

Projection Matrices and Regression Sums of Squares

The projection matrices \mathbf{P} and \mathbf{M}

The multiple regression model for the population, expressed in matrix form, is $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$

Expressed in terms of the estimates, it is $\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}$ (Equation 1)

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

Equation 1 can be written as $\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{u}}$, where $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$

$\hat{\mathbf{y}}$ can be rewritten as follows:

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

Define $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$, so $\hat{\mathbf{y}} = \mathbf{P}\mathbf{y}$

Kutner et al. call this matrix \mathbf{H} , the "hat matrix", because it "puts the hat on" \mathbf{y} .

Properties of the \mathbf{P} matrix

\mathbf{P} depends only on \mathbf{X} , not on \mathbf{y} . In some derivations, we may need different \mathbf{P} matrices that depend on different sets of variables. In this case, the matrix \mathbf{P} may be written with a subscript, indicating the \mathbf{X} matrix that it depends on, as \mathbf{P}_X . In order to reduce the complexity of the notation, I won't do this unless I need to distinguish different \mathbf{P} matrices.

If \mathbf{X} has n rows and k columns, \mathbf{P} is n by n , a square matrix. Note that n is typically much larger than k , so \mathbf{P} is a large matrix. The \mathbf{P} matrix has great theoretical importance, but it is not usually computed.

Property 1 – Multiplying by \mathbf{P} leaves the \mathbf{X} matrix unchanged, $\mathbf{P}*\mathbf{X}=\mathbf{X}$

Proof:

$$\mathbf{P}\mathbf{X} = (\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}')\mathbf{X} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{X}) = \mathbf{X}$$

Property 2 – \mathbf{P} is idempotent, meaning $\mathbf{P}*\mathbf{P}=\mathbf{P}$.

Proof:

$$\begin{aligned} \mathbf{P}*\mathbf{P} &= (\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}')(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}') = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' = \mathbf{P} \end{aligned}$$

Property 3 – \mathbf{P} is symmetric, $\mathbf{P}' = \mathbf{P}$

Proof:

$$\mathbf{P}' = (\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')' = (\mathbf{X}')'[(\mathbf{X}'\mathbf{X})^{-1}]' \mathbf{X}' = \mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}] \mathbf{X}' = \mathbf{P}$$

The central equality, $[(\mathbf{X}'\mathbf{X})^{-1}]' = [(\mathbf{X}'\mathbf{X})^{-1}]$ is true because $\mathbf{X}'\mathbf{X}$ is symmetric, and the inverse of a symmetric matrix is symmetric.

The \mathbf{M} matrix, the residualizer

In the formula $\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{u}}$, I have showed a matrix expression for $\hat{\mathbf{y}}$ in terms of the matrix \mathbf{P} .

There is also an expression for $\hat{\mathbf{u}}$.

$$\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{P}\mathbf{y} = (\mathbf{I} - \mathbf{P})\mathbf{y}$$

The matrix $(\mathbf{I} - \mathbf{P})$ is called the residualizer, because it makes the residuals. It is usually given the symbol \mathbf{M} . So now we have $\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{u}} = \mathbf{P}\mathbf{y} + \mathbf{M}\mathbf{y}$. This decomposes \mathbf{y} into two chunks.

Properties of the \mathbf{M} matrix

Property 1 – Multiplying \mathbf{M} by \mathbf{X} gives a $\mathbf{0}$ matrix, $\mathbf{M}\mathbf{X}=\mathbf{0}$

$$\text{Proof: } \mathbf{M}\mathbf{X} = (\mathbf{I} - \mathbf{P})\mathbf{X} = \mathbf{X} - \mathbf{P}\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}$$

Property 2 – like \mathbf{P} , the matrix \mathbf{M} is idempotent.

$$\begin{aligned} \text{Proof: } \mathbf{M} * \mathbf{M} &= (\mathbf{I} - \mathbf{P})(\mathbf{I} - \mathbf{P}) = \mathbf{I} * \mathbf{I} - \mathbf{I} * \mathbf{P} - \mathbf{P} * \mathbf{I} + \mathbf{P} * \mathbf{P} \\ &= \mathbf{I} - \mathbf{P} - \mathbf{P} + \mathbf{P} = \mathbf{I} - \mathbf{P} = \mathbf{M} \end{aligned}$$

