

Regression Computation Assignment

The assignment is due on Thursday June 4 at 23:59.

Purpose and Overview

The purpose of this assignment is for you to demonstrate your understanding of the mathematics and computational procedures involved in the estimation of basic regression models. Below are several regression models. The first two (**m1** and **m2**) are OLS models explaining life expectancy and the latter two (**m3** and **m4**) are logistic regression models explaining democracy. All data are drawn from the Quality of Government Basic Dataset. The data and codebook are available here: <http://qog.pol.gu.se/data/datadownloads/qogbasicdata>.

```
d <- rio::import("http://www.qogdata.pol.gu.se/data/qog_bas_cs_jan15.csv")

f1 <- une_leb ~ I(gle_cgdpc/1e4) + I(ross_oil_netexp/1e3)
m1 <- lm(f1, data = d)

f2 <- une_leb ~ I(gle_cgdpc/1e4)*factor(chga_demo) +
             I(ross_oil_netexp/1e3) + factor(ht_colonial)
m2 <- lm(f2, data = d)

f3 <- chga_demo ~ al_ethnic
m3 <- glm(f3, data = d, family = binomial(link='logit'))

f4 <- chga_demo ~ al_ethnic + I(log(wdi_landarea)) +
             I(gle_cgdpc/1e4) + I(ross_oil_netexp/1e3)
m4 <- glm(f4, data = d, family = binomial(link='logit'))
```

Your task is to reproduce parts of these analyses using R (or Stata), without the use of the `lm`, `glm`, `lm.fit`, or similar ready-made functions (or, for Stata, without the use of `reg`, `logit`, `glm`, or similar commands). The instructions below outline the code you should produce. You may use R (based on techniques from lecture) or Stata (for relevant technical details on Stata, see Ch. 11 and Appendix A from Cameron and Trivedi 2010).

Submitting Your Assignment

Submit your assignment via email to Thomas (<mailto:tleeper@ps.au.dk>) in the form of a single, complete R syntax file (.R) or Stata (.do) file in your submitted assignment. Each step of the code should be numbered and labeled (in the form of a R or Stata comment) according to the numbering of the steps below. You do not need to explain, describe, or present the output of any analyses.

Ordinary Least Squares Regression

1. Construct a design matrix \mathbf{X} for `m1`. Extract the \mathbf{R} matrix from the **QR** decomposition of \mathbf{X} . Write code to produce this matrix using the formula given in class for the inverse of a 3-by-3 matrix.
2. Create a vector `y` to represent the outcome. Then construct an appropriate design matrix \mathbf{X} for `m2` (as represented by formula `f2`) from the original variables in the dataset, including the specified categorical (factor) variables and interaction term.
Note: Be careful the handling of missing values!
3. Estimate the Ordinary Least Squares (OLS) coefficients in `m2` using matrix notation (see Fox p.155). You may use a QR decomposition if you so choose.
4. Obtain the fitted values (where $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$), from model `m2` with all variables held at their observed values, except:
 - `gle_cgdpc = 10000`
 - `chga_demo = 1`
 - `ross_oil_netexp = 1000`
 - `ht_colonial = 0`
5. Estimate the variance-covariance matrix of the OLS coefficient estimates for `m2` and extract estimated standard errors from that matrix (see Fox p.158).
6. Use the estimated coefficients and standard errors for `m2` to obtain the t -statistics for each coefficient estimate under a null hypothesis $H_0 : \beta_k = 0$. Calculate the one-tailed and two-tailed p-values for each test statistic using the `pt` function, which returns the value of Student's t cumulative distribution function.

Logistic Regression

7. Represent the log-likelihood function for logistic regression as an R function (or Stata function). Recall the general structure of the log-likelihood for a binary outcome variable is:

$$\ln \mathcal{L}(\beta) = \sum_{i=1}^n y_i \ln(\pi_i) + (1 - y_i) \ln(1 - \pi_i)$$

where π_i in the logistic model is given by the inverse link function $\frac{e^{\mathbf{X}_i\beta}}{1 + e^{\mathbf{X}_i\beta}}$.

8. Estimate the coefficients for **m3** using two methods of maximum likelihood estimation: (a) a grid search and (b) an optimization algorithm implemented by **optim**. For the grid search, search for possible slope and intercept coefficients in the range between -5 and +5.
9. Estimate the standard errors for coefficient estimates in **m3** using bootstrapping.
10. Obtain 95% bootstrap confidence intervals for the **m3** coefficient estimates from the bootstrap distribution of the coefficients.
11. Estimate the coefficients for **m4** using **optim**, obtain the standard errors as the diagonal of the matrix inverse of the negative hessian matrix, calculate the z -statistic for each coefficient estimate, and determine its two-tailed p -value from **pnorm** (the normal cumulative distribution function).
12. Using your estimated coefficients for **m3** and **m4**, compute the *average* predicted probability of observing the outcome when **al_ethnic** is set to each of the five-number summary quantiles (0.00, 0.25, 0.50, 0.75, 1.00) of the observed distribution of **al_ethnic** and all other covariates are held at their observed values. Your results should be two vectors, one for each model, each containing five numbers.