

Econ 21410 - Problem Set I

Review of R and Regression Analysis*

April 3, 2014

This homework should be done in LaTeX. The homework will be graded on correctness, but will also heavily weight clarity and professionalism. Being able to produce clean, clear, well documented write-ups and code is an important skill which will be rewarded. It's better to not do some of the harder parts than to turn in an incomprehensible document. Your R script as well as a log-file should be submitted. Alternatively, use knitr to print your code inline in your latex document.

Make sure to write code which is clear and flexible. Read the whole problem before you begin coding. Some parameters will change and the code should be written in a way to make this easy to implement. We will re-use code in this course. Flexibility and documentation now will save you headaches later in the quarter.

SUBMISSION: The homework must be emailed to Oliver and myself by 2p.m. Monday, April the 7th. The email must include a pdf with the filename `lastname_pset1.pdf` and R code called `lastname_pset1_code.R` where "lastname" should be replaced with your last name. The subject of your email should be [ECON 21410: pset1 submission]

1 A Quick Review of R

1. Create a vector $y = \begin{bmatrix} 100 \\ 200 \\ 300 \\ 400 \\ 500 \end{bmatrix}$

2. Set a random seed equal to 1234

3. Create a matrix x which is 5×5 and contains random draws from a normal with mean 100 and variance 10. Display this output in your code (preferably inline with knitr). None of these should require more than a single line of R. These exercises must be calculated in R, not done by hand.

4. Display $x'x$

5. Display $(x'x)^{-1}$

*Please email johneric@uchicago.edu and obrowne@uchicago.edu if you have questions.

6. Calculate the sum of the entries in y
7. Calculate the row sums of the entries of x
8. Return the maximum value in y
9. Return which entry is the maximum value of y
10. Return which entry is the maximum value of x
11. Return $3 * x$
12. Return z a new vector which contains the value of y sorted from low to high
13. Replace the third row of x with 0s and display it

2 A Quick Review of L^AT_EX

1. Display the matrix and vector x and y above in L^AT_EX (no need to include the decimals)
2. Print the symbols α , θ_j , $\lambda_{t,t+1}$, $\gamma^{s,t}$ inline with text.
3. Write on its own centered line:

$$\sum_{t=1}^T \frac{a_t}{b_t} \xrightarrow{p} \infty$$

4. Write $a \neq b$ and $c \geq d$

3 Getting started with Github

1. You should have already made a github.com account and shared your user name with Oliver and myself, if not, do so as soon as you read this.
2. Go to “CompEcon” at github.com/CompEcon. You should be able to see 3 repositories, including a repository with your name in it. Go into this repository and click on README.md. Modify the README.md file to include your name and email address.

4 Regression

Recall what you have learned about ordinary least squares and what we reviewed in class. In this problem you will estimate OLS coefficients and standard errors a number of ways. You will experiment with how sample size, optimizer options, and other parameters change the estimates. For this assignment, we will provide you with code to generate *almost* all of the data you will need to answer this problem.

Here is the code needed to generate the data:

```

# ===== TITLE: computational economics: assignment 1

# AUTHOR: John Eric Humphries

# abstract: problem set on regression for econ 21410

# Date: 2014-03-14

## @knitr code_part1 =====

# ===== Section 0: setup ===== setwd('')
rm(list = ls()) # Clear the workspace
set.seed(21410) # Set random seed

# ===== Section 1: Generating Data =====
n <- 200 # observations
x1 <- rbinom(n, 10, 0.5)
x2 <- rnorm(n, 20, 10)
X <- cbind(1, x1, x2)
eps <- rnorm(n, 0, sqrt(4))
beta <- matrix(cbind(2, 3, 1.4), 3, 1)
Y <- X %*% beta + eps

```

We will consider the model:

$$Y = X\beta + \epsilon$$

Algebraic Approach

1. What is the algebraic formula for the OLS estimate of $\hat{\beta}$ (in matrix notation)?
2. What are the dimensions of Y_i , X_i , β , and ϵ
3. What are the formulas for the standard errors for the three elements in β ?

Least Squares Approach

1. Write out the minimization problem which OLS solves.
2. Take the first order condition and solve for $\hat{\beta}$.¹ Show that you get the algebraic formula from above.

¹If taking derivatives of matrices is causing problems, there will be only a small deduction for doing this for the single-variate case..

Maximum Likelihood Approach

1. Assuming $\epsilon_i \sim N(0, \sigma^2)$, write the (log) likelihood contribution for a single individual
2. Write (log) likelihood for the whole sample.
3. Write the optimization problem which characterizes the maximum likelihood. Highlight why $\hat{\beta}$ estimated this way will be equivalent to the least-squares estimate of $\hat{\beta}$.
4. How many variables are being optimized
5. Are there bounds on any of the variables being optimized? Why?
6. What is the formula for the standard errors of $\hat{\beta}$? (you do not need to explicitly solve for the second derivative, just write the formula down).

