

Homework 4  
Psychology 312

1. *Multiple Regression. The Model and its Estimates.* In standard multiple regression, we fit the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad (1)$$

where  $\mathbf{y}$  is an  $n \times 1$  vector of independent criterion scores,  $\mathbf{X}$  an  $n \times p'$  matrix of scores on  $p'$  predictor terms, including an intercept term. So in the typical case where there are  $p$  predictors and an intercept,  $p' = p + 1$ . The predictor scores in  $\mathbf{X}$  are considered fixed constants. The errors in  $\mathbf{e}$  are independent random variables, with zero means and constant variance  $\sigma^2$ . So we say that

$$\text{Var}(\mathbf{e}) = \sigma^2(\mathbf{I}) \quad (2)$$

Note that since both  $\mathbf{X}$  and  $\boldsymbol{\beta}$  contain only constants, it immediately follows that

$$\text{Var}(\mathbf{y}) = \text{Var}(\mathbf{e}) = \sigma^2(\mathbf{I}) \quad (3)$$

and

$$\mathbb{E}(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta} \quad (4)$$

We'll assume a parameterization  $\mathbf{X}$  of full column rank, including a column of 1s for the intercept if necessary. So  $(\mathbf{X}'\mathbf{X})^{-1}$  exists. The  $p' \times 1$  vector  $\boldsymbol{\beta}$  contains the regression parameters, including the intercept. The vector  $\mathbf{e}$  consists of errors.

Note: these errors are *unobservable*, because  $\boldsymbol{\beta}$  is itself not known. We can *estimate*  $\boldsymbol{\beta}$ , and use the estimate to construct *estimates* of errors. Note again, these estimates are the model *residuals*. The residuals are not the errors — they are estimates of the errors. Although the errors are assumed to be uncorrelated, the residuals, as we shall see, are not.

The least squares estimates of the regression parameters in  $\boldsymbol{\beta}$  are calculated as

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (5)$$

The *predicted scores* are

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \quad (6)$$

$$= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (7)$$

$$= \mathbf{P}_x\mathbf{y} \quad (8)$$

The *residuals* are

$$\begin{aligned}
\hat{e} &= \mathbf{y} - \hat{\mathbf{y}} \\
&= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \\
&= \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\
&= \mathbf{y} - \mathbf{P}_x\mathbf{y} \\
&= (\mathbf{I} - \mathbf{P}_x)\mathbf{y} \\
&= \mathbf{Q}_x\mathbf{y}
\end{aligned} \tag{9}$$

- (a) **(10 points).** While the elements of  $\boldsymbol{\beta}$  are constants, the elements of  $\hat{\boldsymbol{\beta}}$  are random variables, and have an expected value and a variance. All regression programs compute and report estimated *standard errors* for the elements of  $\hat{\boldsymbol{\beta}}$ , i.e., the estimates of regression parameters. These are obtained as follows: (a) Compute an estimate of the variance-covariance matrix of the random vector  $\hat{\boldsymbol{\beta}}$ ; (b) Since the diagonal elements of this estimated variance-covariance matrix are estimated variances, take the square roots of these diagonal elements to produce estimated standard errors.

Given the above, and remembering that  $\mathbf{P}_x$  and  $\mathbf{Q}_x$  are matrices of constants, **prove the following two results.**

i.

$$\mathbb{E}(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta} \tag{10}$$

*Hint.* Combine Equations 4 and 5.

ii.

$$\Sigma_{\hat{\boldsymbol{\beta}}\hat{\boldsymbol{\beta}}} = \text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \tag{11}$$

*Hint.* Use your results from expected value algebra and compute  $\text{Var}(\hat{\boldsymbol{\beta}})$  as  $\text{Var}(\mathbf{A}'\mathbf{y})$ , where  $\mathbf{A}' = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ , remembering the result on the variance of a linear combination that I described as essential to multivariate analysis theory, and using the result of Equation 3.

Note that, since  $\sigma^2$  is not known, it must be estimated. The unbiased estimator is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{e}_i^2}{n - p'} \tag{12}$$

The estimate for the variance-covariance matrix of  $\hat{\boldsymbol{\beta}}$  is therefore

$$\hat{\Sigma}_{\hat{\boldsymbol{\beta}}\hat{\boldsymbol{\beta}}} = \widehat{\text{Var}}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1} \tag{13}$$

- (b) **(10 points).** Start up R, and load in the file *KidIqData.csv*, using the `read.csv` command. Download the R Utility Functions from the R section of the website. They are accessible under the *R Utility Functions* link.

