12: Ordinary Least Square

Naijia Liu

Spring 2025



Partitioned regression

• Partition covariates and coefficients $\mathbf{X}=[\mathbf{X}_1\ \mathbf{X}_2]$ and $\boldsymbol{\beta}=(\beta_1,\beta_2)'$:

$$\mathbf{Y} = \mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2 + \mathbf{e}$$

- Can we find expressions for $\hat{\beta}_1$ and $\hat{\beta}_2$?
- **Residual regression** or Frisch-Waugh-Lovell theorem to obtain $\hat{\beta}_1$:
 - Use OLS to regress ${\bf Y}$ on ${\bf X}_2$ and obtain residuals $\tilde{{\bf e}}_2$
 - \blacktriangleright Use OLS to regress each column of X_1 on X_2 and obtain residuals \tilde{X}_1

• Use OLS to regress
$$\tilde{\mathbf{e}}_2$$
 on $\tilde{\mathbf{X}}_1$

Focus on simple case

- Focus on single covariate model with no intercept: $Y_i = X_i\beta + e_i$
- Let $\mathbf{X} = (X_1, \dots, X_n)$ and recall inner product:

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \sum_{i=1}^{n} X_i Y_i$$

Inner products measure how similar two vectors are.

• Slope in this case:

$$\hat{\beta} = \frac{\sum_{i=1}^{n} X_i Y_i}{\sum_{i=1}^{n} X_i^2} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\langle \mathbf{X}, \mathbf{X} \rangle}$$

• Suppose we add an orthogonal covariate $\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \mathbf{e}$ with $\langle \mathbf{X}, \mathbf{Z} \rangle = 0$:

$$\hat{\beta} = \frac{\langle \mathbf{X}, \mathbf{Y} \rangle}{\langle \mathbf{X}, \mathbf{X} \rangle}, \quad \hat{\gamma} = \frac{\langle \mathbf{Z}, \mathbf{Y} \rangle}{\langle \mathbf{Z}, \mathbf{Z} \rangle}$$

- With exactly orthogonal covariates, multivariate OLS is the same as univariate OLS
 - Only holds in balanced, designed experiments.

Gov 2001

Partitioned regression and partial regression

Adding the intercept

• Consider the OLS slope with an intercept:

$$\hat{\beta} = \frac{\sum_{i=1}^{n} (X_i - \bar{X}) (Y_i - \bar{Y})}{\sum_{i=1}^{n} (X_i - \bar{X})^2} = \frac{\langle \mathbf{X} - \bar{X}\mathbf{1}, \ \mathbf{Y} - \bar{Y}\mathbf{1} \rangle}{\langle \mathbf{X} - \bar{X}\mathbf{1}, \ \mathbf{X} - \bar{X}\mathbf{1} \rangle} = \frac{\langle \mathbf{X} - \bar{X}\mathbf{1}, \ \mathbf{Y} \rangle}{\langle \mathbf{X} - \bar{X}\mathbf{1}, \ \mathbf{X} - \bar{X}\mathbf{1} \rangle}$$

• How can we get this?

- 1. Regress **X** on 1 to get coefficient \overline{X}
- 2. Regress **Y** on residuals from step 1, $\mathbf{X} \overline{X}\mathbf{1}$
- If we wanted to get the coefficient on added variable Z_i , we could repeat this:
 - 1. Regress Z on $\tilde{X} = X \bar{X}1$ and obtain coefficient $\langle Z, \tilde{X} \rangle / \langle \tilde{X}, \tilde{X} \rangle$
 - 2. Regress Y on residual from ...

Why does residual regression work?

• We can find $\hat{\beta}_1$ by nested minimization:

$$\hat{eta}_1 = rg\min_{eta_1} \left(\min_{eta_2} \| \mathbf{Y} - \mathbf{X}_1 eta_1 - \mathbf{X}_2 eta_2 \|^2
ight)$$

First find the minimum of the SSR over β_2 , fixing β_1

- Then find β_1 that minimizes the resulting SSR
- The projection and annihilator matrices are defined only by covariates.

$$\mathbf{M}_2 = I_n - \mathbf{X}_2 (\mathbf{X}_2' \mathbf{X}_2)^{-1} \mathbf{X}_2'$$

• Creates residuals from a regression on or \mathbf{X}_2

• Solving the nested minimization gives:

$$\hat{\beta}_1 = (\mathbf{X}_1'\mathbf{M}_2\mathbf{X}_1)^{-1}(\mathbf{X}_1'\mathbf{M}_2\mathbf{Y})$$

When will β₁ be the same regardless of whether X₂ is included?
If X₁ and X₂ are orthogonal, so X₂X₁ = 0 and M₂X₁ = X₁

Gov 2001

Residual regression

• Define two sets of residuals:

$$\blacktriangleright ~~ {\tilde {\bf X}}_2 = {\bf M}_1 {\bf X}_2 = {\sf residuals} {\sf ~from} {\sf ~regression} {\sf ~of} {\sf ~{\bf X}}_2 {\sf ~on} {\sf ~{\bf X}}_1$$