Property 3 – like \mathbf{P} , the matrix \mathbf{M} is symmetric

$$\text{Proof: } \mathbf{M}' = (\mathbf{I} - \mathbf{P})' = (\mathbf{I}' - \mathbf{P}') = (\mathbf{I} - \mathbf{P}) = \mathbf{M}$$

Property 4 – $\mathbf{P} * \mathbf{M} = \mathbf{0}$, where $\mathbf{0}$ represents a matrix with all 0 elements.

$$\text{Proof: } \mathbf{P} * \mathbf{M} = \mathbf{P}(\mathbf{I} - \mathbf{P}) = \mathbf{P} - \mathbf{P} * \mathbf{P} = \mathbf{P} - \mathbf{P} = \mathbf{0}$$

Property 4a – Property 4 means that for any vector \mathbf{y} , $\mathbf{P}\mathbf{y}$ is orthogonal to $\mathbf{M}\mathbf{y}$.

Proof: Recall that orthogonal vectors have dot product = 0, and that the dot product can be expressed as a transpose product.

$$(\mathbf{Py})'(\mathbf{My}) = \mathbf{y}'\mathbf{P}'\mathbf{My} = \mathbf{y}'\mathbf{PM}\mathbf{y} = \mathbf{y}'(\mathbf{0})\mathbf{y} = \mathbf{0}$$

Now $\mathbf{Py} = \hat{\mathbf{y}}$ and $\mathbf{My} = \hat{\mathbf{u}}$, so we have just proved that $\hat{\mathbf{y}} \cdot \hat{\mathbf{u}} = 0$. This is the algebraic equivalent of the regression geometry.

Sums of Squares

We can use the \mathbf{P} and \mathbf{M} matrices to prove facts about the regression sums of squares. We need some matrix definitions and facts.

The trace of a matrix

The trace of a matrix is the sum of its diagonal elements. In the matrix

$$\begin{pmatrix} 3 & 2 & 0 \\ 0 & -2 & 1 \\ -1 & 5 & 4 \end{pmatrix}, \text{ the trace is } 3 + (-2) + 4 = 5.$$

I will use \mathbf{tr} to indicate the trace of a matrix.

Trace property 1 – If \mathbf{I} is an n by n identity matrix (sometimes denoted by \mathbf{I}_n), then $\mathbf{tr}(\mathbf{I}_n) = n$.

Proof: \mathbf{I}_n has n 1's along the diagonal.

Matrix fact: Note that in general, two matrices \mathbf{A} and \mathbf{B} may not be conformable for multiplication in both orders, i.e. \mathbf{AB} and \mathbf{BA} may not both exist. If they do, and if \mathbf{A} is n by m , then \mathbf{B} must be m by n . If $m \neq n$, the two products will have different dimensions, \mathbf{AB} will be n by n , and \mathbf{BA} will be m by m .

Trace property 2 – If \mathbf{A} and \mathbf{B} are conformable for multiplication in both orders then $\mathbf{tr}(\mathbf{AB}) = \mathbf{tr}(\mathbf{BA})$. This is true even if \mathbf{AB} does not have the same dimensions as \mathbf{BA} .

Proof: Suppose \mathbf{A} is n by m , \mathbf{B} is m by n .

$$\mathbf{tr}(\mathbf{AB}) = \sum_{i=1}^n (\mathbf{AB})_{ii} = \sum_{i=1}^n \sum_{k=1}^m \mathbf{A}_{ik} \mathbf{B}_{ki} \quad \text{Formula 1}$$

$$\mathbf{tr}(\mathbf{BA}) = \sum_{i=1}^m (\mathbf{BA})_{ii} = \sum_{i=1}^m \sum_{k=1}^n \mathbf{B}_{ik} \mathbf{A}_{ki} \quad \text{Formula 2}$$

Looking at the 2 formulas, we can get formula 2 from formula 1 by a) switching the order of summation; b) switching the names of the indices $i \leftrightarrow k$; and c) reversing the order of \mathbf{A} and \mathbf{B} . This last operation is permitted because the products are of individual elements of \mathbf{A} and \mathbf{B} , not the whole matrix, so multiplication is commutative. So Formula 1 = Formula 2, and $\mathbf{tr}(\mathbf{AB}) = \mathbf{tr}(\mathbf{BA})$.

Trace property 3 – If \mathbf{A} and \mathbf{B} are conformable for addition, $\text{tr}(\mathbf{A}+\mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$ and $\text{tr}(\mathbf{A}-\mathbf{B}) = \text{tr}(\mathbf{A}) - \text{tr}(\mathbf{B})$.

