

# Homework Three

Statistics 151a (Linear Models)

Due on 14 October 2015

06 October, 2015

1. Suppose  $u_1, \dots, u_n$  form an orthonormal basis of  $\mathbb{R}^n$ .

a) For every  $y \in \mathbb{R}^n$ , show that the following is true **(1.5 points)**

$$y = \sum_{i=1}^n (u_i^T y) u_i.$$

b) Show that  $\sum_{i=1}^n u_i u_i^T$  equals the  $n \times n$  identity matrix. **(1.5 points)**

c) Show that the squared norm of  $\sum_{i=1}^n c_i u_i$  equals  $\sum_{i=1}^n c_i^2$ . **(1.5 points)**

d) Fix  $1 \leq r \leq n$ . For every  $y \in \mathbb{R}^n$ , show that the projection of  $y$  onto  $sp\{u_1, \dots, u_r\}$  equals  $\sum_{i=1}^r (u_i^T y) u_i$ . Show that the projection of  $y$  onto the orthogonal complement of  $sp\{u_1, \dots, u_r\}$  equals  $\sum_{i>r} (u_i^T y) u_i$ . **(3 points)**.

e) Fix  $1 \leq r \leq n$ . For every  $y \in \mathbb{R}^n$ , show that the squared length of the projection of  $y$  onto  $sp\{u_1, \dots, u_r\}$  equals  $\sum_{i=1}^r (u_i^T y)^2$ . Show that the squared length of the projection of  $y$  onto the orthogonal complement of  $sp\{u_1, \dots, u_r\}$  equals  $\sum_{i>r} (u_i^T y)^2$ . **(1.5 points)**.

2. Consider the simple linear regression model  $y_i = \beta_0 + \beta_1 x_i + e_i$  for  $i = 1, \dots, n$  where  $\mathbb{E}e = 0$  and  $Cov(e) = \sigma^2 I_n$ . Suppose that the explanatory variable values  $x_1, \dots, x_n$  are not all constant.

a) Show that both parameters  $\beta_0$  and  $\beta_1$  are estimable. **(2 points)**

b) Show that the fitted regression line passes through the point  $(\bar{x}, \bar{y})$  where  $\bar{x} = \sum_i x_i / n$  and  $\bar{y} = \sum_i y_i / n$ . **(2 points)**

c) Suppose  $\bar{x} = 0$ . Then show that  $\hat{\beta}_0$  and  $\hat{\beta}_1$  (these represent the least squares estimators) are uncorrelated. **(3 points)**

d) Suppose  $\bar{x} = 0$  and  $e_1, \dots, e_n$  are jointly normal. Show that  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are independent. **(1 point)**

3. Consider the linear model  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + e_i$  for  $i = 1, \dots, n$  where  $\mathbb{E}e = 0$  and  $\text{Cov}(e) = \sigma^2 I_n$ . Let  $\bar{x}_j = \sum_{i=1}^n x_{ij}/n$  denote the sample mean of the  $j$ th explanatory variable for  $j = 1, \dots, p$ .
- Show that  $\beta_0 + \beta_1 \bar{x}_1 + \dots + \beta_p \bar{x}_p$  is estimable. (**2 points**)
  - What is the least squares estimate of  $\beta_0 + \beta_1 \bar{x}_1 + \dots + \beta_p \bar{x}_p$  and why? (**2 points**)
  - What is the variance of the least squares estimate in (b) and how would you estimate it from the regression data? (**2 points**)
4. Do not use R for this problem. Consider the body fat dataset that we used in class. I want to fit the model for *BODYFAT*

$$\beta_0 + \beta_1 AGE + \beta_2 WEIGHT + \beta_3 HEIGHT + \beta_4 (WEIGHT + 3 * HEIGHT) + \beta_5 WRIST + e$$

which I accomplish by the following R code resulting in the output given below:

```
> model = lm(BODYFAT ~ AGE + WEIGHT + HEIGHT + I(WEIGHT + 3*HEIGHT) + WRIST, data = body)
> summary(model)
```

Call:

```
lm(formula = BODYFAT ~ AGE + WEIGHT + HEIGHT + I(WEIGHT + 3*HEIGHT) + WRIST, data = body)
```

Residuals:

Min	1Q	Median	3Q	Max
-20.5918	-3.3673	-0.0016	3.4240	12.8823

Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	47.21461	8.89363	5.309	2.46e-07 ***
AGE	0.20629	0.02807	7.349	2.91e-12 ***
WEIGHT	0.24341	0.01672	14.562	< 2e-16 ***
HEIGHT	-0.44389	0.09706	-4.574	7.59e-06 ***
I(WEIGHT + 3 * HEIGHT)	NA	NA	NA	NA
WRIST	-2.73998	0.55167	-4.967	1.27e-06 ***

---

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

Residual standard error: 5.142 on 247 degrees of freedom

Multiple R-squared: 0.5669, Adjusted R-squared: 0.5599

F-statistic: 80.82 on 4 and 247 DF, p-value: < 2.2e-16

- a) Why does R produce NAs in the output? (**2 points**)
- b) The estimate for  $\beta_2$  is apparently 0.24341. Does this make sense? Explain. (**2 points**)
- c) I decide against including the variable  $WEIGHT + 3 * HEIGHT$  in the model and just intend to fit

Model M: `BODYFAT ~ AGE + WEIGHT + HEIGHT + WRIST`

What is the RSS for this model? Why? (**2 points**)

- d) The model M has too many parameters for my liking; so I decide to consider the following model:

Model m: `BODYFAT ~ AGE + WEIGHT`

which gave me the following R output:

