{whp} + \beta_2 \text{speed} + \beta_3 \text{grade} + \beta_4 \text{whp_speed} + \epsilon$$

where $\text{whp_speed} = \text{whp} * \text{speed}$. What additional information could be obtained from this model and how would you assess whether or not it is statistically important?

(f) (5 marks) A plot of the residuals versus the predicted values and a normal probability plot of the residuals for the regression in part (a) are on the next page. What additional information do these plots give?



6. (a) (6 marks) A simple linear regression is carried out and one point is determined to be an outlier but not an influential point. It is removed from the data and the regression line is fit to the remaining data. Which of the following quantities will differ by a substantial amount between the two regressions? If there is a substantial difference, indicate in which regression it is larger. Indicate if you don't have enough information to say.

i. the slope

ii. the mean square error

iii. R^2

(b) (3 marks) Explain the difference between an outlier and an influential point.

- (c) (2 marks) If you could choose the values of X at which to collect data before performing a simple linear regression analysis, would you prefer that $\sum_{i=1}^n (X_i - \bar{X})^2$ be large or small? Explain.
- (d) (2 marks) An investigator wishes to use multiple regression to predict a variable, Y , from two other variables, X_1 and X_2 . She is also interested in the quantity that is the sum of X_1 and X_2 and includes a third predictor variable in her model, $X_3 = X_1 + X_2$. What problems might she encounter?
- (e) (4 marks) In a study of infant mortality, a regression model was constructed using birth weight (which is a good indicator of the baby's likelihood of survival) as a dependent variable and several independent variables, including the age of the mother, whether the mother smoked or took drugs during pregnancy, the amount of medical attention she had, her income, etc. The R^2 was 11%, but the t -test provided by SAS for the coefficient of each predictor variable had p -value less than 0.01. An obstetrician has asked you to explain the significance of the study as it relates to her practice. What would you say to her?

Simple regression formulae

$$b_1 = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2}$$

$$b_0 = \bar{Y} - b_1 \bar{X}$$

$$\text{Var}(b_1) = \frac{\sigma^2}{\sum(X_i - \bar{X})^2}$$

$$\text{Var}(b_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum(X_i - \bar{X})^2} \right)$$

$$\text{Cov}(b_0, b_1) = -\frac{\sigma^2 \bar{X}}{\sum(X_i - \bar{X})^2}$$

$$\text{SSTO} = \sum(Y_i - \bar{Y})^2$$

$$\text{SSE} = \sum(Y_i - \hat{Y}_i)^2$$

$$\text{SSR} = b_1^2 \sum(X_i - \bar{X})^2 = \sum(\hat{Y}_i - \bar{Y})^2$$

$$\sigma^2\{\hat{Y}_h\} = \text{Var}(\hat{Y}_h)$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right)$$

$$\sigma^2\{\text{pred}\} = \text{Var}(Y_h - \hat{Y}_h)$$

$$= \sigma^2 \left(1 + \frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right)$$

$$\hat{X}_h \pm \frac{t_{n-2, 1-\alpha/2}}{|b_1|} * \text{appropriate s.e.}$$

(valid approximation if $\frac{t^2 s^2}{b_1^2 \sum(X_i - \bar{X})^2}$ is small)

Working-Hotelling coefficient:

$$W = \sqrt{2F_{2, n-2; 1-\alpha}}$$

Regression in matrix terms

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= \text{E}[(\mathbf{X} - \text{E}\mathbf{X})(\mathbf{X} - \text{E}\mathbf{X})'] \\ &= \text{E}(\mathbf{X}\mathbf{X}') - (\text{E}\mathbf{X})(\text{E}\mathbf{X})' \end{aligned}$$

$$\text{Cov}(\mathbf{A}\mathbf{X}) = \mathbf{A}\text{Cov}(\mathbf{X})\mathbf{A}'$$

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

$$\text{Cov}(\mathbf{b}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$$

$$\hat{\mathbf{Y}} = \mathbf{X}\mathbf{b} = \mathbf{H}\mathbf{Y}$$

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$$

$$\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

$$\text{SSR} = \mathbf{Y}'(\mathbf{H} - \frac{1}{n}\mathbf{J})\mathbf{Y}$$

$$\text{SSE} = \mathbf{Y}'(\mathbf{I} - \mathbf{H})\mathbf{Y}$$

$$\text{SSTO} = \mathbf{Y}'(\mathbf{I} - \frac{1}{n}\mathbf{J})\mathbf{Y}$$

$$\begin{aligned} \sigma^2\{\hat{Y}_h\} &= \text{Var}(\hat{Y}_h) \\ &= \sigma^2 \mathbf{X}'_h (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_h \end{aligned}$$

$$\begin{aligned} \sigma^2\{\text{pred}\} &= \text{Var}(Y_h - \hat{Y}_h) \\ &= \sigma^2(1 + \mathbf{X}'_h (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_h) \end{aligned}$$

$$R^2_{\text{adj}} = 1 - (n-1) \frac{\text{MSE}}{\text{SSTO}}$$