

Name:

Student Number:

## STAT 350: Summer Semester 2008

### Midterm 1: Solutions

9 June 2008

Instructor: Richard Lockhart

**Instructions:** This is an open book test. You may use notes, text, other books and a calculator. Your presentations of statistical analysis will be marked for clarity of explanation. I expect you to explain what assumptions you are making and to comment if those assumptions seem unreasonable. In general you need not finish doing arithmetic; I will be satisfied if your answers contain things like

$$27 \pm 1.96\sqrt{247.5/11},$$

but I have to be absolutely convinced you know what arithmetic to do! I want the answers written on the paper. The exam is out of 25.

1. The following story is taken from a data library called DASL: "In 1929, Edwin Hubble investigated the relationship between distance of a galaxy from the earth and the velocity with which it appears to be receding. Galaxies appear to be moving away from us no matter which direction we look." His paper records the distance (in megaparsecs) measured to 24 nebulae which lie outside our galaxy. For each of these nebulae Hubble also measured the velocity of the nebulae away from the earth (negative numbers mean the nebula is moving towards the earth).

Here are the data:

Distance	0.032	0.034	0.214	0.263	0.275	0.275	0.45	0.5
Velocity	170	290	-130	-70	-185	-220	200	290
Distance	0.5	0.63	0.8	0.9	0.9	0.9	0.9	1
Velocity	270	200	300	-30	650	150	500	920
Distance	1.1	1.1	1.4	1.7	2	2	2	2
Velocity	450	500	500	960	500	850	800	1090

In appendix 1 I fit the model:

$$\text{Velocity} = \beta_0 + \beta_1 \text{Distance} + \text{Error}.$$

I give SAS code and output, R code and output and corresponding JMP output. Use that output to answer the following questions:

- (a) Could the true value of the intercept term be 0? (Hubble's model called for the intercept to be 0.) [3 marks]

**Solution:** *We are to test  $H_o : \beta_0 = 0$ . The relevant  $t$  statistic given in the output is  $-0.489$ . The corresponding  $P$ -value is  $0.63$  which is not at all small. Thus we conclude that the intercept might well be 0; in formal terms this null hypothesis would be accepted at any common level of significance.*

- (b) Give a 95% confidence interval for the coefficient of Distance in our model. This coefficient is called the Hubble constant. [3 marks]

**Solution:** *The estimate of  $\beta_1$  is  $454.16$ ; the corresponding estimated standard error is  $75.24$ . There are 22 degrees of freedom for error here so we need a critical point from the  $t$ -tables with 22 degrees of freedom. I get  $2.07$  in my tables. So my confidence interval is*

$$454.16 \pm 2.07 \times 75.24.$$

- (c) The matrix  $(X^T X)^{-1}$  for this model and this data set is given in Appendix 2. Use this to give a 90% confidence interval for the mean recession velocity of all nebulae whose distance from earth is 1 megaparsec. [6 marks]

**Solution:** We are trying to get a confidence interval for  $\beta_0 + \beta_1$ . We estimate this using  $\hat{\beta}_0 + \hat{\beta}_1 = 454.16 - 40.78$ . The associated standard error is

$$\sigma \sqrt{x^t (X^T X)^{-1} x}$$

where  $x^t = [1 \ 1]$ . Our estimate of  $\sigma$  is 232.9. This gives an estimated standard error of

$$ESE = 232.9 \sqrt{0.1283 + 2(-0.0951) + 0.1043}$$

and the required interval is

$$454.16 - 40.78 \pm 1.717 \times ESE$$

You didn't need to work out the numbers more than that but it's obviously ok if you did.

