

Lecture 2: Three Justifications for OLS

Dennis Sun

June 24, 2015

Today, we will study three derivations of the OLS estimator

$$\hat{\boldsymbol{\beta}}^{\text{OLS}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\mathbf{y} - X\boldsymbol{\beta}\|^2 = (X^T X)^{-1} X^T \mathbf{y}. \quad (1)$$

To save space, we will use the shorthand $X^+ \stackrel{\text{def}}{=} (X^T X)^{-1} X^T$. X^+ is called the pseudoinverse of X because it satisfies $X^+ X = I$. (It's not a proper inverse, though, because $xx^+ \neq I$.)

1 As the Best Linear Unbiased Estimator

Recall that our model for the data is

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (2)$$

where $\mathbb{E}[\boldsymbol{\epsilon}|X] = \mathbf{0}$ and $\operatorname{Var}[\boldsymbol{\epsilon}|X] = \sigma^2 I$.

An estimator is any rule for estimating $\boldsymbol{\beta}$ based on the data \mathbf{y} . $\hat{\boldsymbol{\beta}}^{\text{OLS}} = X^+ \mathbf{y}$ is a natural estimator, of course, but there are others. For instance, we could ignore the data and always estimate $\boldsymbol{\beta}$ as $\mathbf{0}$. This is not a very good estimator, but an estimator nonetheless.

What makes $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ a better estimator than $\mathbf{0}$? A criterion that is often used to compare estimators is mean-squared error (MSE):

$$\operatorname{MSE}_{\hat{\boldsymbol{\beta}}}(\boldsymbol{\beta}) = \sum_{j=1}^p \mathbb{E}(\hat{\beta}_j - \beta_j)^2. \quad (3)$$

The MSE measures how far the estimator is from the truth, in expectation. If one estimator has a smaller MSE than another, for all values of $\boldsymbol{\beta}$, then the first estimator is clearly better.

An important fact about the MSE is that it can be decomposed into the sum of a bias and a variance term:

$$\begin{aligned} \operatorname{MSE}_{\hat{\boldsymbol{\beta}}}[\boldsymbol{\beta}] &= \sum_{j=1}^p \mathbb{E}(\mathbb{E}[\hat{\beta}_j] - \beta_j + \hat{\beta}_j - \mathbb{E}[\hat{\beta}_j])^2 \\ &= \sum_{j=1}^p \underbrace{(\mathbb{E}[\hat{\beta}_j] - \beta_j)^2}_{\text{bias}^2} + \underbrace{\mathbb{E}(\hat{\beta}_j - \mathbb{E}[\hat{\beta}_j])^2}_{\operatorname{Var}[\hat{\beta}_j]} + 2 \underbrace{\mathbb{E}[(\mathbb{E}[\hat{\beta}_j] - \beta_j)(\hat{\beta}_j - \mathbb{E}[\hat{\beta}_j])]}_0. \end{aligned}$$

We can decrease MSE by either decreasing the variance or the bias. In particular, it is sometimes desirable to trade a little bias for a massive reduction in variance.

If we only consider estimators that are unbiased (i.e., $E[\hat{\beta}] = \beta$), then the problem of minimizing the MSE becomes one of minimizing the variance. Therefore, the “best” estimator is the one which minimizes the variance.

Unfortunately, it’s not quite true that $\hat{\beta}^{OLS}$ minimizes the variance among all unbiased estimators. (It would be true if we assumed that the errors ϵ are also normally distributed.) However, it *is* true that $\hat{\beta}^{OLS}$ minimizes the variance among all *linear* unbiased estimators. (A *linear estimator* is an estimator of the form $\hat{\beta} = A\mathbf{y}$ for some matrix A .) For this reason, we say that $\hat{\beta}^{OLS}$ is the best linear unbiased estimator (BLUE).

Theorem 1. $\hat{\beta}^{OLS}$ is the best linear unbiased estimator. That is, if $A\mathbf{y}$ is any other linear unbiased estimator, then

$$\text{MSE}_{\beta}[A\mathbf{y}] = \sum_{j=1}^p \text{Var}_{\beta}[(A\mathbf{y})_j] > \sum_{j=1}^p \text{Var}_{\beta}[\hat{\beta}_j^{OLS}] = \text{MSE}_{\beta}[\hat{\beta}^{OLS}]. \quad (4)$$

Proof. First, let’s see what unbiasedness tells us about A . Since $E[A\mathbf{y}] = AX\beta = \beta$ for all β , this means that $AX = I$.

