

Homework 3: The Death and Life of Great American City Scaling Laws

Background: In the previous lectures and lab, we began to look at user-written functions. For this assignment we will continue with a look at fitting models by optimizing error functions, and making user-written functions parts of larger pieces of code.

In lecture, we saw how to estimate the parameter a in a nonlinear model,

$$Y = y_0 N^a + \text{noise}$$

by minimizing the mean squared error

$$\frac{1}{n} \sum_{i=1}^n (Y_i - y_0 N_i^a)^2.$$

We did this by approximating the derivative of the MSE, and adjusting a by an amount proportional to that, stopping when the derivative became small. Our procedure assumed we knew y_0 . In this assignment, we will use a built-in R function to estimate both parameters at once; it uses a fancier version of the same idea.

Because the model is nonlinear, there is no simple formula for the parameter estimates in terms of the data. Also unlike linear models, there is no simple formula for the *standard errors* of the parameter estimates. We will therefore use a technique called **the jackknife** to get approximate standard errors.

Here is how the jackknife works:

- Get a set of n data points and get an estimate $\hat{\theta}$ for the parameter of interest θ .
- For each data point i , remove i from the data set, and get an estimate $\hat{\theta}_{(-i)}$ from the remaining $n - 1$ data points. The $\hat{\theta}_{(-i)}$ are sometimes called the “jackknife estimates”.
- Find the mean $\bar{\theta}$ of the n values of $\hat{\theta}_{(-i)}$
- The jackknife variance of $\hat{\theta}$ is

$$\frac{n-1}{n} \sum_{i=1}^n (\hat{\theta}_{(-i)} - \bar{\theta})^2 = \frac{(n-1)^2}{n} \text{var}[\hat{\theta}_{(-i)}]$$

where var stands for the sample variance. (*Challenge:* can you explain the factor of $(n-1)^2/n$? *Hint:* think about what happens when n is large so $(n-1)/n \approx 1$.)

- The jackknife standard error of $\hat{\theta}$ is the square root of the jackknife variance.

You will estimate the power-law scaling model, and its uncertainty, using the data alluded to in lecture, available in the file `gmp.dat` from lecture, which contains data for 2006.

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.0 --
## v ggplot2 3.3.2      v purrr   0.3.4
## v tibble  3.0.2      v dplyr  1.0.0
## v tidyr   1.1.0      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()

library(conflicted)
gmp <- read.table("data/gmp.dat")
gmp$pop <- round(gmp$gmp/gmp$pcgmp)
conflict_prefer("filter", "dplyr")