 $\blacktriangleright ~~ \tilde{\mathbf{e}}_1 = \mathbf{M}_1 \mathbf{Y} = \text{residuals from regression of } \mathbf{Y} \text{ on } \mathbf{X}_1$

 $\bullet\,$ Then remembering that \mathbf{M}_1 is symmetric and idempotent:

$$\begin{split} \hat{\beta}_2 &= (\mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}_2' \mathbf{M}_1 \mathbf{Y}) \\ &= (\mathbf{X}_2' \mathbf{M}_1 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}_2' \mathbf{M}_1 \mathbf{M}_1 \mathbf{Y}) \\ &= (\tilde{\mathbf{X}}_2' \tilde{\mathbf{X}}_2)^{-1} (\tilde{\mathbf{X}}_2' \tilde{\mathbf{e}}_1) \end{split}$$

- $\hat{\beta}_2$ can be obtained from a regression of $\tilde{\mathbf{e}}_1$ on $\tilde{\mathbf{X}}_2$.
 - \blacktriangleright Same result applies when using ${\bf Y}$ in place of \tilde{e}_1
 - Intuition: residuals are orthogonal
 - Called the Frisch-Waugh-Lovell Theorem
 - Sample version of the results we saw for the linear projection

Outliers, leverage points, and influential observations

- Least square heavily penalizes large residuals.
- Implies a just a few unusual observations can be extremely influential.
 - Dropping them leads to large changes in the estimated $\hat{\beta}$.
 - ▶ Not all "unusual" observations have the same effect, though.
- Useful to categorize:
 - 1. Leverage point: extreme in one X direction
 - 2. **Outlier**: extreme in the *Y* direction
 - 3. Influence point: extreme in both directions

Butterfly did it?

Wand et al, 2001 APSR



Butterfly did it?



Leverage point definition



- Values that are extreme in the X dimension
- That is, values far from the center of the covariate distribution

Leverage values

$$\hat{\mathbf{Y}} = \mathbf{P}\mathbf{Y} \implies \hat{Y}_i = \sum_{j=1}^n h_{ij}Y_j$$

 $\blacktriangleright \ h_{ij} =$ importance of observation j is for the fitted value \hat{Y}_i

- Leverage/hat values: h_{ii} diagonal entries of the hat matrix
- With a simple linear regression, we have

$$h_{ii} = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{j=1}^n (X_j - \bar{X})^2}$$

• \rightsquigarrow how far *i* is from the center of the *X* distribution

• Rule of thumb: examine hat values greater than 2(k+1)/n

Outlier definition



• An **outlier** is far away from the center of the Y distribution.

• Intuitively: a point that would be poorly predicted by the regression.

Gov 2001

Detecting outliers

- Want values poorly predicted? Look for big residuals, right?
 - Problem: we use i to estimate $\hat{\beta}$ so \hat{Y} aren't valid predictions.
 - \blacktriangleright unit might pull the regression line toward itself \rightsquigarrow small residual
- Better: leave-one-out prediction errors,

1. Regress $\mathbf{Y}_{(-\mathit{i})}$ on $\mathbf{X}_{(-\mathit{i})}\text{,}$ where these omit unit $\mathit{i}\text{:}$

$$\hat{\beta}_{(-i)} = \left(\mathbf{X}_{(-i)}'\mathbf{X}_{(-i)}\right)^{-1}\mathbf{X}_{(-i)}'\mathbf{Y}_{(-i)}$$

- 2. Calculate predicted value of Y_i using that regression: $\tilde{Y}_i = \mathbf{X}'_i \hat{\beta}_{(-i)}$
- 3. Calculate prediction error: $\tilde{e}_i = Y_i \tilde{Y}_i$
- Simple closed-form expressions:

$$\hat{\beta}_{(-i)} = \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}_i\hat{e}_i, \qquad \tilde{e}_i = \frac{\hat{e}_i}{1 - h_{ii}}$$

Influence points



- An **influence point** is one that is both an outlier and a leverage point.
- Extreme in both the X and Y dimensions

Overall measures of influence

• Influence of *i* can be measured by change in predictions:

$$\hat{Y}_i - \tilde{Y}_i = h_{ii}\tilde{e}_i$$

- How much does excluding *i* from the regression change its predicted value?
- ► Equal to "leverage × outlier-ness"
- Lots of diagnostics exist, but are mostly heuristic.
 - Does removing the point change a coefficient by a lot?

Limitations of the standard tools



- What happens when there are two influence points?
- Red line drops the red influence point
- Blue line drops the blue influence point

What to do about outliers and influential units?

- Is the data corrupted?
 - Fix the observation (obvious data entry errors)
 - Remove the observation
 - Be transparent either way
- Is the outlier part of the data generating process?
 - ▶ Transform the dependent variable (log(y))
 - Use a method that is robust to outliers (robust regression, least absolute deviations)