

Linear Models Lecture 3: Estimation of the Linear Projection Model

Robert Gulotty

University of Chicago

February 20, 2026

Deriving the OLS estimator

- Recall, geometry can work in populations or samples.
- In the following slides we will derive
 - the OLS slope estimator $\hat{\beta}$,
 - the OLS estimator for error variance,
 - Decomposition of OLS (Analysis of Variance), with R^2
- Next time, the FWL Theorem.

Review: Modeling the Stochastic Element in the Population

- The CEF is the regression of Y on X , but requires knowing the distribution of (Y, X) , a fact about Populations.
- The linear projection model predicts Y as a linear function of $X = (X_1, X_2, \dots, 1)$, approximating the CEF, but still a fact about Populations.

$$\beta = [\mathbb{E}(XX')]^{-1}\mathbb{E}[XY]$$

- Estimation involves using a model of a sample to make inferences about the population.
- A statistic is random variable constructed from the sample, can be used for description or inference.
- An estimator is a statistic used to estimate population parameters.

Samples

- We estimate the projection coefficient β with measurements of (Y, X) , a *sample*.
- $\{(Y_i, X_i) : i \in 1, \dots, n\}$, is the sample, governed by the distribution of (Y, X) .
- From the next slide forward, X is going to be a matrix, and we will use lower case for vectors e.g. e .

Identically Distributed.

- Main approach is to assume homogeneity: $\{(Y_i, \mathbf{x}_i) : i \in 1, \dots, n\}$ are identically distributed.

$$\begin{aligned} \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} &= \beta_1 \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} + \beta_2 \begin{bmatrix} X_{21} \\ \vdots \\ X_{2n} \end{bmatrix} + \beta_3 \begin{bmatrix} X_{31} \\ \vdots \\ X_{3n} \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix} \\ &= \begin{bmatrix} 1 & X_{21} & X_{31} \\ \vdots & \vdots & \vdots \\ 1 & X_{2n} & X_{3n} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix} \\ \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \end{aligned}$$

- Note, $E(Y_1) \neq E(Y_2)$.

Least Squares Estimator

- The linear projection coefficient β is the minimizer of the expected squared error

$$S(\beta) = \mathbb{E}[(Y - X'\beta)^2]$$

- The *moment estimator* of $S(\beta)$ is the sample average:

$$\hat{S}(\beta) = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i'\beta)^2$$

- $\sum_{i=1}^n (Y_i - \mathbf{x}_i'\beta)^2$ is called $\text{SSE}(\beta)$, the **sum of squared errors**.
- The estimator $\hat{\beta}$ is the minimizer of $\hat{S}(\beta)$, as well as the minimizer of the SSE.
- $\hat{\beta}$ is called the least squares estimator, sometimes written $\hat{\beta}_n$.

Solving the minimization problem

- The form of $\hat{\beta}$ is the sample analogue to the population linear projection coefficient β .
- The solution will require some additional concepts associated with matrix/multivariate calculus.
 - Quadratic form of a matrix, connecting matrices with polynomial equations.
 - Definiteness of matrices, connecting matrices to convex functions.
 - Generalizing derivatives to cover matrices.

Finite Sample Assumptions of OLS estimator

- Assumption 1: The random variables $\{(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)\}$ are independent and identically distributed.
- Assumption 2: $Y = \mathbf{X}'\beta + e$, where $\mathbb{E}(e|\mathbf{X}) = 0$.
- Assumption 3: $\mathbb{E}(\mathbf{X}\mathbf{X}') > 0$ is invertible (with probability 1).
- Assumption 4*: $\mathbb{E}[e^2|\mathbf{X}] = \sigma^2(\mathbf{X}) = \sigma^2$.

Derivation of Least Squares Estimator

- The functional form of the least squares estimator can be derived via multivariate calculus.
- We look for a minimum of $SSE(\beta)$.
- Take first derivatives to calculate first-order conditions and locate the critical values.
- Then we check we have the global minimum, by using the fact that positive definite matrices represent convex functions.

