

STAT 341: Tutorial 4 – Practice with Implicit Attributes

Friday January 31, 2020

Part I: Not all Estimating Equations Arise from Derivatives

With this example I want to make clear that although many of the estimating equations we encounter arise as derivatives from some objective function, the *don't have to*.

Consider the implicitly defined attribute $\theta = (\mu, \sigma)$ where μ and σ are respectively measures of center and spread in the population $\mathcal{P} = \{y_1, y_2, \dots, y_N\}$. Find $\hat{\theta} = (\hat{\mu}, \hat{\sigma})$, the solution to the following system of equations

$$\psi(\theta; \mathcal{P}) = \mathbf{0} \quad E(Y) = \mu$$

given by

$$\begin{aligned} \textcircled{1} &\rightarrow \left[\begin{array}{c} (\frac{1}{N} \sum_{u \in \mathcal{P}} y_u) - \mu \\ (\frac{1}{N} \sum_{u \in \mathcal{P}} y_u^2) - \mu^2 - \sigma^2 \end{array} \right] = \left[\begin{array}{c} 0 \\ 0 \end{array} \right] \\ \textcircled{2} &\rightarrow \end{aligned} \quad E[Y^2] = \mu^2 + \sigma^2$$

Note that this system of equations arose by equating sample moments with population moments. This method of estimation is referred to as the Method of Moments.

$$\text{From } \textcircled{1} \text{ we have } \frac{1}{N} \sum_{u \in \mathcal{P}} y_u = \mu \Rightarrow \mu = \bar{y}$$

Substituting $\mu = \bar{y}$ into $\textcircled{2}$ yields:

$$\left(\frac{1}{N} \sum_{u \in \mathcal{P}} y_u^2 \right) - \bar{y}^2 = \sigma^2$$

$$\frac{\sum_{u \in \mathcal{P}} y_u^2 - N \bar{y}^2}{N} = \sigma^2$$

$$\frac{\sum_{u \in \mathcal{P}} (y_u - \bar{y})^2}{N} = \sigma^2$$

$$\therefore \sigma = \sqrt{\frac{\sum_{u \in \mathcal{P}} (y_u - \bar{y})^2}{N}}$$

$$\therefore \hat{\theta} = (\hat{\mu}, \hat{\sigma})$$

$$= \left(\bar{y}, \sqrt{\frac{\sum (y_u - \bar{y})^2}{N}} \right)$$

Part II: Least Absolute Deviations Regression

In class we have talked a lot about estimating $\theta = (\alpha, \beta)$ the intercept and slope associated with the simple linear regression

$$y_u = \alpha + \beta(x_u - \bar{x}) + r_u$$

And we have done this in a variety of different ways by altering the objective function. One such possible objective function is the *absolute error loss* function which gives rise to **least absolute deviations (LAD)** regression:

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{u \in \mathcal{P}} |y_u - \alpha - \beta(x_u - \bar{x})|$$

lad (-) from L1 pack package

In this part we will use gradient descent to fit this model to the **Animals** data. But first we need to define **rho** and **gradient** functions in R which requires that we first compute the gradients by hand.

(a) Let

$$\rho(\theta; \mathcal{P}) = \sum_{u \in \mathcal{P}} |y_u - \alpha - \beta(x_u - \bar{x})| = \sum_{u \in \mathcal{P}} |r_u|$$

Calculate the gradient vector $\mathbf{g} = \nabla \rho(\theta; \mathcal{P})$.

$$\mathbf{g} = \nabla \rho(\theta; \mathcal{P}) = \begin{bmatrix} \frac{\partial \rho}{\partial \alpha} \\ \frac{\partial \rho}{\partial \beta} \end{bmatrix}$$

where $r_u = y_u - \alpha - \beta(x_u - \bar{x})$

where $\frac{\partial \rho}{\partial \alpha} = \sum_{u \in \mathcal{P}} \frac{\partial |r_u|}{\partial r_u} \times \frac{\partial r_u}{\partial \alpha} = \sum_{u \in \mathcal{P}} \frac{r_u}{|r_u|} \times (-1)$

Think of $f(x) = |x|$

as $f(x) = \sqrt{x^2}$

$$f'(x) = \frac{x}{|x|}$$

$$\approx - \sum_{u \in \mathcal{P}} \operatorname{sign}(r_u)$$

$$\operatorname{sign}(r_u) = \begin{cases} 1 & r_u > 0 \\ 0 & r_u = 0 \\ -1 & r_u < 0 \end{cases}$$

$$\begin{aligned} \frac{\partial \rho}{\partial \beta} &= \sum_{u \in \mathcal{P}} \frac{\partial |r_u|}{\partial r_u} \times \frac{\partial r_u}{\partial \beta} \\ &= - \sum_{u \in \mathcal{P}} \frac{r_u}{|r_u|} \times (x_u - \bar{x}) \end{aligned}$$