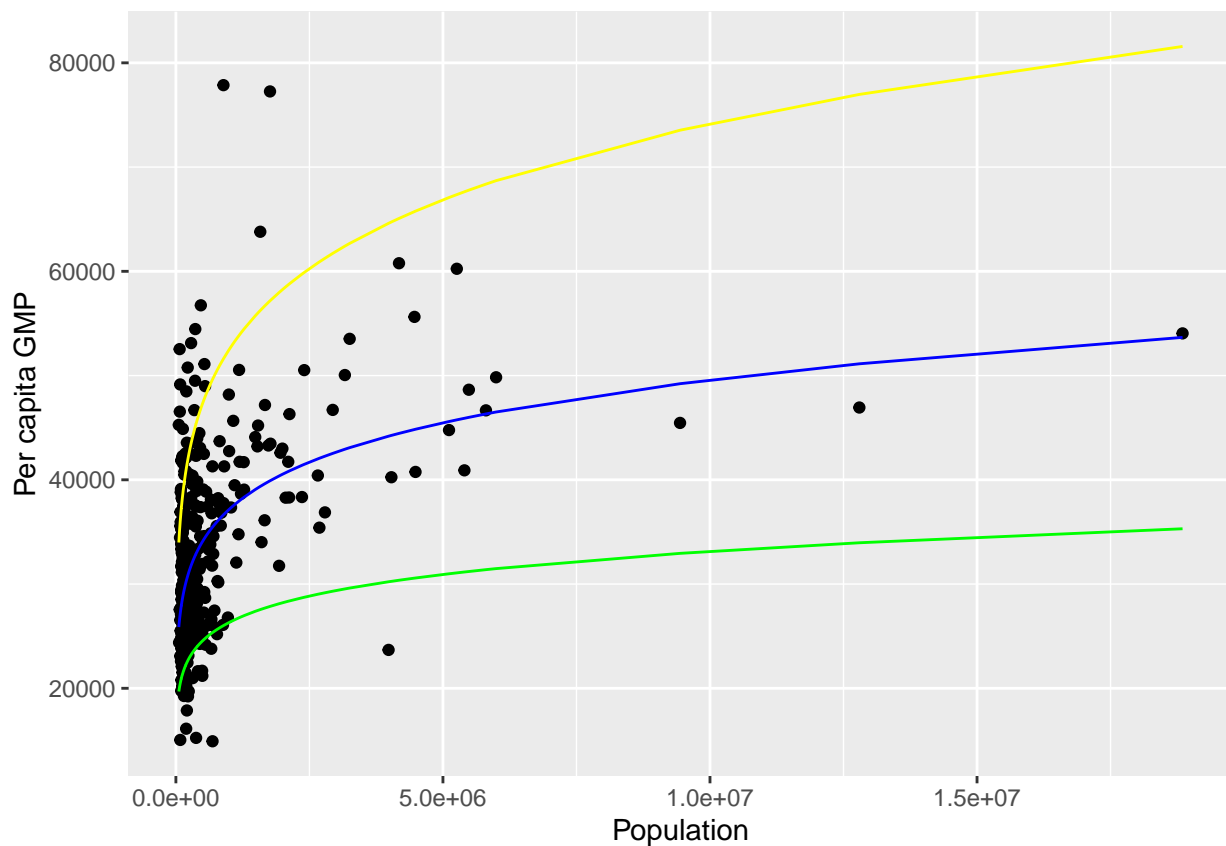
## [conflicted] Will prefer dplyr::filter over any other package
conflict_prefer("select", "dplyr")

## [conflicted] Will prefer dplyr::select over any other package
```

1. First, plot the data as in lecture, with per capita GMP on the y-axis and population on the x-axis. Add the curve function with the default values provided in lecture. Add two more curves corresponding to $a = 0.1$ and $a = 0.15$; use the `col` option to give each curve a different color (of your choice).

The plot code is as follows.

```
gmp %>% ggplot() + geom_point(aes(x = pop, y = pcgmp)) +
  labs(x = "Population", y = "Per capita GMP") +
  geom_line(aes(x = pop, y = 6611*(pop)^(1/8)), col = "blue") +
  geom_line(aes(x = pop, y = 6611*(pop)^(0.1)), col = "green") +
  geom_line(aes(x = pop, y = 6611*(pop)^(0.15)), col = "yellow")
```



2. Write a function, called `mse()`, which calculates the mean squared error of the model on a given data set. `mse()` should take three arguments: a numeric vector of length two, the first component standing for y_0 and the second for a ; a numerical vector containing the values of N ; and a numerical vector containing the values of Y . The function should return a single numerical value. The latter two arguments should have as the default values the columns `pop` and `pcgmp` (respectively) from the `gmp` data frame from lecture. Your function may not use `for()` or any other loop. Check that, with the default data, you get the following values.

The `mse()` code is as follows. And after checking with the given example, the code is right.

```
mse <- function(para, N = gmp$pop, Y = gmp$pcgmp)
{
  return(mean((Y - para[1]*(N)^(para[2]))^2))
}
```

```
mse(c(6611,0.15))
```

```
## [1] 207057513
```

```
# [1] 207057513
```

```
mse(c(5000,0.10))
```

```
## [1] 298459914
```

```
# [1] 298459915
```

4. R has several built-in functions for optimization, which we will meet as we go through the course. One of the simplest is `nlm()`, or non-linear minimization. `nlm()` takes two required arguments: a function, and a starting value for that function. Run `nlm()` three times with your function `mse()` and three starting value pairs for y_0 and a as in