Calculus Derivation of Least Squares Estimator

$$\begin{aligned} SSE(\beta) &= (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) && \text{(Definition of } \mathbf{e} \text{)} \\ &= (\mathbf{y}' - \beta'\mathbf{X}')(\mathbf{y} - \mathbf{X}\beta) \\ &= \mathbf{y}'\mathbf{y} - \beta'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\beta + \beta'\mathbf{X}'\mathbf{X}\beta && \text{(Transpose rules)} \end{aligned}$$

(Distributive property)

$$\min_{\beta} SSE(\beta) = \min_{\beta} \underbrace{\mathbf{y}'\mathbf{y}}_{\text{constant}} - \underbrace{2\mathbf{y}'\mathbf{X}\beta}_{\text{linear in } \beta} + \underbrace{\beta'\mathbf{X}'\mathbf{X}\beta}_{\text{quadratic form}}$$

Calculus Derivation of Least Squares Estimator

$$\begin{aligned}\frac{\partial \mathbf{y}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{X}'\mathbf{Y} - \mathbf{y}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{\partial \boldsymbol{\beta}} &= \mathbf{0} - \mathbf{X}'\mathbf{Y} - \mathbf{Y}'\mathbf{X} + (\mathbf{X}'\mathbf{X} + \mathbf{X}'\mathbf{X})\hat{\boldsymbol{\beta}} \\ &= -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \quad (\text{Symmetric}) \\ -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{0} \quad (\text{find critical point}) \\ \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{X}'\mathbf{y}. \quad (\text{normal equations}) \\ \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (\text{If } \mathbf{X}'\mathbf{X} \text{ can be inverted (A3).})\end{aligned}$$

Second order condition

Take second derivative, and check that it is positive

$$\frac{\partial -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\beta}}{\partial \hat{\beta}'} = 2(\mathbf{X}'\mathbf{X})'$$

Which is positive so long as $\mathbf{X}'\mathbf{X}$ is positive definite (same assumption as A1).
 $\hat{\beta}$ minimizes the sum of squared errors, assuming \mathbf{X} is of full rank.

OLS in R: The Formula in Action

```
library(carData)
X <- cbind(1, Prestige$education, Prestige$income)
y <- Prestige$prestige

# The matrix formula from the slides:
solve(t(X) %*% X) %*% t(X) %*% y

# What lm() computes:
coef(lm(prestige ~ education + income, data=Prestige))
```

The two are identical: `lm()` implements $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.

Projection Matrix

- There is a geometric interpretation for the minimization problem that uses the idea of "projection".
- The following matrix is called a "projection matrix"

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

- Properties of \mathbf{P}
 - $\mathbf{P}\mathbf{X} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} = \mathbf{X}$
 - If $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$, then $\mathbf{P}\mathbf{X}_1 = \mathbf{X}_1$.
 - $\mathbf{P} = \mathbf{P}'$, that is, \mathbf{P} is symmetric.
 - $\mathbf{P}\mathbf{P} = \mathbf{P}$, that is, \mathbf{P} is idempotent.
 - $\mathbf{P}\mathbf{y} = \mathbf{X}\hat{\beta} = \hat{\mathbf{y}}$, called the fitted value. This is why \mathbf{P} is called the "hat matrix".

Projection onto column of 1

- If $X = \mathbf{1}$, then $P_1 = \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}' = \frac{1}{n}\mathbf{1}\mathbf{1}'$
- So if we project y onto P_1 , we get a vector repeating the sample mean:

$$P_1y = \frac{1}{n}\mathbf{1}\mathbf{1}'y = \mathbf{1} \cdot \frac{1}{n} \sum_{i=1}^n Y_i = \mathbf{1}\bar{Y}$$

- **Takeaway:** The simplest regression (intercept only) projects y onto the sample mean. Every additional regressor refines this projection.
- The annihilator $M_1y = y - \mathbf{1}\bar{Y}$ *demeans* y —it removes the part explained by the intercept.

Orthogonal Projection

- The following matrix is called an "orthogonal projection matrix", or the Annihilator Matrix.