$$= - \sum_{u \in \mathcal{P}} \operatorname{sign}(r_u) \times (x_u - \bar{x})$$

$$\vec{g} = - \sum_{u \in \mathcal{P}} \operatorname{sign}(r_u) \begin{bmatrix} 1 \\ x_u - \bar{x} \end{bmatrix}$$

$$\frac{r_u}{|r_u|} = \begin{cases} 1 & r_u > 0 \\ \text{undef.} & r_u = 0 \\ -1 & r_u < 0 \end{cases}$$

$$= - \sum_{u \in \mathcal{P}} \operatorname{sign}(r_u) \vec{z}_u$$

where $\vec{z}_u = \begin{bmatrix} 1 \\ x_u - \bar{x} \end{bmatrix}$

- (b) Write *factory functions* `createLADrho(x,y)` and `createLADgradient(x,y)` which take in as input only the data and which return as output the least absolute deviations objective function and the corresponding gradient function, respectively

```
→ createLADrho <- function(x, y) {  
  ## local variable  
  xbar <- mean(x)  
  ## Return this function  
  function(theta) {  
    alpha <- theta[1]  
    beta <- theta[2]  
    sum(abs(y - alpha - beta * (x - xbar)))  
  }  
}  
  
→ createLADgradient <- function(x, y) {  
  ## local variables  
  xbar <- mean(x)  
  function(theta) {  
    alpha <- theta[1]  
    beta <- theta[2]  
    ru = y - alpha - beta * (x - xbar)  
    -1 * c(sum(sign(ru)), sum(sign(ru) * (x - xbar)))  
  }  
}
```

- (c) Using the `gradientDescent` function (from class) together with the `gridLineSearch` and `testConvergence` functions (from class) as well as `rho` and `gradient` functions created by your factory functions from part (b), find the LAD estimates $\hat{\theta} = (\hat{\alpha}, \hat{\beta})$ for the Animals data.

```
library(robustbase)
rho <- createLADrho(x = log(Animals2$body), y = log(Animals2$brain))
g <- createLADgradient(x = log(Animals2$body), y = log(Animals2$brain))

res.manual <- gradientDescent(theta = c(0, 0), rhoFn = rho, gradientFn = g,
  lineSearchFn = gridLineSearch, testConvergenceFn = testConvergence,
  maxIterations = 5000, tolerance = 1e-20, relative = TRUE)

print(res.manual)
```

\$theta
[1] 3.3589387 0.7368431

\$converged
[1] TRUE

\$iteration
[1] 26

\$fnValue
[1] 48.01975

Handwritten notes:
 - Red arrows point from $\hat{\alpha}$ to 3.3589387 and from $\hat{\beta}$ to 0.7368431.
 - A red arrow points from $\hat{\theta}$ to the entire output line for \$theta.
 - A red arrow points from "min of p at $\hat{\theta}$ " to the \$fnValue output line.

- (d) In class we performed LAD regression using the `lad` function from the `L1pack` package. Let's confirm that what we found above agrees with the output of this other function.

```
library(L1pack)
res.L1pack <- lad(log(Animals2$brain) ~ I(log(Animals2$body) - mean(log(Animals2$body))))
print(res.L1pack)
```

Call:
lad(formula = log(Animals2\$brain) ~ I(log(Animals2\$body) - mean(log(Animals2\$body))))
Converged in 9 iterations

Coefficients:
(Intercept)
3.3541
I(log(Animals2\$body) - mean(log(Animals2\$body)))
0.7373