- (d) When I fit a cubic polynomial model

$$V = \beta_0 + \beta_1 D + \beta_2 D^2 + \beta_3 D^3 + \epsilon$$

to this data I find the error sum of squares is 1062087 whereas the error sum of squares in the straight line model above is 1193442. What is the value of the  $F$  statistic for testing the hypothesis that  $\beta_2 = \beta_3 = 0$  and what are the relevant degrees of freedom? [3 marks]

**Solution:** The  $F$  statistic is a ratio of Extra Sum of Squares over its degrees of freedom to mean squared error in the full model. Thus

$$F = \frac{(1193442 - 1062087)/2}{1062087/(24 - 4)}$$

2. Consider the sand / fibre / plaster hardness example from class. In that example there were 18 batches of plaster – 2 batches made from each combination of 3 sand contents  $S$  and 3 fibre contents,  $F$ . I regressed Hardness on  $S, S^2, F, F^2$  and  $SF$ .

(a) If I now regressed on  $S, S^2, F, F^2, SF, S^3, S^2F, SF^2$  and  $F^3$  would the error sum of squares:

- i. go up,
- ii. go down,
- iii. stay the same
- iv. or is it not possible to tell without further information?

[1 mark for the answer and 2 for the explanation]

**Solution:** *Adding columns to a design matrix always makes the error sum of squares go down. When you estimate the beta vector for the model with extra columns one possible choice is the estimate for the original model with all the new components set to 0. So you can always get the same fitted value as in the restricted model. But you can also wiggle the new coefficients and the old ones. Since the least squares estimates make the error sum of squares as small as possible you must do at least as well as in the restricted model and unless the  $Y$  vector is very precisely set you will do better. In marking I am looking for this idea. I want the answer “go down” but I will accept ‘not possible to tell ... if you provide the explanation above.*

(b) Suppose instead I made another 18 batches of plaster again making 2 batches of each of the combinations of  $S$  and  $F$ . If I now used all 36 data points to fit the same model as in the start of this problem would the error sum of squares:

- (a) be higher than when using the original 18 data points,
- (b) be lower than when using the original 18 data points,
- (c) be the same
- (d) or is it not possible to tell without further information?

[1 mark for the answer and 2 for the explanation]

**Solution:** *The error sum of squares will go up. It is the sum of squared fitted errors and the sum of squares of fitted errors for the 18 data points cannot be improved by changing  $\hat{\beta}$  since that is the least squares estimates. Unless all the new data points have residuals exactly equal to 0 the ESS will go up. If you said ‘not possible to tell’ your explanation needs to be very good.*

- (c) Consider the situation described in part (b). When we fit the model to the original 18 data points we will get a certain standard error for the estimated coefficient of  $F$ . This standard error will have the form  $C\sigma$  for some constant  $C$ . When we fit the same model to all 36 data points will the the new standard error be higher or lower than  $C\sigma$ ? [1 mark for the answer, 1 for the explanation and a bonus mark for saying what the new standard error will be]

**Solution:** *The new matrix  $X^T X$  will simply be multiplied by 2 so the matrix  $(X^T X)^{-1}$  will be half as big. The standard error will now be  $C\sigma/\sqrt{2}$ . I will accept for full marks an explanation that says that collecting more data will increase the precision of all our estimates, thus decreasing the standard errors.*

3. In a study of a large number of sets of identical twins the correlation between IQ of the first born twin and IQ of the second born twin was found to be 0.85. Mean IQs were found to be similar in the two groups, as were standard deviations. When researchers picked out those pairs where the first born twin had an IQ over 130 (about 2% of the sets of twins) they found that the second born twin had a somewhat lower IQ than the first born twin on average. The opposite happened when they took sets of twins where the first born had an IQ below 70. Is this to be expected? Explain. [2 marks for the explanation]

**Solution:** *Yes this is to be expected. It is called the regression effect. The least squares line has the equation*

$$\frac{y - \bar{y}}{s_y} = r \frac{x - \bar{x}}{s_x}$$

*Since the correlation coefficient is less than 1 the predicted values are closer to the mean for the  $y$ s than the  $x$ s are and that is just what happened. No further explanation is required. Notice that since the means and standard deviations are the same being closer in standard units is the same as being closer in the original units.*