Load in the functions as shown below. The actual filename is *Steiger R Library Functions.txt*. You will probably find it convenient to attach the data file. Treat the variable `kid.iq` as the dependent variable  $\mathbf{y}$ , and treat the variables `mom.iq` and `mom.grad` as the predictors in a matrix  $\mathbf{X}$ . Verify that the file has 500 observations on 3 variables. Add an intercept variable (a column of 1's) to  $\mathbf{X}$ . Here is an example of how to do this.

```
data <- read.csv("KidIqData.csv")
attach(data)
dim(data)

## [1] 500   3

names(data)

## [1] "kid.iq"    "mom.iq"    "mom.grad"

## Create an intercept variable
source("Steiger R Library Functions.txt")
one <- UnitVector(500)
X <- cbind(one,mom.iq,mom.grad)
y <- kid.iq
```

The variable `mom.grad` is a binary variable coded 0-1, depending on whether the mother graduated from high school. Using the `lm` function, predict `kid.iq` from `mom.iq` and `mom.grad` and save the result as a `lm` object. (Note, if you don't know how to do this, examine the online lecture notes on multiple regression and chapters 3 and 4 from Gelman and Hill.) Apply the `summary` function to the object. You should see values for the intercept and regression coefficients corresponding to  $\beta' = [32.83, 0.66, 14.99]$ . Notice also that there are standard errors reported for the coefficients. The standard errors are, respectively, 4.891, 0.049, 1.948. Notice also that the program prints a “Residual Standard Error” of 14.46. This is the square root of the quantity defined in Equation 12. Here is the actual output.

```
## 
## Call:
## lm(formula = kid.iq ~ mom.iq + mom.grad)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -37.183  -9.282  -0.201  10.247  37.993
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 32.83107   4.89084   6.713 5.22e-11 ***
```

```

## mom.iq      0.66067   0.04934 13.390 < 2e-16 ***
## mom.grad    14.98590  1.94786  7.694 7.76e-14 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.46 on 497 degrees of freedom
## Multiple R-squared: 0.4631, Adjusted R-squared: 0.4609
## F-statistic: 214.3 on 2 and 497 DF, p-value: < 2.2e-16

```

**Calculate  $\hat{\beta}$  from Equation 5, and verify that it agrees with the R output.**

- (c) **(10 points).** Once you have the fit object, it is easy to extract the residuals. Just apply the `residuals` function to the fit object you have saved. **Calculate the estimated residual variance** defined in Equation 12. Check that its square root is equal to the “residual standard error” shown in the R output. (*Hint.* Remember the definition of  $p'$ .)

The predicted scores can be generated in R by applying the `predict` function to the linear model object. Of course, we can also calculate them as  $\mathbf{P}_x \mathbf{y}$ . My support routines contain the functions `P` and `Q` to allow you to create projectors easily and directly. Compute the predicted scores using `predict`, and extract the first predicted score. Compare it to the first element of the array calculated in R using the as `y.hat <- P(X) %*% y`. Then compare the value of the first element of the vector of residuals produced by the `residuals` function, and verify that `(Q(X) %*% y)[1]` is in fact equal to the first residual.

- (d) **(10 points).** Using the estimated residual variance you computed in the preceding part, and also using the result of Equation 13, **compute the estimated standard errors** of the elements of  $\hat{\beta}$ , and verify that they agree with the R output.
- (e) **(10 points).** The squared multiple correlation  $R^2$  between  $\mathbf{y}$  and the non-intercept variables in  $\mathbf{X}$  can be calculated in a number of ways. One way is from a formula analogous to the population formula given in the lecture slides on Key Regression Algebra, i.e.

$$R^2 = \frac{s_{yx} S_{xx}^{-1} s_{xy}}{s_y^2} = \frac{s'_{xy} S_{xx}^{-1} s_{xy}}{s_y^2} \quad (14)$$

Another way is to use the linear model fit object, compute the correlation between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ , and square it, i.e.

$$R^2 = r_{y,\hat{y}}^2 \quad (15)$$

Compute  $R^2$  both ways with R (remember, the first calculation will require that you drop the intercept column from  $\mathbf{X}$ ), and verify that the result agrees with the value of 0.4631 printed in the summary table from the `lm` object.