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

- \mathbf{M} and \mathbf{X} are orthogonal:

$$\begin{aligned} \mathbf{MX} &= (\mathbf{I}_n - \mathbf{P})\mathbf{X} \\ &= \mathbf{X} - \mathbf{PX} \\ &= \mathbf{X} - \mathbf{X} \\ &= 0 \end{aligned}$$

- We can define the residuals from projecting onto a column of 1s: $\mathbf{M}_1\mathbf{y} = \mathbf{y} - \mathbf{1}\bar{Y}$, this demeans \mathbf{y} .

Relationship between residuals and disturbances

$$\begin{aligned}\hat{\mathbf{e}} &= \mathbf{y} - \hat{\mathbf{y}} \\ &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{y} - \mathbf{P}\mathbf{y} \\ &= \mathbf{M}\mathbf{y} \\ &= \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \mathbf{e}) \\ &= \mathbf{M}\mathbf{X}\boldsymbol{\beta} + \mathbf{M}\mathbf{e} \\ &= \mathbf{M}\mathbf{e}\end{aligned}$$

Projection and Residuals in R

```
mod <- lm(prestige ~ education + income, data=Prestige)
X <- cbind(1, Prestige$education, Prestige$income)
P <- X %*% solve(t(X) %*% X) %*% t(X)
M <- diag(nrow(X)) - P

# P*y = fitted values
all.equal(as.vector(P %*% y), fitted(mod)) # TRUE

# M*y = residuals
all.equal(as.vector(M %*% y), resid(mod)) # TRUE
```

P and **M** are not abstract notation—they are the matrices R uses to compute `fitted()` and `resid()`.

Why Does Projection Help?

- We want to estimate the systematic part $\mathbf{X}\beta$. Two candidates:
 - 1 Use \mathbf{y} itself (just report the raw data).
 - 2 Use $\mathbf{P}\mathbf{y} = \hat{\mathbf{y}}$ (the projection onto the column space of \mathbf{X}).
- Compare their variances as estimators of $\mathbf{X}\beta$:

$$\text{Var}(\mathbf{y}|\mathbf{X}) = \sigma^2 \mathbf{I}$$

$$\text{Var}(\mathbf{P}\mathbf{y}|\mathbf{X}) = \sigma^2 \mathbf{P}$$

- The difference is $\sigma^2(\mathbf{I} - \mathbf{P}) = \sigma^2 \mathbf{M}$, which is positive semi-definite.
- So $\text{Var}(\mathbf{y}) \geq \text{Var}(\mathbf{P}\mathbf{y})$: projection *removes noise* without distorting the signal, because $\mathbf{P}\mathbf{X}\beta = \mathbf{X}\beta$.

Variance of the regression errors

- We want to measure the precision of our regression estimates.
- The error (or disturbance) of an observation is the deviation of a value from its theoretical mean, cannot be observed.
- The coefficients inherit that uncertainty from the theoretical model.
- Our estimates have additional uncertainty, related to the fact we only have a sample.
- Residuals are the difference between observations and estimates, can be observed.

The bane of statistics

- The mean squared error, or MSE, is calculated on the computed residuals, not the unobservable errors.
- Residuals, although observed, have a distribution that is not identical to the population disturbances.
- We will try to stick to disturbances and residuals to make distinctions.

Estimation of Error Variance

- Natural estimator: plug residuals \hat{e}_i into the sample variance formula:

$$\hat{\sigma}^2 = n^{-1} \hat{\mathbf{e}}' \hat{\mathbf{e}} = n^{-1} (\mathbf{M}\mathbf{y})' (\mathbf{M}\mathbf{y}) = n^{-1} \mathbf{y}' \mathbf{M} \mathbf{y} = n^{-1} \mathbf{e}' \mathbf{M} \mathbf{e}$$

where the last step uses $\mathbf{M} = \mathbf{M}'\mathbf{M}$ and $\mathbf{M}\mathbf{X} = 0$.

- But $\hat{\sigma}^2$ underestimates the true variance $\tilde{\sigma}^2 = n^{-1} \mathbf{e}' \mathbf{e}$:

$$\tilde{\sigma}^2 - \hat{\sigma}^2 = n^{-1} \mathbf{e}' \mathbf{e} - n^{-1} \mathbf{e}' \mathbf{M} \mathbf{e} = n^{-1} \mathbf{e}' \mathbf{P} \mathbf{e} \geq 0$$

because \mathbf{P} is positive semi-definite.