Proof: These can be seen easily from the definitions of matrix addition and the trace.

Trace property 4 – Suppose that matrix \mathbf{X} is n by k , and matrix $\mathbf{X}'\mathbf{X}$ has an inverse.

Then $\text{tr}(\mathbf{P}_\mathbf{X}) = k$.

Proof:

$\text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') = \text{tr}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}) = \text{tr}(\mathbf{I}_k) = k$. Here \mathbf{I}_k is the k by k identity matrix, which has k 1's along the diagonal.

Trace property 5 – $\text{tr}(\mathbf{M}) = n - k$.

Proof:

$$\text{tr}(\mathbf{M}) = \text{tr}(\mathbf{I}_n - \mathbf{P}) = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{P}) = n - k$$

Expectation of sum of squares of residuals

Now we have the tools to show that if $\text{var}(u) = \sigma^2$, then the expectation of the sum of squares of the residuals is $(n - k)\sigma^2$.

First, we'll look at the sum of squares of the residuals. This can be expressed as $\sum_{i=1}^n \hat{u}_i^2 = \hat{\mathbf{u}}' \hat{\mathbf{u}}$,

and rewritten as

$$\hat{\mathbf{u}}' \hat{\mathbf{u}} = (\mathbf{M}\mathbf{y})' \mathbf{M}\mathbf{y} = \mathbf{y}' \mathbf{M}\mathbf{y} = (\mathbf{X}\boldsymbol{\beta} + \mathbf{u})' \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \mathbf{u})$$

Recall that $\mathbf{M}\mathbf{X}=\mathbf{0}$. In the last expression, there are 4 terms when we multiply it out:

$$(\mathbf{X}\boldsymbol{\beta} + \mathbf{u})' \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \mathbf{u}) = \boldsymbol{\beta}' \mathbf{X}' \mathbf{M}\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}' \mathbf{X}' \mathbf{M}\mathbf{u} + \mathbf{u}' \mathbf{M}\mathbf{X}\boldsymbol{\beta} + \mathbf{u}' \mathbf{M}\mathbf{u}.$$

All except the last term have either $\mathbf{M}\mathbf{X}$ or $\mathbf{X}'\mathbf{M}$ or both. But $\mathbf{X}'\mathbf{M}=(\mathbf{M}\mathbf{X})'=\mathbf{0}$, so all terms except the last drop out. Putting it together, $\hat{\mathbf{u}}' \hat{\mathbf{u}} = \mathbf{u}' \mathbf{M}\mathbf{u}$.

Now, \mathbf{u} is a vector of the errors. We don't know the actual values of \mathbf{u} , but we're still allowed to do algebra with the vector. In particular, $\text{tr}(\mathbf{u}' \mathbf{M}\mathbf{u}) = \text{tr}(\mathbf{M}\mathbf{u}\mathbf{u}')$. Since we have showed the matrix expression, $\hat{\mathbf{u}}' \hat{\mathbf{u}} = \mathbf{u}' \mathbf{M}\mathbf{u}$, we must have $\text{tr}(\hat{\mathbf{u}}' \hat{\mathbf{u}}) = \text{tr}(\mathbf{u}' \mathbf{M}\mathbf{u}) = \text{tr}(\mathbf{M}\mathbf{u}\mathbf{u}')$.

If you look carefully at this last set of equations, you may notice something strange. On the left side is $\text{tr}(\hat{\mathbf{u}}'\hat{\mathbf{u}})$. But $\hat{\mathbf{u}}$ is an n by 1 vector, so $\hat{\mathbf{u}}'\hat{\mathbf{u}}$ is a 1 by 1 matrix, a scalar. Its trace is simply its value. On the other side, $\mathbf{M}\mathbf{u}\mathbf{u}'$ is an n by n matrix multiplied by a n by 1 vector multiplied by a 1 by n vector. The whole product is n by n . But we're not saying the matrices on the two sides are equal, just their traces. Now we're ready to show that the expectation of the sum of squares of the residuals is $(n-k)\sigma^2$.

