

Underfitting and Overfitting

Underfitting

Suppose the true model is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta} + \boldsymbol{\varepsilon},$$

where $\boldsymbol{\eta}$ is an unknown fixed vector and $\boldsymbol{\varepsilon}$ satisfies the GMM.

Suppose we incorrectly assume that

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

This example of misspecifying the model is known as underfitting.

Note that η may equal $W\alpha$ for some design matrix W whose columns could contain explanatory variables excluded from X .

What are the implications of underfitting?

Find $E(\mathbf{c}'\hat{\boldsymbol{\beta}})$.

Find $E(\hat{\sigma}^2)$.

Derivation of $E(\mathbf{c}'\hat{\boldsymbol{\beta}})$

$$\begin{aligned}E(\mathbf{c}'\hat{\boldsymbol{\beta}}) &= E(\mathbf{a}'\mathbf{X}\hat{\boldsymbol{\beta}}) \\&= E(\mathbf{a}'\mathbf{P}_X\mathbf{y}) \\&= \mathbf{a}'\mathbf{P}_X E(\mathbf{y}) \\&= \mathbf{a}'\mathbf{P}_X(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}) \\&= \mathbf{a}'\mathbf{X}\boldsymbol{\beta} + \mathbf{a}'\mathbf{P}_X\boldsymbol{\eta} \\&= \mathbf{c}'\boldsymbol{\beta} + \mathbf{a}'\mathbf{P}_X\boldsymbol{\eta}.\end{aligned}$$

$\mathbf{c}'\hat{\boldsymbol{\beta}}$ is biased for $\mathbf{c}'\boldsymbol{\beta}$ unless $\mathbf{a}'\mathbf{P}_X\boldsymbol{\eta} = 0$.

Note that if η is close to $\mathbf{0}$, the bias $\mathbf{a}'\mathbf{P}_X\eta$ may be small.

If $\eta \in \mathcal{C}(\mathbf{X})^\perp = \mathcal{N}(\mathbf{X}')$, then

$$\begin{aligned}\mathbf{X}'\eta = \mathbf{0} &\Rightarrow \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\eta = \mathbf{0} \\ &\Rightarrow \mathbf{a}'\mathbf{P}_X\eta = 0.\end{aligned}$$

As an example of this last point, suppose we fit a multiple regression but omit one explanatory variable.

Suppose for our sample of n observations, the vector $\mathbf{x}^* = \mathbf{x} - \bar{x}\mathbf{1}$ contains the values of the missing variable centered so that the sample mean is zero.

If the sample covariance of the missing variable x with each of the included variables x_1, \dots, x_p is 0, then the LSE of the multiple regression coefficients $\hat{\beta}$ will still be unbiased for β even though x is excluded $\because X'x^* = \mathbf{0}$.

Derivation of $E(\hat{\sigma}^2)$

$$\begin{aligned}(n-r)E(\hat{\sigma}^2) &= E(\mathbf{y}'(\mathbf{I} - \mathbf{P}_X)\mathbf{y}) \\ &= (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta})'(\mathbf{I} - \mathbf{P}_X)(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}) + \text{tr}((\mathbf{I} - \mathbf{P}_X)\sigma^2\mathbf{I}) \\ &= \boldsymbol{\eta}'(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta} + \sigma^2(n-r) \\ &\because \mathbf{X}'(\mathbf{I} - \mathbf{P}_X) = \mathbf{0} \quad \text{and} \quad (\mathbf{I} - \mathbf{P}_X)\mathbf{X} = \mathbf{0}. \\ \therefore E(\hat{\sigma}^2) &= \frac{\boldsymbol{\eta}'(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}}{n-r} + \sigma^2.\end{aligned}$$

Note that

$$\begin{aligned}\boldsymbol{\eta}'(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta} &= \boldsymbol{\eta}'(\mathbf{I} - \mathbf{P}_X)'(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta} \\ &= \|(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}\|^2 \\ &= \|\boldsymbol{\eta} - \mathbf{P}_X\boldsymbol{\eta}\|^2.\end{aligned}$$

Thus, $E(\hat{\sigma}^2) = \sigma^2$ iff

$$\begin{aligned}(I - P_X)\eta = \mathbf{0} &\iff \eta \in \mathcal{N}(I - P_X) \\ &\iff \eta \in \mathcal{C}(P_X) = \mathcal{C}(X) \\ &\iff \exists \alpha \ni X\alpha = \eta \\ &\iff \exists \alpha \ni E(\mathbf{y}) = X\beta + \eta \\ &\qquad\qquad\qquad = X\beta + X\alpha \\ &\qquad\qquad\qquad = X(\beta + \alpha) \\ &\iff E(\mathbf{y}) \in \mathcal{C}(X).\end{aligned}$$

Example 1

Consider an experiment with two experimental units (mice in this case) for each of two treatments.