- (f) **(5 points).** Verify in one line of R that the predicted and error scores from the linear model analysis have a correlation of zero.
2. *Eigenvalues and Eigenvectors.* Load in the data from the *AthleticsData.csv* file.
- (a) **(10 points).** Compute the correlation matrix for the AthleticsData. Call it  $\mathbf{R}_{yy}$ . Then compute the eigenvectors and eigenvalues of the correlation matrix using the `eigen` command. Create a matrix  $\mathbf{V}$  containing the eigenvectors and a diagonal matrix  $\mathbf{D}$  containing the eigenvalues on the diagonal. Verify in R that
- $$\mathbf{R}_{yy} = \mathbf{V}\mathbf{D}\mathbf{V}' = \mathbf{F}\mathbf{F}'$$
- up to rounding error by computing the square root of the sum of squared differences between  $\mathbf{R}_{yy}$  and  $\mathbf{V}\mathbf{D}\mathbf{V}'$ .
- (b) **(10 points).** Compute  $\mathbf{F} = \mathbf{V}\mathbf{D}^{1/2}$ . The matrix is sometimes called the “principal component pattern.” Install the `psych` library on your system, load the library, and compute the raw principal component pattern with the command
- ```
principal(AthleticsData,nfactors=9,rotate="none").
```
- Compare the output with your  $\mathbf{F}$ . Explain any discrepancy.
- (c) **(10 points).** Extract the first 3 columns of your  $\mathbf{F}$ . Call the extracted columns  $\mathbf{G}$ . Compute  $\hat{\mathbf{R}} = \mathbf{G}\mathbf{G}'$ , and compare it to  $\mathbf{R}_{yy}$  as follows. Extract the elements of  $\mathbf{R}_{yy}$  into an  $81 \times 1$  vector, and, likewise, extract the elements of  $\hat{\mathbf{R}}$  as an  $81 \times 1$  vector. Compute the correlation between the two vectors and do a scatterplot. Add an identity line (slope 1, intercept 0) to the plot to aid interpretation. Notice anything interesting in the scatterplot? Is there some way you could “clean up” the scatterplot to improve it as a way of comparing the two matrices?
3. **(15 points).** In the class notes on Key Regression Algebra, we discussed a “regression component” system of the form

$$\mathbf{y} = \mathbf{F}\mathbf{x} + \mathbf{e}$$

in which all variables are in deviation score form. We define  $\Sigma = E(\mathbf{y}\mathbf{y}')$ , and assume  $\mathbf{x} = \mathbf{B}'\mathbf{y}$ , and  $\mathbf{F}$  is a set of least squares linear regression

weights for predicting  $\mathbf{y}$  from  $\mathbf{x}$ . In the notes, we showed that  $\mathbf{F}$  can be computed as

$$\mathbf{F} = \boldsymbol{\Sigma} \mathbf{B} (\mathbf{B}' \boldsymbol{\Sigma} \mathbf{B})^{-1} \quad (16)$$

This means that given  $\mathbf{B}$  and  $\boldsymbol{\Sigma}$ , the system is completely determinate.

- (a) Using simple (but careful and involved) substitution, show that, given  $\mathbf{F}$  and  $\boldsymbol{\Sigma}$ ,  $\mathbf{B}$  is completely determined and may be calculated as

$$\mathbf{B} = \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{F}' \boldsymbol{\Sigma}^{-1} \mathbf{F})^{-1} \quad (17)$$

(*Hint:* simply substitute Equation 16 into the right side of Equation 17 and show it reduces to  $\mathbf{B}$ .)

- (b) Principal components analysis is a special case of regression component analysis, so we may apply the rules of the latter to the former. Earlier, we demonstrated that if  $\boldsymbol{\Sigma} = \mathbf{V} \mathbf{D} \mathbf{V}'$  is an Eckart-Young decomposition of  $\boldsymbol{\Sigma}$  into eigenvectors and eigenvalues, that the principal components pattern may be written as  $\mathbf{F} = \mathbf{V} \mathbf{D}^{1/2}$ . Using Equation 17, and your knowledge of matrix factorization and symmetric square roots, prove that, if  $\mathbf{F} = \mathbf{V} \mathbf{D}^{1/2}$ , then  $\mathbf{B} = \mathbf{V} \mathbf{D}^{-1/2}$ .
- (c) Prove that, if  $\mathbf{x} = \mathbf{B}' \mathbf{y}$ , and  $\mathbf{B} = \mathbf{V} \mathbf{D}^{-1/2}$ , then  $\text{Var}(\mathbf{x}) = \mathbf{E}(\mathbf{x}\mathbf{x}') = \mathbf{I}$ .