Now, let’s write the variance of $A\mathbf{y}$ in terms of the variance of $\hat{\beta}^{OLS}$:

$$\begin{aligned} \text{Var}[A\mathbf{y}] &= \text{Var}[A\mathbf{y} - \hat{\beta}^{OLS} + \hat{\beta}^{OLS}] \\ &= \text{Var}[\hat{\beta}^{OLS}] + \text{Var}[A\mathbf{y} - \hat{\beta}^{OLS}] + 2 \text{Cov}[A\mathbf{y} - \hat{\beta}^{OLS}, \hat{\beta}^{OLS}]. \end{aligned}$$

The final term is zero:

$$\begin{aligned} \text{Cov}[A\mathbf{y} - \hat{\beta}^{OLS}, \hat{\beta}^{OLS}] &= \text{Cov}[(A - X^+)\mathbf{y}, X^+\mathbf{y}] \\ &= (A - X^+) \text{Var}(\mathbf{y})(X^+)^T \\ &= \sigma^2(A - X^+)(X^+)^T \\ &= \sigma^2(A - (X^T X)^{-1} X^T)(X(X^T X)^{-1}) \\ &= \sigma^2 \underbrace{(AX - I)}_0 (X^T X)^{-1} \\ &= 0. \end{aligned}$$

Therefore, we have established that:

$$\text{Var}[A\mathbf{y}] = \text{Var}[\hat{\beta}^{OLS}] + \text{Var}[A\mathbf{y} - \hat{\beta}^{OLS}]. \quad (5)$$

The MSE is the sum of the diagonal entries of the variance matrix (i.e., the trace). By taking the trace of both sides of (5), we see that

$$\text{MSE}[A\mathbf{y}] = \text{tr}(\text{Var}[A\mathbf{y}]) = \text{tr}(\text{Var}[\hat{\beta}^{OLS}]) + \text{tr}(\text{Var}[A\mathbf{y} - \hat{\beta}^{OLS}]) = \text{MSE}[\hat{\beta}^{OLS}] + \text{something positive},$$

so $\hat{\beta}^{OLS}$ must have the smallest MSE among all estimators of the form $A\mathbf{y}$. \square

2 As the Maximum Likelihood Estimator

If we assume that the errors ϵ are normal, then under model (2), $\mathbf{y} \sim N(X\boldsymbol{\beta}, \sigma^2 I)$. One general strategy for estimating parameters is maximum likelihood estimation. It yields estimators that are asymptotically efficient (meaning that they achieve the smallest possible variance as $n \rightarrow \infty$).

Theorem 2. *The maximum likelihood estimator (MLE) of $\boldsymbol{\beta}$ when $\mathbf{y} \sim N(X\boldsymbol{\beta}, \sigma^2 I)$ is $\hat{\boldsymbol{\beta}}^{OLS}$.*

Proof. \mathbf{y} has a multivariate normal distribution with likelihood

$$L(\boldsymbol{\beta}) = \frac{1}{(2\pi)^{n/2}(\det(\sigma^2 I))^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{y} - X\boldsymbol{\beta})^T (\sigma^2 I)^{-1} (\mathbf{y} - X\boldsymbol{\beta}) \right\}.$$

We typically take logs before attempting to maximize the likelihood to make things easier.

$$\log L(\boldsymbol{\beta}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{y} - X\boldsymbol{\beta}\|^2. \quad (6)$$

Notice that the first term is just a constant (it doesn't depend on $\boldsymbol{\beta}$). So maximizing (6) is equivalent to maximizing $-\frac{1}{2\sigma^2} \|\mathbf{y} - X\boldsymbol{\beta}\|^2$. To find the value of $\boldsymbol{\beta}$ that maximizes this, we can alternatively minimize $\|\mathbf{y} - X\boldsymbol{\beta}\|^2$. But this is simply $\hat{\boldsymbol{\beta}}^{OLS}$. □

Remark 1. *In this example, there was a closed-form for the MLE because it reduced to a least-squares problem. However, most MLEs cannot be written down in closed form and have to be computed iteratively using optimization algorithms. All optimization algorithms work basically as follows: you start at a initial guess of the parameter and then iteratively move in a direction of increasing likelihood until you cannot increase the likelihood anymore. Two of the most commonly used methods (which we will use later in the course) are gradient descent and Newton's method.*

3 As a Sample Estimate of the MMSE Predictor

Suppose \mathbf{x}_0 and y_0 are random variables. We might want to predict y_0 using some function $\hat{f}(\mathbf{x}_0)$. A reasonable requirement is that this function minimize (mean square) prediction error, i.e.,

$$\operatorname{argmin}_f \mathbb{E}(y_0 - f(\mathbf{x}_0))^2.$$

It turns out that the minimum MSE (MMSE) predictor is the conditional expectation of y_0 given \mathbf{x}_0 .