We might assume the GMM holds with

$$E(\mathbf{y}) = E \left(\begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \end{bmatrix} \right) = \mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix} .$$

Example 1 (continued)

Suppose the person who conducted the experiment neglected to mention that, in each treatment group, one of the experimental units was male and the other was female.

Example 1 (continued)

Then the true model may require

$$\begin{aligned} E(\mathbf{y}) &= \begin{bmatrix} \mu + \tau_1 + \alpha/2 \\ \mu + \tau_1 - \alpha/2 \\ \mu + \tau_2 + \alpha/2 \\ \mu + \tau_2 - \alpha/2 \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix} + \begin{bmatrix} \alpha/2 \\ -\alpha/2 \\ \alpha/2 \\ -\alpha/2 \end{bmatrix} \\ &= \mathbf{X}\boldsymbol{\beta} + \begin{bmatrix} 1/2 \\ -1/2 \\ 1/2 \\ -1/2 \end{bmatrix} \begin{bmatrix} \alpha \end{bmatrix} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}\boldsymbol{\alpha} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}. \end{aligned}$$

Example 1 (continued)

If we analyze the data assuming the GMM with $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$, determine

- 1 $E(\widehat{\tau_1 - \tau_2})$, and
- 2 $E(\hat{\sigma}^2)$.

Example 1 (continued)

From slide 6,

$$\begin{aligned} E(\widehat{\tau_1 - \tau_2}) &= \tau_1 - \tau_2 + \mathbf{a}'\mathbf{P}_X\boldsymbol{\eta} \\ &= \tau_1 - \tau_2 + \mathbf{a}'\mathbf{P}_X \begin{bmatrix} 1/2 \\ -1/2 \\ 1/2 \\ -1/2 \end{bmatrix} \begin{bmatrix} \alpha \end{bmatrix} \\ &= \tau_1 - \tau_2 + \mathbf{a}'\mathbf{0} \begin{bmatrix} \alpha \end{bmatrix} \\ &= \tau_1 - \tau_2. \end{aligned}$$

Thus, the LSE of $\tau_1 - \tau_2$ is unbiased in this case.

Example 1 (continued)

From slide 10,

$$\begin{aligned}E(\hat{\sigma}^2) &= \frac{\boldsymbol{\eta}'(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}}{n - r} + \sigma^2 \\&= \frac{\boldsymbol{\eta}'(\mathbf{I}\boldsymbol{\eta} - \mathbf{P}_X\boldsymbol{\eta})}{n - r} + \sigma^2 \\&= \frac{\boldsymbol{\eta}'(\boldsymbol{\eta} - \mathbf{0})}{n - r} + \sigma^2 \\&= \frac{\boldsymbol{\eta}'\boldsymbol{\eta}}{n - r} + \sigma^2 \\&= \frac{\alpha^2}{4 - 2} + \sigma^2.\end{aligned}$$

Thus, $\hat{\sigma}^2$ is biased for σ^2 in this case.

Example 2

Once again consider an experiment with two experimental units (mice) for each of two treatments.

Suppose we assume the GMM holds with

$$E(\mathbf{y}) = E \left(\begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \end{bmatrix} \right) = \mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix} .$$

Example 2 (continued)

Suppose the person who conducted the experiment neglected to mention that both experimental units in treatment group 1 were female and that both experimental units in treatment group 2 were male.

Example 2 (continued)

Then the true model may require

$$\begin{aligned} E(\mathbf{y}) &= \begin{bmatrix} \mu + \tau_1 + \alpha/2 \\ \mu + \tau_1 + \alpha/2 \\ \mu + \tau_2 - \alpha/2 \\ \mu + \tau_2 - \alpha/2 \end{bmatrix} = \begin{bmatrix} \mu + \tau_1 \\ \mu + \tau_1 \\ \mu + \tau_2 \\ \mu + \tau_2 \end{bmatrix} + \begin{bmatrix} \alpha/2 \\ \alpha/2 \\ -\alpha/2 \\ -\alpha/2 \end{bmatrix} \\ &= \mathbf{X}\boldsymbol{\beta} + \begin{bmatrix} 1/2 \\ 1/2 \\ -1/2 \\ -1/2 \end{bmatrix} \begin{bmatrix} \alpha \end{bmatrix} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}\boldsymbol{\alpha} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}. \end{aligned}$$

Example 2 (continued)

If we analyze the data assuming the GMM with $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$, determine

- 1 $E(\widehat{\tau_1 - \tau_2})$, and
- 2 $E(\hat{\sigma}^2)$.