## Appendix 1

### R Code and output

```
> fit <- lm(Velocity~Distance,data=d)
```

Call:

```
lm(formula = Velocity ~ Distance, data = d)
```

Coefficients:

(Intercept)	Distance
-40.78	454.16

```
> summary(fit)
```

Call:

```
lm(formula = Velocity ~ Distance, data = d)
```

Residuals:

Min	1Q	Median	3Q	Max
-397.96	-158.10	-13.16	148.09	506.63

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-40.78	83.44	-0.489	0.63
Distance	454.16	75.24	6.036	4.48e-06 ***

---

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

Residual standard error: 232.9 on 22 degrees of freedom

Multiple R-Squared: 0.6235, Adjusted R-squared: 0.6064

F-statistic: 36.44 on 1 and 22 DF, p-value: 4.477e-06

```
> anova(fit.lin)
```

Analysis of Variance Table

Response: Velocity

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Distance	1	1976648	1976648	36.438	4.477e-06 ***
Residuals	22	1193442	54247		

SAS Code and output

```

data Hubble;
  infile 'Hubble.dat' firstobs=2;
  input Distance Velocity;
run;
proc glm;
  model Velocity = Distance;
run;

```

The GLM Procedure

```

Number of Observations Read      24
Number of Observations Used      24

```

Dependent Variable: Velocity

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	1976648.259	1976648.259	36.44	<.0001
Error	22	1193442.366	54247.380		

Corr Tot 23 3170090.625

R-Square	Coeff Var	Root MSE	Velocity Mean
0.623531	62.42162	232.9107	373.1250

Source	DF	Type I SS	Mean Square	F Value	Pr > F
Distance	1	1976648.259	1976648.259	36.44	<.0001

Parameter	Estimate	Standard Error	t Value	Pr >  t
Intercept	-40.7836491	83.43886994	-0.49	0.6298
Distance	454.1584409	75.23710535	6.04	<.0001



# JMP output

The next 3 pages show output from JMP when Velocity is taken as Y and Distance as X:

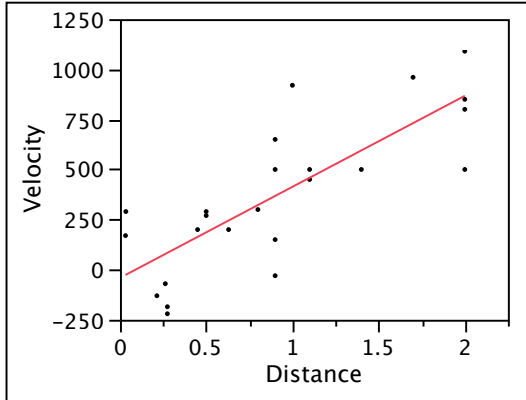
Hubble: Fit Least Squares

Page 1 of 3

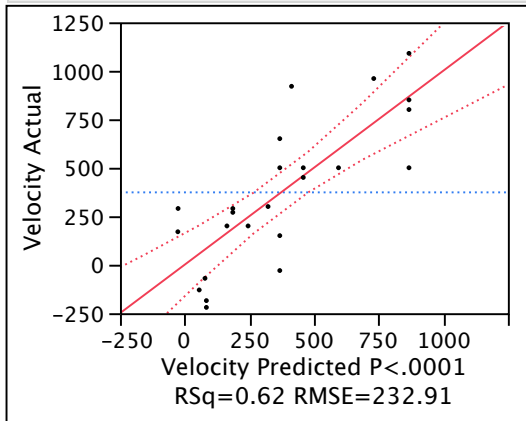
## Response Velocity

### Whole Model

#### Regression Plot



#### Actual by Predicted Plot



#### Summary of Fit

RSquare	0.623531
RSquare Adj	0.606418
Root Mean Square Error	232.9107
Mean of Response	373.125
Observations (or Sum Wgts)	24