We have $\text{tr}(\hat{\mathbf{u}}'\hat{\mathbf{u}}) = \text{tr}(\mathbf{M}\mathbf{u}\mathbf{u}')$, so $E(\text{tr}(\hat{\mathbf{u}}'\hat{\mathbf{u}})) = E(\text{tr}(\mathbf{M}\mathbf{u}\mathbf{u}'))$. Now a trace is just a fancy way of doing a sum of terms, so we can move the expectation inside the trace operation. In addition, the matrix \mathbf{M} is just a function of the \mathbf{X} matrix, and we are treating \mathbf{X} as non-random, so we can move the expectation inside the product with \mathbf{M} , giving

$$E(\text{tr}(\hat{\mathbf{u}}'\hat{\mathbf{u}})) = \text{tr}(\mathbf{M} * E(\mathbf{u}\mathbf{u}'))$$

Now we'll look at $E(\mathbf{u}\mathbf{u}')$. Recall that $\mathbf{u}\mathbf{u}'$ is an n by n square matrix. On the diagonal will be terms like u_i^2 . The off diagonal elements will be terms like $u_i u_j$. Assuming that the errors are all independent and identically distributed, with variance σ^2 , $E(u_i^2) = \sigma^2$ and $E(u_i u_j) = 0$.

So $E(\mathbf{u}\mathbf{u}')$ will be a matrix with σ^2 for each diagonal element and 0 for each off-diagonal element. In other words, $E(\mathbf{u}\mathbf{u}') = \sigma^2 \mathbf{I}_n$.

Finally, $\mathbf{M} * E(\mathbf{u}\mathbf{u}') = \mathbf{M}(\sigma^2 \mathbf{I}_n) = \sigma^2 \mathbf{M}$, so $\text{tr}(\mathbf{M} * E(\mathbf{u}\mathbf{u}')) = \text{tr}(\sigma^2 \mathbf{M}) = \sigma^2 \text{tr}(\mathbf{M}) = (n-k)\sigma^2$.

If you go back over the steps, you'll see that we've showed that $E(\text{SSR}) = (n-k)\sigma^2$, where SSR is the sum of squares of the residuals. That means that $\text{SSR}/(n-k)$ is an unbiased estimator of σ^2 .

Expectation of model sum of squares

Consider the regression model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$. Now assume that the true $\boldsymbol{\beta}$ is $\mathbf{0}$ (a vector where all components are 0). Then all that's left of the regression model is $\mathbf{y} = \mathbf{u}$. In this case, what's the expectation of the model (regression) sum of squares?

First, estimate the model and get $\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}$. Note that although $\boldsymbol{\beta} = \mathbf{0}$, in general $\hat{\boldsymbol{\beta}} \neq \mathbf{0}$.

Second, calculate the model sum of squares:

$$\text{Model sum of squares} = \sum_{i=1}^n \hat{y}_i^2 = \hat{\mathbf{y}}'\hat{\mathbf{y}}$$

Next, plug in $\hat{y} = \mathbf{P}y$. Model sum of squares = $\mathbf{y}'\mathbf{P}'\mathbf{P}y = \mathbf{y}'\mathbf{P}y = \mathbf{u}'\mathbf{P}u$.

Note that $\mathbf{u}'\mathbf{u}$ is the sum of squares of the errors, which has expectation $n\sigma^2$.

Also, $\mathbf{u}'\mathbf{u} = \mathbf{u}'\mathbf{I}u = \mathbf{u}'(\mathbf{P} + \mathbf{M})u = \mathbf{u}'\mathbf{P}u + \mathbf{u}'\mathbf{M}u$.

So $\mathbf{u}'\mathbf{P}u = \mathbf{u}'\mathbf{u} - \mathbf{u}'\mathbf{M}u$, and $E(\mathbf{u}'\mathbf{P}u) = E(\mathbf{u}'\mathbf{u}) - (\mathbf{u}'\mathbf{M}u) = n\sigma^2 - (n - k)\sigma^2 = k\sigma^2$

In words, the expectation of the model sum of squares is the number of regressors times the variance of the errors, if the complete null hypothesis, $\boldsymbol{\beta}=\mathbf{0}$, is true.

Sums of squares when β_0 is not 0

The derivation given above assumes that all the β 's are 0, including β_0 . More realistically, we are interested in the situation where β_0 may be non-zero, and want to test all the other β 's. Note that the derivation of the expectation of the residual sum of squares did not assume anything about the value of the β 's, so the expectation of the residual sum of squares is still $(n - k)\sigma^2$, where k is the total number of regressors, including the constant. With some additional work (not shown here), we can show that the effect of β_0 can be eliminated if we make the model sum of squares be the sum of deviations of the predicted values from the overall mean, rather than simply the sum of squares of predicted values. In this case, the expectation of the model sum of squares is $(k - 1)\sigma^2$. Notice that $k - 1$ is simply the number of regressors, excluding the constant.