Example 2 (continued)

From slide 6,

$$\begin{aligned} E(\widehat{\tau_1 - \tau_2}) &= \tau_1 - \tau_2 + \mathbf{a}'\mathbf{P}_X\boldsymbol{\eta} \\ &= \tau_1 - \tau_2 + \mathbf{a}'\mathbf{P}_X \begin{bmatrix} 1/2 \\ 1/2 \\ -1/2 \\ -1/2 \end{bmatrix} \begin{bmatrix} \alpha \end{bmatrix} \\ &= \tau_1 - \tau_2 + \mathbf{a}' \begin{bmatrix} 1/2 \\ 1/2 \\ -1/2 \\ -1/2 \end{bmatrix} \begin{bmatrix} \alpha \end{bmatrix} = \tau_1 - \tau_2 + \alpha. \end{aligned}$$

Example 2 (continued)

Note that $\widehat{\tau_1 - \tau_2} = \bar{y}_{1\cdot} - \bar{y}_{2\cdot}$.

The previous slide shows that $\bar{y}_{1\cdot} - \bar{y}_{2\cdot}$ is not an unbiased estimator of the difference between treatment effects.

However, $\bar{y}_{1\cdot} - \bar{y}_{2\cdot}$ is an unbiased estimator of the difference between the means of the two treatment groups; i.e.,

$$E(\bar{y}_{1\cdot} - \bar{y}_{2\cdot}) = (\mu + \tau_1 + \alpha/2) - (\mu + \tau_2 - \alpha/2) = \tau_1 - \tau_2 + \alpha.$$

Part of the difference may be due to treatment, but part may be due to sex of the mice.

Example 2 (continued)

From slide 10,

$$\begin{aligned}E(\hat{\sigma}^2) &= \frac{\boldsymbol{\eta}'(\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}}{n - r} + \sigma^2 \\&= \frac{\boldsymbol{\eta}'(\mathbf{I}\boldsymbol{\eta} - \mathbf{P}_X\boldsymbol{\eta})}{n - r} + \sigma^2 \\&= \frac{\boldsymbol{\eta}'(\boldsymbol{\eta} - \boldsymbol{\eta})}{n - r} + \sigma^2 \\&= \sigma^2.\end{aligned}$$

Thus, $\hat{\sigma}^2$ is unbiased for σ^2 in this case.

Example 2 (continued)

Because $\boldsymbol{\eta} \in \mathcal{C}(X)$, both assumptions

$$E(\mathbf{y}) = X\boldsymbol{\beta} \quad \text{and} \quad E(\mathbf{y}) = X\boldsymbol{\beta} + \boldsymbol{\eta}$$

are equivalent to $E(\mathbf{y}) \in \mathcal{C}(X)$.

Thus, even though we ignore sex of the mice, our model for the mean is correct.

The only mistake we would make is to assume that the difference in means for the treatment groups is due only to treatment rather than to a combination of treatment and sex.

Overfitting

Now suppose we consider the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2] \quad \text{and} \quad \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} \quad \ni \quad \mathbf{X}\boldsymbol{\beta} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2.$$

Furthermore, suppose that (unknown to us) $\mathbf{X}_2\boldsymbol{\beta}_2 = \mathbf{0}$.

In this case, we say that we are overfitting.

Note that we are fitting a model that is more complicated than it needs to be.

To examine the impact of the overfitting, consider the case where $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$ is of full-column rank.

If we were to fit the simpler and correct model $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}$, the LSE of $\boldsymbol{\beta}_1$ is $\tilde{\boldsymbol{\beta}}_1 = (\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{y}$. Then

$$\begin{aligned} E(\tilde{\boldsymbol{\beta}}_1) &= (\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1E(\mathbf{y}) \\ &= (\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_1\boldsymbol{\beta}_1 \\ &= \boldsymbol{\beta}_1. \end{aligned}$$

$$\begin{aligned}\text{Var}(\tilde{\beta}_1) &= (\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\text{Var}(\mathbf{y})\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1} \\ &= \sigma^2(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1} \\ &= \sigma^2(\mathbf{X}'_1\mathbf{X}_1)^{-1}.\end{aligned}$$