Implementing OLS

Using the derivations from the previous section we will now construct a function to estimate $\hat{\beta}$ various ways. The end goal will be to construct a function which takes three inputs X , Y , and “method = c(0,1,2,3)”, where method will specify which method the code should use to estimate $\hat{\beta}$. If necessary, you may create 3 new functions rather than 1 general purpose function with little penalty.

- method=0 should call the 'lm()' function, which is R's default method for performing OLS.
- method=1 should estimate β using the algebraic approach.
- method=2 should estimate β using the least squares approach.
- method=3 should estimate β using MLE.

Methods 2 and 3 require a use of numeric optimization to solve the optimization problem. R can call world-class optimization packages, but for now we will use “optim()”, the standard multivariate optimization command.

Output: Produce a 4 x 3 matrix where each column contains the estimates of β from one of the four methods.

Optimizers

The optim() command has an option which takes a number of optimizers. Re-estimate the model using the “BFGS”, “CG”, and “SANN” methods. The BFGS method uses gradient decent. As part of this method it numerically approximates the gradient and Hessian. You can improve the numerical accuracy of the optimization by providing the analytic gradient.

In three different plots, plot the three elements of $\hat{\beta}$ for sample sizes {8, 20, 50, 100, 1000, 10000} (estimates on Y axis, sample size on X axis, one parameter per plot.) Plot a line for (a) the algebraic approach (b) the least squares approach (c) the least squares approach using “BFGS” (d) least squares using “CG” (e) least squares using “SANN” (f) least squares using “BFGS” and providing the analytic gradient.² Briefly discuss the results (no more than 5 sentences).

²Part f is quite a bit harder than parts a - e

Standard Errors

Use the “hessian=T” option in `optim()` to return the hessian of the optimization. Use it to construct standard errors for the coefficients in the MLE approach using the “BFGS” optimizer. Make a plot for each parameter with the standard error on the Y axis and the sample sizes on the x-axis (only one line needed for each plot). Discuss briefly how standard error change with sample size (no more than 5 sentences).

5 Probit Estimation

Suppose we have the same model is before:

$$Y = X\beta + \epsilon$$

, but now Y is the gain in utility from a particular choice (say the choice to work or not). Inherently, we do not observe utility, but only observe if the individual chooses to work or not. We will use a “threshold model”, where an individual chooses to work if their unobserved utility is above a particular threshold:

$$D = \begin{cases} 1, & \text{if } Y \geq T, \\ 0, & \text{if } Y < T. \end{cases}$$

We will set $T = 0$. Using our linear-in-parameters form above, we can rewrite the decision process as:

$$D = \begin{cases} 1, & \text{if } \epsilon \geq (-X\beta), \\ 0, & \text{if } \epsilon < (-X\beta). \end{cases}$$

Assume $\epsilon \sim N(0, 1)$ and answer the following questions:

1. Why can we set $T = 0$ and $\sigma_\epsilon^2 = 1$ without consequence?
2. Write out the likelihood contribution for an individual
3. Write out the likelihood contribution for the population
4. Modify the simulated data above to create a new dataset for a binary outcome. To do this, replace β_0 to be -45 rather than 2. Also adjust the variance of the error to be 1.
5. Create a function which calculates the likelihood of the data.
6. Use “optim” to solve for $\hat{\beta}$. How close are your estimates to the true parameters?
7. Estimate \hat{D} using $\hat{\beta}$ and the observed estimates. What proportion of the decisions do you get correct? Calculate this proportion with a sample size of 10, 25, 50, 100, and 5000.

Potential Side Projects

- Read about the bootstrap algorithm for calculating standard errors. Add an option to calculate the standard errors of one of the OLS methods above using bootstrap. Write no more than a page summarizing how the bootstrap algorithm is implemented. Upload that description as a page on the wiki for the final half-point. (2.5 points)

- Complete your problem set in knitr. Create a sample knitr file that can be shared with your classmates and write a page on using knitr in the class wiki. (1.5 points)
- Rewrite at least a portion of regression code above in Julia³ or C++ using the Rcpp package. Compare how long the new code takes to run in comparison with your R code. (3 points)
- Rewrite a portion of regression code in python. (1.5 points)
- Add a fifth method to your code to estimate the regression using Generalized Method of Moments. Include a short write-up of how you implemented GMM. (2.5 points)
- Make a meaningful contribution to the class wiki, start an issue and ask a valuable question, provide a detailed and useful answer to a classmate's question. Include 2-3 sentences in your homework stating your contributions. (1 point)

³Julia is a very promising new programming language for statistical computing. It is still very new, but I believe it may eventually be a quality replacement for R or python and some early investment now could be beneficial later. It is fast, has a simple syntax, is open source, and has a large community for such a young language.