```
nlm(mse, c(y0=6611,a=1/8))
```

```
## $minimum
```

```
## [1] 61857060
```

```
##
```

```
## $estimate
```

```
## [1] 6611.0000000 0.1263177
```

```
##
```

```
## $gradient
```

```
## [1] 50.048639 -9.983778
```

```
##
```

```
## $code
```

```
## [1] 2
```

```
##
```

```
## $iterations
```

```
## [1] 3
```

```
nlm(mse, c(y0=5000,a=1/8),check.analyticals=FALSE)
```

```
## $minimum
```

```
## [1] 62521484
```

```
##
## $estimate
## [1] 5000.0000004    0.1475909
##
## $gradient
## [1] -1030.494171    -2.473593
##
## $code
## [1] 2
##
## $iterations
## [1] 6
```

```
nlm(mse, c(y0=6611,a=0.1))
```

```
## $minimum
## [1] 61857060
##
## $estimate
## [1] 6611.0000003    0.1263177
##
## $gradient
## [1] 50.04683 -166.46832
##
## $code
## [1] 2
##
## $iterations
## [1] 6
```

What do the quantities `minimum` and `estimate` represent? What values does it return for these?

`minimum` means the value of the estimated minimum of given function, such as, `mse()`. And `estimate` means the point at which the minimum value of given function, such as `mse()`, is obtained, for example, \hat{y}_0 and \hat{a} .

In the examples above, we have three starting value pairs for y_0 and a .

When $y_0=6611$, $a=1/8$, the `minimum` is 61857060, the `estimate` is $\hat{y}_0 = 6611.0000000$ and $\hat{a} = 0.1263177$.

When $y_0=5000$, $a=1/8$, the `minimum` is 62521484, the `estimate` is $\hat{y}_0 = 5000.0000004$ and $\hat{a} = 0.1475909$.

When $y_0=6611$, $a=0.1$, the `minimum` is 61857060, the `estimate` is $\hat{y}_0 = 6611.0000003$ and $\hat{a} = 0.1263177$.

5. Using `nlm()`, and the `mse()` function you wrote, write a function, `plm()`, which estimates the parameters y_0 and a of the model by minimizing the mean squared error. It should take the following arguments: an initial guess for y_0 ; an initial guess for a ; a vector containing the N values; a vector containing the Y values. All arguments except the initial guesses should have suitable default values. It should return a list with the following components: the final guess for y_0 ; the final guess for a ; the final value of the MSE. Your function must call those you wrote in earlier questions (it should not repeat their code), and the appropriate arguments to `plm()` should be passed on to them.

What parameter estimate do you get when starting from $y_0 = 6611$ and $a = 0.15$? From $y_0 = 5000$ and $a = 0.10$? If these are not the same, why do they differ? Which estimate has the lower MSE?

The `plm()` function is given as follows.

```
plm <- function(para, N=gmp$pop, Y=gmp$pcgmp)
{
  nlm_result <- nlm(mse, para)
  return(list(y0=nlm_result$estimate[1], a=nlm_result$estimate[2], minMSE=nlm_result$minimum))
}
```

When running the `plm()` with different `para`, we can get the following answer.

When $y_0=6611$, $a=0.15$, the minimum is 61857060, the estimate is $\hat{y}_0 = 6611$ and $\hat{a} = 0.1263182$.

When $y_0=5000$, $a=0.1$, the minimum is 62521484, the estimate is $\hat{y}_0 = 5000$ and $\hat{a} = 0.1475913$.

I think the reason is that there are different extreme points near the starting point. And the first pair of $y_0=6611$, $a=0.15$ has lower MSE.

```
plm(c(y0=6611,a=0.15))
```

```
## $y0
## [1] 6611
##
## $a
## [1] 0.1263182
##
## $minMSE
## [1] 61857060
```

```
plm(c(y0=5000,a=0.1))
```

```
## $y0
## [1] 5000
##
## $a
## [1] 0.1475913
##
## $minMSE
## [1] 62521484
```

7._Convince yourself the jackknife can work_.

a. Calculate the mean per-capita GMP across cities, and the standard error of this mean, using the built-in functions `mean()` and `sd()`, and the formula for the standard error of the mean you learned in your intro. stats. class (or looked up on Wikipedia...).

