

Applying Matrix Rules to Derive the OLS estimator

To get the OLS estimator $\hat{\beta}_{ols} = (X'X)^{-1}X'Y$ we find the parameter vector $\beta = [\beta_1, \beta_2, \dots, \beta_k]$ which minimizes the squared error $\sum_{i=1}^N u_i^2 = \mathbf{u}'\mathbf{u} = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$. This is the objective function, $f(\beta)$, of our optimization problem. The first order conditions to solve this problem require setting the partial derivatives of f with respect to $\beta_1, \beta_2, \dots, \beta_k$ equal to zero. In vector notation this means

$$\frac{df}{d\beta} \equiv \left[\frac{\partial f}{\partial \beta_1}, \frac{\partial f}{\partial \beta_2}, \dots, \frac{\partial f}{\partial \beta_k} \right]' = \mathbf{0}, \text{ where } \mathbf{0} \equiv [0, 0, \dots, 0]' \text{ is the zero vector.}$$

Note that according to general convention $\frac{df}{d\beta}$ is a column vector with the partial derivatives arranged in a way that the i th component of β , β_i , corresponds to the i th component of $\frac{df}{d\beta}$, $\frac{df}{d\beta_i}$. Also, the dimension of $\frac{df}{d\beta}$ is the same as that of β . This is true, in general, when taking the derivative of a scalar valued function with respect to a vector. Alternatively, if we had a vector valued function that evaluated to an $N \times 1$ vector, differentiating with respect to an $M \times 1$ vector would result in an $M \times N$ matrix called the Jacobian (or $N \times M$ according to another convention). Fortunately, this is not a concern for this problem because $f(\beta) = \sum_{i=1}^N u_i^2$ is a scalar. When this is so, $\frac{df}{d\beta}$ is also called the gradient vector. An alternative notation for the gradient is $\nabla f = [f_1, f_2, \dots, f_k]'$

Before differentiation, it's useful to simplify our objective function. We do so, first, by applying the distributive property of matrices (i.e. $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{A}\mathbf{B} + \mathbf{A}\mathbf{C}$ which is not in general equal to $\mathbf{B}\mathbf{A} + \mathbf{C}\mathbf{A}$ because matrix multiplication does NOT obey a commutative property).

$$f(\beta) = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{y}'\mathbf{y} - \beta'X'\mathbf{y} - \mathbf{y}'X\beta + \beta'X'X\beta$$

Two other matrix rules which were applied for the above manipulation are $(\mathbf{A} + \mathbf{B})' = \mathbf{A}' + \mathbf{B}'$ and $(\mathbf{A}\mathbf{B})' = \mathbf{B}'\mathbf{A}'$. We can see the term $\mathbf{y}'\mathbf{y}$ has no β in it and so it is not relevant for the optimization problem. Another simplification can be made by observing $\mathbf{y}'X\beta = (\mathbf{y}'X\beta)' = \beta'X'\mathbf{y}$. This is, of course, not true in general of transposes, but follows in this case because this matrix product evaluates to a scalar (and the transpose of a scalar is the same scalar). Now our optimization problem simplifies to.

$$\min_{\beta} -2\beta'X'\mathbf{y} + \beta'X'X\beta$$

We can already see the $X'\mathbf{y}$ and the $X'X$ from the $\hat{\beta}_{ols}$ formula in the objective function above. But because the required vector differentiation may still be somewhat involved for the inexperienced, we will show differentiation for each of the above terms separately.

1. Differentiating $2\beta'X'y$ with respect to β

By letting $g = \beta'X'y$ we can easily see that

$$\frac{d2\beta'X'y}{d\beta} = \begin{bmatrix} \frac{\partial(2g)}{\partial\beta_1} \\ \vdots \\ \frac{\partial(2g)}{\partial\beta_k} \end{bmatrix} = \begin{bmatrix} 2\frac{\partial g}{\partial\beta_1} \\ \vdots \\ 2\frac{\partial g}{\partial\beta_k} \end{bmatrix} = 2 \begin{bmatrix} \frac{\partial g}{\partial\beta_1} \\ \vdots \\ \frac{\partial g}{\partial\beta_k} \end{bmatrix} = 2 \frac{d\beta'X'y}{d\beta}$$

You should be able to see that this is a general property of vector differentiation so that $\frac{\partial c.A}{\partial B} = c \frac{\partial A}{\partial B}$ for any scalar c , and vectors B and $A = A(B)$.

Next it useful to simplify notation by letting the $k \times 1$ vector $X'y = s = [s_1, s_2, \dots, s_k]'$. This gives us

$$\beta'X'y = \beta's = [\beta_1, \dots, \beta_k] \begin{bmatrix} s_1 \\ \vdots \\ s_k \end{bmatrix} = \beta_1 s_1 + \beta_2 s_2 + \dots + \beta_k s_k$$

Now we can easily differentiate by computing the partial derivative components of the gradient

$$\frac{d\beta's}{d\beta} \equiv \left[\frac{\partial\beta's}{\partial\beta_1}, \dots, \frac{\partial\beta's}{\partial\beta_k} \right]' = [s_1, s_2, \dots, s_k]' = X'y$$

Similarly, we can show for any two column vectors of the same dimension that $\frac{dA'B}{dB} = \frac{dB'A}{dB} = A$