If we were to fit the full model

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

that is correct but more complicated than it needs to be, then the LSE of $\boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix}$ is

$$\begin{aligned} \begin{bmatrix} \hat{\boldsymbol{\beta}}_1 \\ \hat{\boldsymbol{\beta}}_2 \end{bmatrix} &= ([\mathbf{X}_1, \mathbf{X}_2]'[\mathbf{X}_1, \mathbf{X}_2])^{-1} [\mathbf{X}_1, \mathbf{X}_2]'\mathbf{y} \\ &= \begin{bmatrix} \mathbf{X}'_1\mathbf{X}_1 & \mathbf{X}'_1\mathbf{X}_2 \\ \mathbf{X}'_2\mathbf{X}_1 & \mathbf{X}'_2\mathbf{X}_2 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}'_1\mathbf{y} \\ \mathbf{X}'_2\mathbf{y} \end{bmatrix}. \end{aligned}$$

If $X_1'X_2 = \mathbf{0}$, then

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} &= \begin{bmatrix} X_1'X_1 & \mathbf{0} \\ \mathbf{0} & X_2'X_2 \end{bmatrix}^{-1} \begin{bmatrix} X_1'y \\ X_2'y \end{bmatrix} \\ &= \begin{bmatrix} (X_1'X_1)^{-1}X_1'y \\ (X_2'X_2)^{-1}X_2'y \end{bmatrix} = \begin{bmatrix} \tilde{\beta}_1 \\ (X_2'X_2)^{-1}X_2'y \end{bmatrix}. \end{aligned}$$

Now suppose $\mathbf{X}'_1\mathbf{X}_2 \neq \mathbf{0}$.

$$\begin{aligned} E\left(\begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}\right) &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{y}) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \\ &= \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}. \end{aligned}$$

Thus, $E(\hat{\beta}_1) = \beta_1$.

$$\begin{aligned}\text{Var}(\hat{\beta}) &= \text{Var} \left(\begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} \right) \\ &= \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma^2 \begin{bmatrix} \mathbf{X}'_1\mathbf{X}_1 & \mathbf{X}'_1\mathbf{X}_2 \\ \mathbf{X}'_2\mathbf{X}_1 & \mathbf{X}'_2\mathbf{X}_2 \end{bmatrix}^{-1} .\end{aligned}$$

By Exercise A.72,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}\mathbf{E}^{-1}\mathbf{C}\mathbf{A}^{-1} & \mathbf{A}^{-1}\mathbf{B}\mathbf{E}^{-1} \\ -\mathbf{E}\mathbf{C}\mathbf{A}^{-1} & \mathbf{E}^{-1} \end{bmatrix},$$

where $\mathbf{E} = \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}$.

Thus, $\text{Var}(\hat{\beta}_1)$ is σ^2 times

$$\begin{aligned}
& (\mathbf{X}'_1\mathbf{X}_1)^{-1} + (\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_2(\mathbf{X}'_2\mathbf{X}_2 - \mathbf{X}'_2\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_2)^{-1}\mathbf{X}'_2\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1} \\
& = (\mathbf{X}'_1\mathbf{X}_1)^{-1} + (\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_2(\mathbf{X}'_2(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{X}_2)^{-1}\mathbf{X}'_2\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}.
\end{aligned}$$

Thus,

$$\text{Var}(\hat{\beta}_1) - \text{Var}(\tilde{\beta}_1) = \sigma^2(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_2(\mathbf{X}'_2(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{X}_2)^{-1}\mathbf{X}'_2\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}.$$

In a homework problem, you will show that

$$\text{Var}(\hat{\beta}_1) - \text{Var}(\tilde{\beta}_1) \text{ is NND.}$$

Thus, one cost of overfitting is increased variability of estimators of regression coefficients.

How is estimation of σ^2 affected?

Let $r_1 = \text{rank}(\mathbf{X}_1)$ and $r_2 = \text{rank}(\mathbf{X}_2)$.

If we fit a simpler model $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}$, then

$$\tilde{\sigma}^2 = \frac{\mathbf{y}'(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}}{n - r_1} \quad \text{and}$$

$$E(\mathbf{y}'(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1})\mathbf{y}) = (n - r_1)\sigma^2$$

$$\Rightarrow E(\tilde{\sigma}^2) = \sigma^2.$$

If we overfit with the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, then

$$\hat{\sigma}^2 = \frac{\mathbf{y}'(\mathbf{I} - \mathbf{P}_X)\mathbf{y}}{n - r} \quad \text{and}$$

$$E(\hat{\sigma}^2) = \sigma^2.$$

Thus, overfitting does not lead to biased estimation of σ^2 .

However, as we will see later in the course, overfitting leads to a loss of degrees of freedom ($n - r < n - r_1$), which can lead to a loss of power for testing hypotheses about β .