Degrees of freedom: 65 total; 63 residual
Scale estimate: 1.04454

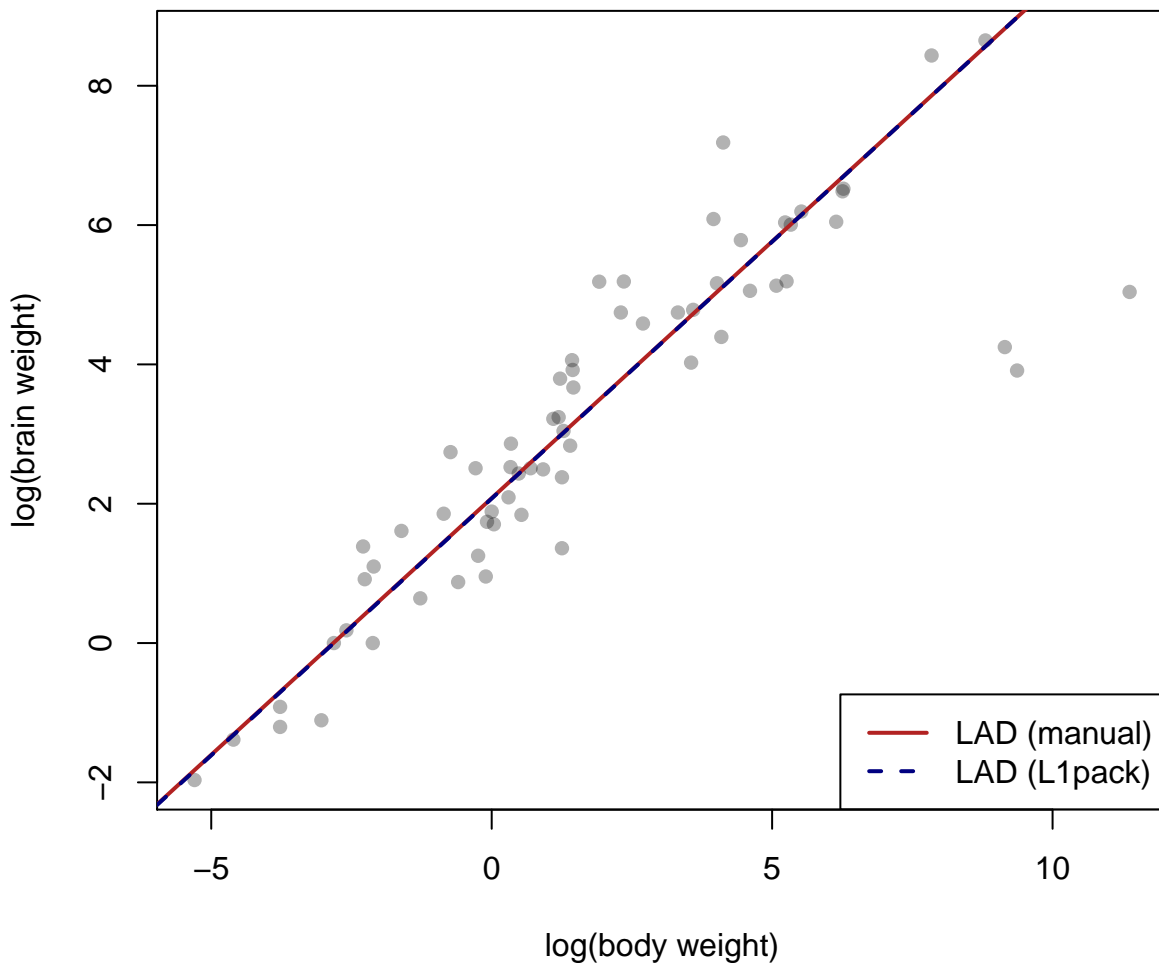
- (e) Construct a scatter plot of $\log(\text{brain weight})$ versus $\log(\text{body weight})$ for the Animals data and plot both lines of best fit – the one we determined manually and the one calculated using `lad` – to see if the difference is material. Use a legend to distinguish among the lines.

```
plot(x = log(Animals2$body), y = log(Animals2$brain), main = "", xlab = "log(body weight)",
     ylab = "log(brain weight)", pch = 16, col = adjustcolor("black", alpha.f = 0.3))

abline(a = res.manual$theta[1] - res.manual$theta[2] * mean(log(Animals2$body)),
       b = res.manual$theta[2], col = "firebrick", lwd = 2, lty = 1)

abline(a = res.L1pack$coef[1] - res.L1pack$coef[2] * mean(log(Animals2$body)),
       b = res.L1pack$coef[2], col = "navyblue", lwd = 2, lty = 2)

legend("bottomright", legend = c("LAD (manual)", "LAD (L1pack)"), col = c("firebrick",
                                "navyblue"), lty = 1:2, lwd = 2)
```



(f) Can we apply the Newton-Raphson Method to this problem? In other words, can we determine the LAD estimate $\hat{\theta} = (\hat{\alpha}, \hat{\beta})$ via the Newton-Raphson Method? **NO!!**

$$\vec{g} = \begin{bmatrix} \frac{\partial p}{\partial \alpha} \\ \frac{\partial p}{\partial \beta} \end{bmatrix} = \vec{\Psi}(\vec{\theta}; P) = \begin{bmatrix} \psi_1 \\ \psi_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

We want $\hat{\theta}$ that solves this

We also need $\vec{\Psi}'(\vec{\theta}; P) = \begin{bmatrix} \frac{\partial \psi_1}{\partial \alpha} & \frac{\partial \psi_1}{\partial \beta} \\ \frac{\partial \psi_2}{\partial \alpha} & \frac{\partial \psi_2}{\partial \beta} \end{bmatrix}$

$= \begin{bmatrix} \frac{\partial^2 p}{\partial \alpha^2} & \frac{\partial^2 p}{\partial \alpha \partial \beta} \\ \frac{\partial^2 p}{\partial \beta \partial \alpha} & \frac{\partial^2 p}{\partial \beta^2} \end{bmatrix}$ ← Hessian of $p(\vec{\theta}; P)$

$= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ Since $\frac{\partial p}{\partial \alpha}$ and $\frac{\partial p}{\partial \beta}$ are both constants with respect to α and β .

We need $[\vec{\Psi}'(\vec{\theta}; P)]^{-1}$. Since this does not exist we can't apply NR.