Although, not necessary to derive $\hat{\beta}_{ols}$ we can see that the components of s work out to the sums

$$X'y = \begin{bmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{N1} & \dots & x_{Nk} \end{bmatrix}' \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_{11} & \dots & x_{N1} \\ \vdots & \ddots & \vdots \\ x_{1k} & \dots & x_{Nk} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_{11}y_1 + \dots + x_{N1}y_N \\ \vdots \\ x_{1k}y_1 + \dots + x_{Nk}y_N \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N x_{i1}y_i \\ \vdots \\ \sum_{i=1}^N x_{ik}y_i \end{bmatrix} = \begin{bmatrix} s_1 \\ \vdots \\ s_k \end{bmatrix}$$

2. Differentiating $\beta'X'X\beta$ with respect to β

Again to simplify notation we let the $k \times k$ matrix $X'X = M$ with components defined as

$$X'X = \begin{bmatrix} x_{11} & \dots & x_{N1} \\ \vdots & \ddots & \vdots \\ x_{1k} & \dots & x_{Nk} \end{bmatrix} \begin{bmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{N1} & \dots & x_{Nk} \end{bmatrix}' = \begin{bmatrix} \sum_{i=1}^N x_{i1}x_{i1} & \dots & \sum_{i=1}^N x_{i1}x_{ik} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_{ik}x_{i1} & \dots & \sum_{i=1}^N x_{ik}x_{ik} \end{bmatrix} = \begin{bmatrix} m_{11} & \dots & m_{1k} \\ \vdots & \ddots & \vdots \\ m_{k1} & \dots & m_{kk} \end{bmatrix}$$

Then $\beta'X'X\beta$ works out to

$$\begin{aligned}
\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} &= [\beta_1, \dots, \beta_k] \begin{bmatrix} m_{11} & \cdots & m_{1k} \\ \vdots & \ddots & \vdots \\ m_{k1} & \cdots & m_{kk} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} = [\beta_1, \dots, \beta_k] \begin{bmatrix} m_{11}\beta_1 + \cdots + m_{1k}\beta_k \\ \vdots \\ m_{k1}\beta_k + \cdots + m_{kk}\beta_k \end{bmatrix} \\
&= m_{11}\beta_1^2 + \beta_1 \sum_{i=1, i \neq 1}^k m_{1i}\beta_i \\
&\quad + m_{22}\beta_2^2 + \beta_2 \sum_{i=1, i \neq 2}^k m_{2i}\beta_i \\
&\quad + \cdots \\
&\quad + m_{kk}\beta_k^2 + \beta_k \sum_{i=1, i \neq k}^k m_{ki}\beta_i
\end{aligned}$$

Where the notation $\sum_{i=1, i \neq j}^k a_i$ means to take the sum of all a_1, \dots, a_n excluding a_j . For clarity, when $k = 3$

$$\begin{aligned}
\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} &= m_{11}\beta_1^2 + \beta_1(m_{12}\beta_2 + m_{13}\beta_3) + m_{22}\beta_2^2 + \beta_2(m_{21}\beta_1 + m_{23}\beta_3) + m_{33}\beta_3^2 + \beta_3(m_{31}\beta_1 + m_{32}\beta_2) \\
&= (m_{11}\beta_1^2 + m_{12}\beta_1\beta_2 + m_{13}\beta_1\beta_3) + (m_{21}\beta_2\beta_1 + m_{22}\beta_2^2 + m_{23}\beta_2\beta_3) + (m_{31}\beta_3\beta_1 + m_{32}\beta_3\beta_2 + m_{33}\beta_3^2)
\end{aligned}$$

The brackets in the last line are present to help you see the pattern in m_{ij} coefficients and β_i, β_j pairs.

$$\text{Then } \frac{\partial \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{\partial \beta_1} = 2m_{11}\beta_1 + m_{12}\beta_2 + m_{13}\beta_3 + m_{21}\beta_2 + m_{31}\beta_3 = 2(m_{11}\beta_1 + m_{12}\beta_2 + m_{13}\beta_3)$$

The last equality follows from the symmetry of $\mathbf{X}'\mathbf{X} = \mathbf{M}$ meaning $m_{12} = m_{21}$ and $m_{13} = m_{31}$.

We can show the same for the remaining partial derivatives and generalize for any k , giving us

$$\frac{d\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{d\boldsymbol{\beta}} = \begin{bmatrix} 2(m_{11}\beta_1 + m_{12}\beta_2 + \cdots + m_{1k}\beta_k) \\ 2(m_{21}\beta_1 + m_{22}\beta_2 + \cdots + m_{2k}\beta_k) \\ \vdots \\ 2(m_{k1}\beta_1 + m_{k2}\beta_2 + \cdots + m_{kk}\beta_k) \end{bmatrix} = \mathbf{2M}\boldsymbol{\beta} = \mathbf{2X}'\mathbf{X}\boldsymbol{\beta}$$

This rule of vector differentiation can be also generalized for any $N \times 1$ column vector \mathbf{B} , and $N \times N$ symmetric matrix \mathbf{A} , so that $\frac{d\mathbf{B}'\mathbf{A}\mathbf{B}}{d\mathbf{B}} = \mathbf{2AB}$. Matrix products of the form $\mathbf{B}'\mathbf{A}\mathbf{B}$ are called quadratic forms.

3) Putting it all together

From sections 1) and 2) we see that are first order conditions becomes $\frac{df}{d\boldsymbol{\beta}} = -\mathbf{2X}'\mathbf{y} + \mathbf{2X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$

Dividing both sides by 2 and rearranging we have $\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\mathbf{y}$. Multiply both sides by $(\mathbf{X}'\mathbf{X})^{-1}$ to get

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