- Projection “absorbs” some of the error into fitted values, so residuals understate the true noise.

Estimation of σ_e^2

- Our estimator will be $s_{\hat{e}}^2 = \hat{\mathbf{e}}' \hat{\mathbf{e}} / (N - K)$.
- Proof of unbiasedness of $s_{\hat{e}}^2$:

$$\hat{\mathbf{e}}' \hat{\mathbf{e}} = \mathbf{e}' \mathbf{M}' \mathbf{M} \mathbf{e} = \mathbf{e}' \mathbf{M} \mathbf{e} \quad (\mathbf{M} \text{ symmetric, idempotent})$$

$$\begin{aligned} E[\hat{\mathbf{e}}' \hat{\mathbf{e}} | \mathbf{X}] &= E[\mathbf{e}' \mathbf{M} \mathbf{e} | \mathbf{X}] \\ &= E[\text{tr}(\mathbf{e}' \mathbf{M} \mathbf{e}) | \mathbf{X}] \quad (\text{scalar} = \text{its own trace}) \\ &= E[\text{tr}(\mathbf{M} \mathbf{e} \mathbf{e}' | \mathbf{X})] \\ &= \text{tr}(\mathbf{M} E[\mathbf{e} \mathbf{e}' | \mathbf{X}]) \quad (\mathbf{M} \text{ fixed given } \mathbf{X}) \\ &= \text{tr}(\mathbf{M} \sigma^2 \mathbf{I}) = \sigma^2 \text{tr}(\mathbf{M}) = \sigma_e^2 (N - K) \quad (\text{A4, tr}(\mathbf{M}) = N - K) \end{aligned}$$

Standard Deviation and Standard Error

- Random variables have a mean and a standard deviation, the latter quantifies variability.
- Statistics (like the sample mean) are calculated from a random sample to make inferences.
- Statistics have a sampling distribution with a mean and standard deviation.
- The standard deviation of an estimated statistic is called the "standard error".
- $\sqrt{s_e^2}$ is called the (computed) standard error of the regression
- $\sqrt{s_e^2}$ is a biased estimator for σ , but it is **consistent**.
- Similarly, the standard deviation of $\hat{\beta}$ is estimated using the standard errors.

That is, we have a tool for estimating parameters which itself has parameters to estimate!

Interpretation of Residual Standard Error

- summary(lm()) reports s_e as the Residual standard error or "sigma"
- STATA reports this quantity as "Root Mean Squared Error".
- s_e is on the same scale as the outcome.
- If $y > 0$, it can make sense to calculate an *average prediction error rate*: $\frac{s_e}{\bar{y}}$.
- Example: if we are predicting houses that are on average 1 million dollars, being generally off by \$50k is less bad than if we are predicting houses that are \$100k.

Variance of $\hat{\beta}$

$$\begin{aligned} V(\hat{\beta} | \mathbf{X}) &= V((\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} | \mathbf{X}) && \text{(Normal Equation)} \\ &= V(\mathbf{A}\mathbf{y} | \mathbf{X}) && \text{(Pick } \mathbf{A} \text{ as placeholder)} \\ &= \mathbf{A} V(\mathbf{y} | \mathbf{X}) \mathbf{A}' && \text{(Var}(cX) = c^2 \text{Var}(X)) \\ &= \mathbf{A} \Sigma \mathbf{A}' \\ &= \mathbf{A} (\sigma^2 \mathbf{I}) \mathbf{A}' && \text{(Assumption A4)} \\ &= \sigma^2 \mathbf{A} \mathbf{A}' \\ &= \sigma^2 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}' \mathbf{X})^{-1} \\ \hat{V}(\hat{\beta} | \mathbf{X}) &= s_e^2 (\mathbf{X}' \mathbf{X})^{-1} && \text{(Standard Error of } \hat{\beta}) \end{aligned}$$

Standard Errors in R: Connecting to the Formula

```
mod <- lm(prestige ~ education + income, data=Prestige)

# The formula from the slides: s^2 (X'X)^{-1}
sigma(mod)^2 * solve(t(X) %*% X)

# What R reports:
vcov(mod)

# Standard errors = sqrt of the diagonal
sqrt(diag(vcov(mod)))
coef(summary(mod))[, "Std. Error"] # identical
```

Every standard error in `summary(lm())` comes from $\sqrt{\text{diag}(s_e^2(\mathbf{X}'\mathbf{X})^{-1})}$.