Call:

`lm(formula = BODYFAT ~ AGE + WEIGHT, data = body)`

Residuals:

Min	1Q	Median	3Q	Max
-15.3171	-4.3293	0.2917	3.9898	18.5237

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-18.37392	2.57545	-7.134	1.06e-11 ***
AGE	0.18269	0.02853	6.403	7.54e-10 ***
WEIGHT	0.16271	0.01224	13.298	< 2e-16 ***

---

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

Residual standard error: 5.696 on 249 degrees of freedom

Multiple R-squared: 0.4642, Adjusted R-squared: 0.4599

F-statistic: 107.9 on 2 and 249 DF, p-value: < 2.2e-16

Find the  $p$ -value for testing the model  $m$  against the model  $M$ . If you do not have a calculator that can calculate the  $p$ -value, write the answer in terms of the  $F$ -statistic. (**4 points**).

5. Do not use R for this problem. For the Bodyfat dataset used in class, consider the linear model

$$\text{BODYFAT} = \beta_0 + \beta_1 \text{AGE} + \beta_2 \text{WEIGHT} + \beta_3 \text{HEIGHT} + \beta_4 \text{THIGH} + e.$$

If  $X$  denotes the  $X$ -matrix for this regression, then R tells me that  $(X^T X)^{-1}$  equals

$$\begin{pmatrix} 3.740212022 & -5.908839e-03 & 6.662131e-03 & -3.218478e-02 & -4.048954e-02 \\ -0.005908839 & 3.238651e-05 & -1.222844e-05 & 3.416435e-05 & 7.148358e-05 \\ 0.006662131 & -1.222844e-05 & 2.632523e-05 & -4.483900e-05 & -1.292477e-04 \\ -0.032184784 & 3.416435e-05 & -4.483900e-05 & 3.866749e-04 & 1.944136e-04 \\ -0.040489539 & 7.148358e-05 & -1.292477e-04 & 1.944136e-04 & XXXXXX \end{pmatrix}$$

The regression summary given by R is as follows:

Call:

```
lm(formula = BODYFAT ~ AGE + WEIGHT + HEIGHT + THIGH, data)
```

Residuals:

Min	1Q	Median	3Q	Max
-17.3699	-3.9361	-0.0351	3.6796	16.0833

Coefficients:

	Estimate	Std. Error	t value
(Intercept)	-1.07425	10.30553	-0.104
AGE	0.18901	0.03033	6.233
WEIGHT	0.12373	XXXXXX	XXXXXX
HEIGHT	-0.46074	0.10478	-4.397
THIGH	XXXXXX	0.14952	2.444

---

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

Residual standard error: XXXXXX on 247 degrees of freedom

Multiple R-squared: 0.5349

F-statistic: XXXXXX on 4 and 247 DF, p-value: < 2.2e-16

- Fill the six missing values (one in the  $(X^T X)^{-1}$  matrix and five in the R summary; all indicated by XXXXX) above. (6 points)
- Based on this dataset and the above linear model, I want to predict the bodyfat percentage for a new individual who is 30 years of age, weighs 180 lbs, is 72 inches tall and who thigh circumference is 60 cm. For this consider the following output:

```
> x0 = data.frame(AGE = 30, WEIGHT = 180, HEIGHT = 72, THIGH = 60)
> predict(M, x0, interval = "confidence")
      fit      lwr      upr
XXXXXX 14.56726 XXXXXX
```

```
> predict(M, x0, interval = "prediction")
      fit      lwr      upr
XXXXXX XXXXXXX XXXXXXX
```

Fill in the four missing values above. **(5 points)**

- c) Consider the following R output for testing Model 1 against Model 2 where

Analysis of Variance Table

Model 1: BODYFAT ~ I(AGE + THIGH) + WEIGHT + HEIGHT

Model 2: BODYFAT ~ AGE + WEIGHT + HEIGHT + THIGH

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	XXX	XXXXX				
2	XXX	XXXXX	XX	XXXXXXX	1.6203	0.2042

Fill in the six missing values. **(6 points)**

6. Determine whether each of following statements is true or false. Provide reasons in each case. **(11 points - 1 point for each question. No point will be awarded if no reason is provided.)**

- In simple linear regression (i.e., when there is only one explanatory variable), the slope of the regression line can never be larger than one.
- Again consider simple linear regression. Suppose that the response and explanatory variable values are standardized to have mean zero and unit standard deviation. Then the slope of the regression line can never be larger than one.
- The residual standard error always increases when explanatory variables are removed from the linear model.
- Any linear function of  $\beta = (\beta_0, \dots, \beta_p)$  is estimable when the matrix  $X^T X$  is invertible.
- Because of the assumptions underlying the linear model, the residuals  $\hat{e}_1, \dots, \hat{e}_n$  all have the same variance.
- If the normality assumption is violated, then the vector of residuals and the vector of fitted values may not be orthogonal.
- If the normality assumption is violated, then the vector of residuals and the vector of fitted values may not be uncorrelated.
- If the normality assumption is violated, then the vector of residuals and the vector of fitted values may

not be independent.

- i) A small  $p$ -value for the F-statistic in the regression summary validates the linear model.
- j) An archaeologist fits a regression model rejecting the hypothesis that  $\beta_2 = 0$  after getting a  $p$ -value less than 0.005. This must mean that  $\beta_2$  must be large.
- k) An archaeologist fits a regression model rejecting the hypothesis that  $\beta_2 = 0$  after getting a  $p$ -value less than 0.005. This must mean that  $\hat{\beta}_2$  must be large.