## Hubble: Fit Least Squares

### Response Velocity

#### Whole Model

##### Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	1	1976648.3	1976648	36.4377
Error	22	1193442.4	54247	<b>Prob &gt; F</b>
C. Total	23	3170090.6		<.0001*

##### Lack Of Fit

Source	DF	Sum of Squares	Mean Square	F Ratio
Lack Of Fit	13	722504.9	55577.3	<b>Prob &gt; F</b>
Pure Error	9	470937.5	52326.4	0.4768
Total Error	22	1193442.4		<b>Max RSq</b> 0.8514

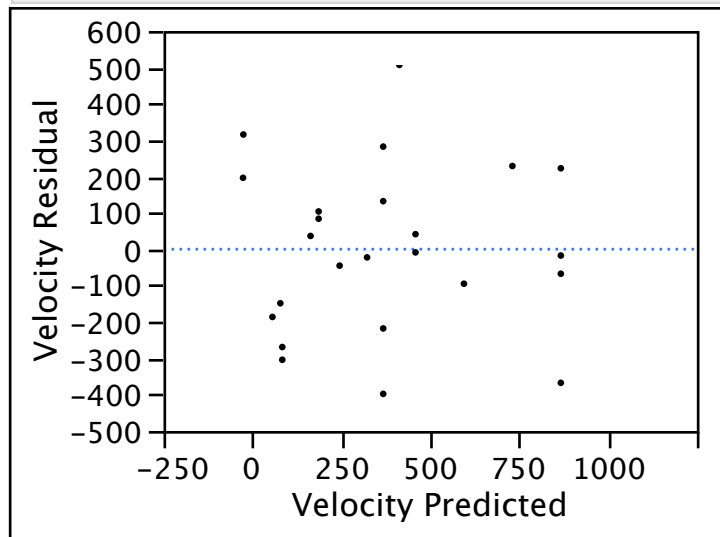
##### Parameter Estimates

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	-40.78365	83.43887	-0.49	0.6298
Distance	454.15844	75.23711	6.04	<.0001*

##### Effect Tests

Source	Nparm	DF	Sum of Squares	F Ratio	Prob > F
Distance	1	1	1976648.3	36.4377	<.0001*

##### Residual by Predicted Plot

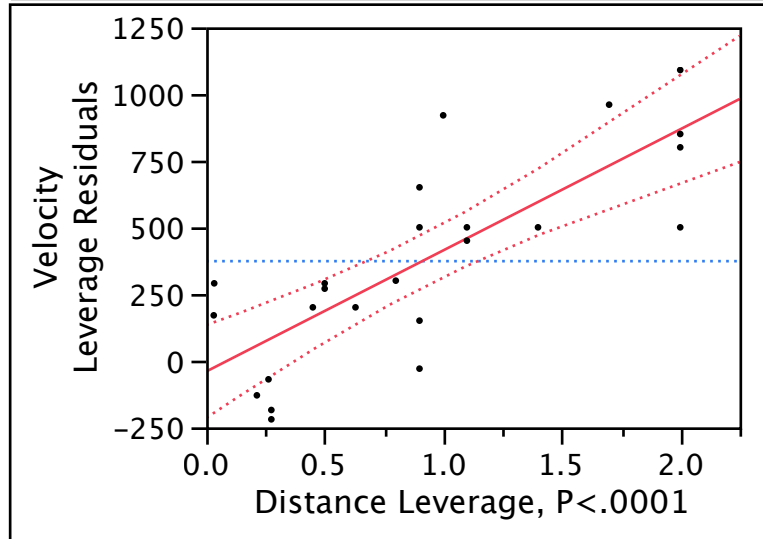


## Hubble: Fit Least Squares

Response Velocity

Distance

Leverage Plot



## Appendix 2

We have

$$(X^T X)^{-1} = \begin{bmatrix} 0.12833882 & -0.09510043 \\ -0.09510043 & 0.10434830 \end{bmatrix}$$

1a		3
1b		3
1c		6
1d		3
2a		3
2b		3
2c		2
3		2
Total		25