

Practice Exam 1 Solutions, ST518/590

1. Summary statistics from the heights, x , and weights, y , of $n = 101$ football players are given below.

$$\begin{array}{ll} \bar{x} = 73.5 \text{ inches} & \bar{y} = 226.0 \text{ lbs} \\ s_x^2 = 5.8 \text{ inches}^2 & \implies S_{xx} = \sum (x_i - \bar{x})^2 = ? \quad \boxed{580} \\ s_y^2 = 1639.1 \text{ lbs}^2 & \implies SS[Total] = \sum (y_i - \bar{y})^2 = ? \quad \boxed{163910} \\ s_{xy} = 62.3 \text{ lb-inches} & \end{array}$$

- (a) Report r_{xy} , the sample correlation coefficient for weight and height among these $n = 101$ players.

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{62.3}{\sqrt{5.8 * 1639.1}} = 0.64.$$

- (b) Suppose that these 101 players can be regarded as a random sample from a bivariate population of interest. Briefly specify a simple regression model for Y_1, \dots, Y_n in which the mean weight of players in the population is a linear function of height, x_i .

$$Y_i = \beta_0 + \beta_1 x_i + E_i \text{ for } i = 1, 2, \dots, 101.$$

$$\text{where } E_i \stackrel{iid}{\sim} N(0, \sigma^2).$$

- (c) Fill in the spaces after the question marks above with the total sum of squares of the weights, $SS[Total]$ and the sum of squares for heights, S_{xx} . Use the summary statistics to obtain the least squares (LS) estimate of the slope. Specify the units.

$$SS[Total] = \sum (y_i - \bar{y})^2 = (n - 1)s_y^2 = 163910$$

$$\sum (x_i - \bar{x})^2 = (n - 1)s_x^2 = 580$$

$$\hat{\beta} = \frac{s_{xy}}{s_x^2} = \frac{62.3}{5.8} = 10.7(\text{lbs/inches}).$$

- (d) True or false: in SLR, the least squares (LS) estimate of slope satisfies

$$\hat{\beta}_1 = r_{xy} \frac{s_x}{s_y}.$$

(The truth is that $\hat{\beta}_1 = r_{xy} \frac{s_y}{s_x}$.)

- (e) True or false: $SS[Reg] = \hat{\beta}_1^2 \sum (x_i - \bar{x})^2$
 (f) True or false: $SS[E] = \sum (\hat{y}_i - \bar{y})^2$. *The truth is that $SS[E] = \sum (y_i - \hat{y}_i)^2$.*
 (g) True or false: under repeated sampling in the SLR model, the estimated slope and intercept will be positively correlated in this context.

(h) Complete the ANOVA table below:

Source	<i>df</i>	Sum of squares	Mean square
Regression	1	66902	66902
Error	99	97008	980
Total	100	163910	

(i) The observed F -ratio is greater than the 99.99 th percentile of the F distribution with 1 and 99 numerator and denominator degrees of freedom. With regard to the association of y and x ,

- i. ... this provides strong evidence of a linear association. (or false.)
- ii. ... this provides evidence of a strong linear association. (True or)

(j) Evaluating the least squares regression line at $x = 73.5$ yields

$$\hat{\mu}(73.5) = \hat{\beta}_0 + \hat{\beta}_1(73.5) = 225.5$$

- i. Report an estimated standard error for the estimated mean weight among the population of players of height = 73.5 inches.

$$SE(\hat{\mu}(73.5)) = \sqrt{MSE(\frac{1}{101})} = 3.11lbs.$$

- ii. Estimate the standard deviation among this population of players with $x = 73.5$ inches.

$$\hat{\sigma} = \sqrt{MSE} = 31.3lbs.$$

2. The data below taken from Dickey's website correspond to temperatures taken in January, 1992.

	Mon		Tues		Wed	
	H	L	H	L	H	L
seattle	51	44	52	44	59	47
boston	29	12	32	29	44	28
richmond	47	23	55	40	51	28
louisville	53	30	37	29	46	24
lubbock	46	40	48	42	54	42
omaha	31	26	36	22	47	31
sanfran	56	47	65	49	65	47
philly	36	18	46	27	46	26
cincinnati	50	25	36	29	41	30
phoenix	74	49	75	48	75	48
miami	72	68	77	71	79	72
milwaukee	31	23	33	26	35	26
dallas	50	47	53	47	56	44
burlingvt	20	-2	28	03	39	24
buffalo	34	18	34	28	32	26
charlotte	49	38	60	41	59	41
bismark	27	-5	43	15	47	17
elpaso	61	34	64	33	64	32
rapidcity	46	20	62	25	57	41

Consider predicting the high temperature thi for a given day using a linear function of as many as 4 independent variables: the low and high temp from two days ago ($lo2, hi2$) and yesterday's low and high temp (ylo, yhi). The correlations of thi with the four predictors, along with p - values for tests of 0 correlation are given below:

Pearson Correlation Coefficients, N = 19
 Prob > |r| under H0: Rho=0

	thi
hi2	0.84769 <.0001
lo2	0.77497 <.0001
yhi	0.94572 <.0001
ylo	0.77438 <.0001
tlo	0.83471 <.0001

Several multiple regression models were considered. Sums of squares appear below. Parameter estimates are given for the full model.