```
mean(gmp$pcgmp)
```

```
## [1] 32922.53
```

```
sd(gmp$pcgmp)/sqrt(length(gmp$pcgmp))
```

```
## [1] 481.9195
```

Given the data, we can get the mean value of per-capita GMP is 32922.53, and standard error of the mean is 481.9195.

The formula for the standard error is as following.

$$SD_of_mean = \sqrt{Var(\bar{X})} = \sqrt{Var(\frac{1}{n} \sum_{i=1}^n X_i)} = \frac{1}{\sqrt{n}} SD(X)$$
$$SD = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}}$$

b. Write a function which takes in an integer `i`, and calculate the mean per-capita GMP for every city *except* city number `i`.

The needed function is as following.

```
mean_except_i <- function(i, Y=gmp$pcgmp)
{
  return(mean(Y[-i]))
}
```

c. Using this function, create a vector, `jackknifed.means`, which has the mean per-capita GMP where every city is held out in turn. (You may use a `for` loop or `sapply()`.)

The vector, `jackknifed.means`, is obtained by following code.

```
jackknifed.means <- sapply(seq(1, length(gmp$pcgmp)), mean_except_i)
```

d. Using the vector `jackknifed.means`, calculate the jack-knife approximation to the standard error of the mean. How well does it match your answer from part (a)?

```
jack_mean <- mean(jackknifed.means)
n <- length(jackknifed.means)
jack_var <- (n-1) * sum((jackknifed.means-jack_mean)^2) / n
jack_sd <- sqrt(jack_var)
jack_sd
```

```
## [1] 481.9195
```

From the above, we can know the jack-knife approximation to the standard error of the mean is 481.9195, and the answer from part (a) is also 481.9195.

So the jackknife matches well.

8. Write a function, `plm.jackknife()`, to calculate jackknife standard errors for the parameters y_0 and a . It should take the same arguments as `plm()`, and return standard errors for both parameters. This function should call your `plm()` function repeatedly. What standard errors do you get for the two parameters?

```
plm.jackknife <- function(para, N=gmp$pop, Y=gmp$pcgmp)
{
  plm_except_i <- function(i, para, N, Y)
  {
    return(plm(para, N[-i], Y[-i])[c(1,2)])
  }
  jack_para <- sapply(seq(1, length(gmp$pcgmp)), plm_except_i, para=para, N=N, Y=Y)
  para_y0 <- sapply(seq(1,length(jack_para),2), m<-function(i, a){return(a[[i]])}, a=jack_para)
  para_a <- sapply(seq(2,length(jack_para),2), m<-function(i, a){return(a[[i]])}, a=jack_para)
  jack_sd <- function(jack_para)
  {
    jack_mean <- mean(jack_para)
    n <- length(jack_para)
    jack_var <- (n-1) * sum((jack_para-jack_mean)^2) / n
    jack_sd <- sqrt(jack_var)
    return(jack_sd)
  }
  para_y0_sd <- jack_sd(para_y0)
  para_a_sd <- jack_sd(para_a)
  return(c(para_y0_sd, para_a_sd))
}
plm.jackknife(c(6611,0.15))
```

```
## [1] 0 0
```

9. The file `gmp-2013.dat` contains measurements for for 2013. Load it, and use `plm()` and `plm.jackknife` to estimate the parameters of the model for 2013, and their standard errors. Have the parameters of the model changed significantly?

```
gmp2013 <- read.table("data/gmp-2013.dat")
gmp2013$pop <- round(gmp2013$gmp/gmp2013$pcgmp)
estimate2013 <- plm(c(6611,0.15), N=gmp2013$pop, Y=gmp2013$pcgmp)
para <- c(estimate2013[[1]], estimate2013[[2]])
estimate2013
```

```
## $y0
## [1] 6611
##
## $a
## [1] 0.1263182
##
## $minMSE
## [1] 61857060
```

```
plm.jackknife(para, N=gmp2013$pop, Y=gmp2013$pcgmp)
```

```
## [1] 0 0
```