Decomposition

- M makes the least squares residuals:

$$My = y - Py = y - X\hat{\beta} = \hat{e}$$

- We can rewrite y in terms of

$$y = Py + My = \hat{y} + \hat{e}$$

- This decomposition is orthogonal,

$$\hat{y}'\hat{e} = (PY)'(MY) = y'PMY = 0$$

$$\begin{aligned} Y'Y &= (\hat{y} + \hat{e})'(\hat{y} + \hat{e}) = (\hat{Y}'\hat{Y} + \hat{Y}'\hat{e} + \hat{e}'\hat{Y} + \hat{e}'\hat{e}) \\ &= \hat{Y}'\hat{Y} + \hat{e}'\hat{e} \end{aligned}$$

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \hat{e}_i^2$$

Analysis of Variance

- Subtract \bar{Y} from both sides of the decomposition:

$$\mathbf{y} - \mathbf{1}\bar{Y} = (\hat{\mathbf{y}} - \mathbf{1}\bar{Y}) + \hat{\mathbf{e}}$$

- Take inner products. The cross terms vanish when \mathbf{X} contains a constant:

$$(\hat{\mathbf{y}} - \mathbf{1}\bar{Y})' \hat{\mathbf{e}} = \hat{\mathbf{y}}' \hat{\mathbf{e}} - \bar{Y} \mathbf{1}' \hat{\mathbf{e}} = 0 - 0 = 0$$

- So we get the **analysis of variance** (ANOVA) decomposition:

$$\begin{aligned}\|\mathbf{y} - \mathbf{1}\bar{Y}\|^2 &= \|\hat{\mathbf{y}} - \mathbf{1}\bar{Y}\|^2 + \|\hat{\mathbf{e}}\|^2 \\ \underbrace{\sum(Y_i - \bar{Y})^2}_{SST} &= \underbrace{\sum(\hat{Y}_i - \bar{Y})^2}_{SSR \text{ (regression)}} + \underbrace{\sum \hat{e}_i^2}_{SSE \text{ (error)}}\end{aligned}$$

Goodness-of-Fit: R^2

- Dividing through by SST:

$$1 = \frac{\text{SSR}}{\text{SST}} + \frac{\text{SSE}}{\text{SST}}$$
$$R^2 \equiv \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}$$

- R^2 is the proportion of the variance in Y that is linearly explained by \mathbf{X} .
- $0 \leq R^2 \leq 1$: equals 0 when \mathbf{X} explains nothing; equals 1 when $\hat{e}_i = 0$ for all i .
- We can always make $R^2 = 1$ by adding enough linearly independent columns to \mathbf{X} (one per observation). This does not mean the model is good.

R^2 in Practice

- R^2 measures *descriptive fit*, not causal validity. A high R^2 does not mean the coefficients are unbiased; a low R^2 does not mean the model is useless.
- Typical values vary by field:
 - Cross-sectional micro data (e.g. earnings regressions): $R^2 \approx 0.2\text{--}0.4$
 - Aggregate time-series (e.g. GDP growth models): $R^2 \approx 0.7\text{--}0.9$
 - Experimental data with noisy outcomes: R^2 can be very low and the treatment effect still precisely estimated.
- In R: `summary(lm(y ~ x))$r.squared` reports R^2 .
- The **adjusted** $R^2 = 1 - \frac{n-1}{n-k}(1 - R^2)$ penalizes adding regressors and can decrease when a variable adds no explanatory power.

Naming Conventions: A Warning

- Different textbooks use SSE and SSR with *opposite* meanings:

	Hansen / this course	Some other texts
$\sum(\hat{Y}_i - \bar{Y})^2$	SSR (regression)	ESS or SSR
$\sum \hat{e}_i^2$	SSE (error)	RSS or SSE

- The underlying math is always the same: $SST = (\text{explained}) + (\text{unexplained})$.
- When reading papers or other textbooks, check which convention is in use.