Theorem 3. *The MMSE predictor is the conditional expectation $f(\mathbf{x}_0) = \mathbb{E}[y_0 | \mathbf{x}_0]$.*

Proof. Add and subtract $E[\mathbf{y}|\mathbf{x}]$, then expand:

$$\begin{aligned} E(y_0 - f(\mathbf{x}_0))^2 &= E(y_0 - E[y_0|\mathbf{x}_0] + E[y_0|\mathbf{x}_0] - f(\mathbf{x}_0))^2 \\ &= E(y_0 - E[y_0|\mathbf{x}_0])^2 + (E[y_0|\mathbf{x}_0] - f(\mathbf{x}_0))^2 + 2E[(y_0 - E[y_0|\mathbf{x}_0])(E[y_0|\mathbf{x}_0] - f(\mathbf{x}_0))] \end{aligned}$$

We can see that the last term is zero by conditioning further on \mathbf{x}_0 . To avoid notational overload, we will denote $E[y_0|\mathbf{x}_0]$ by $\mu(\mathbf{x}_0)$:

$$\begin{aligned} E[(y_0 - \mu(\mathbf{x}_0))(\mu(\mathbf{x}_0) - f(\mathbf{x}_0))] &= E\left[E[(y_0 - \mu(\mathbf{x}_0))(\mu(\mathbf{x}_0) - f(\mathbf{x}_0)) | \mathbf{x}_0]\right] \\ &= E\left[(\mu(\mathbf{x}_0) - f(\mathbf{x}_0)) E[y_0 - \mu(\mathbf{x}_0) | \mathbf{x}_0]\right] \\ &= E\left[(\mu(\mathbf{x}_0) - f(\mathbf{x}_0)) \underbrace{(E[y_0|\mathbf{x}_0] - \mu(\mathbf{x}_0))}_0\right] \\ &= 0. \end{aligned}$$

Thus, we have shown that

$$E(y_0 - f(\mathbf{x}_0))^2 = E(y_0 - E[y_0|\mathbf{x}_0])^2 + (E[y_0|\mathbf{x}_0] - f(\mathbf{x}_0))^2.$$

How should we choose f to minimize this? The first term on the right-hand side does not depend on f , so we have no control over it. But we can force the second term to zero by choosing $f(\mathbf{x}_0) = E[y_0|\mathbf{x}_0]$. Hence, this must be the choice that minimizes the MSE. \square

What does this conditional expectation look like? It depends on the distribution. But if (\mathbf{x}_0, y_0) are jointly normal, then $E[y_0|\mathbf{x}_0]$ is linear. In particular, suppose

$$\begin{pmatrix} \mathbf{x}_0 \\ y_0 \end{pmatrix} \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}\right). \quad (7)$$

Then, the conditional expectation and variance are given by

$$E[y_0|\mathbf{x}_0] = \Sigma_{yx}\Sigma_{xx}^{-1}\mathbf{x}_0 \quad (8)$$

$$\text{Var}[y_0|\mathbf{x}_0] = \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}. \quad (9)$$

Unfortunately, we usually cannot compute Σ_{xx} , Σ_{xy} , or Σ_{yy} in practice, as these are unknown population quantities. But they can be estimated from data (\mathbf{x}_i, y_i) :

$$\hat{\Sigma}_{xx} = \frac{1}{n}X^T X \quad (10)$$

$$\hat{\Sigma}_{xy} = \frac{1}{n}X^T \mathbf{y} \quad (11)$$

$$\hat{\Sigma}_{yy} = \frac{1}{n}\mathbf{y}^T \mathbf{y}. \quad (12)$$

Substituting these into (8), we estimate the MMSE predictor function as

$$\hat{f}(\mathbf{x}_0) = \frac{1}{n}\mathbf{y}^T X \left(\frac{1}{n}X^T X\right)^{-1} \mathbf{x}_0.$$

If you don't recognize this yet, try canceling the $\frac{1}{n}$'s and transposing it, which should be the same (since this is a 1×1 matrix). You should obtain:

$$\hat{f}(\mathbf{x}_0) = \mathbf{x}_0^T (X^T X)^{-1} X^T \mathbf{y} = \mathbf{x}_0^T \hat{\boldsymbol{\beta}}^{\text{OLS}}.$$

This is of course how you would predict y_0 once you had the OLS coefficients.