- Model 1 $\mu(ylo, yhi, lo2, hi2) = \beta_0 + \beta_1 lo2 + \beta_2 hi2 + \beta_3 ylo + \beta_4 yhi$
- Model 2 $\mu(ylo, yhi, lo2, hi2) = \beta_0 + \beta_3 ylo + \beta_4 yhi$
- Model 3 $\mu(ylo, yhi, lo2, hi2) = \beta_0 + \beta_2 hi2 + \beta_4 yhi$
- Model 4 $\mu(ylo, yhi, lo2, hi2) = \beta_0 + \beta_4 yhi$

Model: MODEL1

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	4	2624.48822	656.12206	36.43	<.0001
Error	14	252.14336	18.01024		
Corrected Total	18	2876.63158			

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	15.67260	4.25081	3.69	0.0024
yhi	1	0.71275	0.13444	5.30	0.0001
ylo	1	-0.12905	0.20062	-0.64	0.5305
hi2	1	-0.00512	0.17243	-0.03	0.9768
lo2	1	0.21464	0.17800	1.21	0.2479

Model: MODEL2

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	2	2590.67200	1295.33600	72.48	<.0001
Error	16	285.95958	17.87247		
Corrected Total	18	2876.63158			

Model: MODEL3

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	2	2593.51335	1296.75667	73.28	<.0001
Error	16	283.11823	17.69489		
Corrected Total	18	2876.63158			

Model: MODEL4

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	2572.83737	2572.83737	143.97	<.0001
Error	17	303.79421	17.87025		

- (a) True or false: the regression coefficient β_4 has the same interpretation in all four models.

In MLR, the regression coefficients are partial. That is, they describe the mean change in y when controlling for all other variables in the model. Since the models are different, the variables being controlled are different.

- (b) Note that $F(0.05, 2, 16) = 3.63$ and $F(0.05, 1, 18) = 4.41$ and $F(.05, 2, 14) = 3.74$.
- By comparison of a reduced model with Model 1, use an F -ratio to test the hypothesis that today's expected high temperature does not depend on temperatures from two days ago, given the other predictors.

$$F = \frac{R(\beta_1, \beta_2 | \beta_0, \beta_3, \beta_4)/2}{MSE_{full}} = \frac{(2624 - 2590)/2}{18.0} = 0.9.$$

Not significant: the model which does not include temps from two days ago is plausible.

- By comparison of a reduced model with Model 1, use an F -ratio to test the hypothesis that today's expected high temperature does not depend on either of last two low temperatures, given the other predictors.

$$F = \frac{R(\beta_1, \beta_3 | \beta_0, \beta_2, \beta_4)/2}{MSE_{full}} = \frac{(2624 - 2593)/2}{18.0} = 0.86.$$

Not significant: the model which does not include either low temp is plausible.

- Test $H_0 : \beta_{2.4} = 0$ using level $\alpha = 0.05$. Briefly describe what's being tested here. Report and interpret the appropriate partial coefficient of determination. You may use $F(.05, 1, 16) = 4.49$. *This is a test of whether the partial slope β_2 is 0 after controlling for x_4 . That is, it compares Models 3 and 4. After controlling for yesterday's high temp, is it plausible that the high from two days ago is not linearly associated with the response?*

$$F = \frac{R(\beta_2 | \beta_0, \beta_4)}{MSE_3} = \frac{2593.5 - 2572.8}{17.7} = 1.2$$

Not significant: the model with yesterday's high is not improved by adding the high from two days ago.

$$r_{y2.4}^2 = \frac{R(\beta_2 | \beta_0, \beta_4)}{SS[Total] - R(\beta_4 | \beta_0)} = \frac{2593.5 - 2572.8}{303.8} = 0.068.$$

- (c) Provide a matrix formulation of Model 1. Make sure to identify X, Y, β and an error vector in the model, specifying the dimension of each. Provide matrix expressions for $\hat{\beta}$ and its estimated variance-covariance matrix $\hat{\Sigma}$.

$$\begin{matrix} Y & = & X & \beta & + & E \\ (19 \times 1) & & (19 \times 5) & (5 \times 1) & & (19 \times 1) \end{matrix}$$

β is a vector of unknown regression coefficients, and $E \sim N_{19}(0, \sigma^2 I_{19})$. Further, $\hat{\beta} = (X'X)^{-1}X'Y$ and $\hat{\Sigma} = (X'X)^{-1}\hat{\sigma}^2$.

- (d) Consider Model 1. Provide a matrix expression for an estimate of the mean mean high temperature among days where the preceding measurements were $(lo2, hi2, ylo, yhi) = (50, 30, 60, 50)$. Provide an expression for the standard error as well.

Let $x_0 = (1, 50, 30, 60, 50)$. With $\hat{\beta} = (X'X)^{-1}X'Y$, the estimated mean is

$$\hat{\mu}(x_0) = x_0\hat{\beta}$$

with standard error $SE(\hat{\mu}(x_0)) = \sqrt{x_0\hat{\Sigma}x_0'} = \sqrt{x_0(X'X)^{-1}x_0'\hat{\sigma}^2}$

- (e) True or false: There is no evidence of any linear association between today's high temperature and the high temperature from two days ago.

The truth is that the sample correlation between today's high and that from two days ago is $r = 0.77$ $r = 0.85$ which differs significantly from 0 ($p < 0.0001$).