Foundations of Linear Regression

3. Properties and Assumptions

GOVT 6029 - Spring 2021

Cornell University

Is $(X'X)^{-1}X'y$ a good estimate of β ?

Is (X'X)⁻¹X'y a good estimate of β? Would another estimator be better?

Is (X'X)⁻¹X'y a good estimate of β? Would another estimator be better? What would an alternative be? Is (X'X)⁻¹X'y a good estimate of β? Would another estimator be better? What would an alternative be? Maybe minimizing the sum of absolute errors? Is (X'X)⁻¹X'y a good estimate of β? Would another estimator be better? What would an alternative be? Maybe minimizing the sum of absolute errors? Or something nonlinear? Is (X′X)⁻¹X′y a good estimate of β?

Would another estimator be better?

What would an alternative be?

Maybe minimizing the sum of absolute errors?

Or something nonlinear?

First we'll have to decide what makes an estimator good.

Bias

Bias The meaning of bias in statistics is more specific (and at times at variance) with plain English.

It does <u>not</u> mean subjectivity.

Bias The meaning of bias in statistics is more specific (and at times at variance) with plain English.

It does not mean subjectivity.

Is the estimate $\hat{\beta}$ provided by the model expected to equal the true value β ?

If not, how far off is it?

Bias The meaning of bias in statistics is more specific (and at times at variance) with plain English.

It does not mean subjectivity.

Is the estimate $\hat{\beta}$ provided by the model expected to equal the true value β ?

If not, how far off is it?

This is the **bias**, $E(\hat{\beta} - \beta)$

Bias The meaning of bias in statistics is more specific (and at times at variance) with plain English.

It does not mean subjectivity.

Is the estimate $\hat{\beta}$ provided by the model expected to equal the true value β ?

If not, how far off is it?

This is the **bias**, $E(\hat{\beta} - \beta)$

Although it seems "obvious" on face that we always prefer an unbiased estimator if one is available we also want the estimate to be close to the truth most of the time

They usually still miss the truth by some amount, But the direction in which they miss is not systematic or known ahead of time.

Unbiased estimates can be useless.

They usually still miss the truth by some amount, But the direction in which they miss is not systematic or known ahead of time.

Unbiased estimates can be useless. Why? One unbiased estimate of the time of day:

They usually still miss the truth by some amount, But the direction in which they miss is not systematic or known ahead of time.

Unbiased estimates can be useless. Why?

One unbiased estimate of the time of day:a random draw from the numbers 0–24. Utterly useless.

They usually still miss the truth by some amount, But the direction in which they miss is not systematic or known ahead of time.

Unbiased estimates can be useless. Why?

One unbiased estimate of the time of day:a random draw from the numbers 0–24. Utterly useless.

Biased estimates are not <u>necessarily</u> terrible.

A biased estimate of the time of day: a clock that is 2 minutes fast.

Efficiency:

Measures of efficiency answer the question: How much do we miss the truth by on average?

Efficiency thus incorporates both bias and variance of estimator.

Measures of efficiency answer the question: How much do we miss the truth by on average?

Efficiency thus incorporates both bias and variance of estimator.

A biased est with low variance may be "better" than an unbiased high var est

Measures of efficiency answer the question: How much do we miss the truth by on average?

Efficiency thus incorporates both bias and variance of estimator.

A biased est with low variance may be "better" than an unbiased high var est

Some examples:

Some examples:

Unbiased? Efficient?

Stopped clock. Random clock. Clock that is "a lot fast" Clock that is "a little fast" A well-run atomic clock

Some examples:

	Unbiased?	Efficient?
Stopped clock.	No	No
Random clock.	Yes	No
Clock that is "a lot fast"	No	No
Clock that is "a little fast"	No	Yes
A well-run atomic clock	Yes	Yes

To measure efficiency, we use **mean squared error**:

MSE =
$$E\left[\left(\beta - \hat{\beta}\right)^2\right]$$

= $Var(\hat{\beta}) + Bias(\hat{\beta}|\beta)^2$

 \sqrt{MSE} is how much you miss the truth by on average

In most cases, we want to use the estimator that minimizes MSE

We will be especially happy when this is also an unbiased estimator

But it won't always be

Consistency:

Consistency: An estimator that converges to the truth as the number of observations grows

Consistency: An estimator that converges to the truth as the number of observations grows Formally, $E(\hat{\beta} - \beta) \rightarrow 0$ as $N \rightarrow \infty$ Of great concern to many econometricians Not as great a concern in political science Consistency: An estimator that converges to the truth as the number of observations grows

Formally, $\mathrm{E}(\hat{eta} - eta)
ightarrow 0$ as $N
ightarrow \infty$

Of great concern to many econometricians

Not as great a concern in political science (as a thought experiment, $N \rightarrow \infty$ doesn't help much when the observations are, say, industrialized countries)

Consistency: An estimator that converges to the truth as the number of observations grows

Formally, $\mathrm{E}(\hat{eta} - eta)
ightarrow 0$ as $N
ightarrow \infty$

Of great concern to many econometricians

Not as great a concern in political science (as a thought experiment, $N \rightarrow \infty$ doesn't help much when the observations are, say, industrialized countries)

We will be mainly concerned with efficiency, secondarily with bias, and hardly at all with consistency

Two things that can go wrong:

- omitted variable bias
- specification bias

Two things that can go wrong:

- omitted variable bias
- specification bias

Another important source of bias is **selection**:

If we select observations non-randomly from the world, we may get biased estimates of means, regression coefficients, and other quantities Another important source of bias is **selection**:

If we select observations non-randomly from the world, we may get biased estimates of means, regression coefficients, and other quantities Another important source of bias is **selection**:

If we select observations non-randomly from the world, we may get biased estimates of means, regression coefficients, and other quantities

• Average children per marriage is 2.5. How many were in your family growing up? Are these numbers different? Who is "left out" in the second sample? Another important source of bias is **selection**:

If we select observations non-randomly from the world, we may get biased estimates of means, regression coefficients, and other quantities

- Average children per marriage is 2.5. How many were in your family growing up? Are these numbers different? Who is "left out" in the second sample?
- In testimony to NY state senate, motorcyclists testified that in their (multiple) crashes, helmets would not have prevented injuries. Who didn't testify?
- Regression example: Selection on the observed variables
Selection bias



Suppose we conducted a survey & asked people their income (x) and conservatism (y)

With the full range of respondents, we find a strong relationship

Selection bias



But suppose high income (or highly conservative) people decline to answer

Then we run a regression on the red dots only.

And get a result biased towards 0.

Selection bias



 \rightarrow Try to maximize variance of covariates, and avoid selecting on response variables

Most selection is unintentional, so think hard about sources of selection bias

Even if your data are sampled without bias from the population of interest, and your model correctly specified, several data problems can violate the linear regression assumptions Even if your data are sampled without bias from the population of interest, and your model correctly specified, several data problems can violate the linear regression assumptions

In order of declining severity:

Perfect collinearity Endogeneity of covariates Heteroskedasticity Serial correlation Non-normality

Lots of new jargon. Let's work through it.

Perfect collinearity occurs when $X^\prime X$ is singular; ie, the determinant $|X^\prime X|=0$

Happens when two or more columns of **X** are linearly dependent on each other

Perfect collinearity occurs when $X^\prime X$ is singular; ie, the determinant $|X^\prime X|=0$

Happens when two or more columns of **X** are linearly dependent on each other

Common causes: including a variable twice, or a variable and itself times a constant

Perfect collinearity occurs when $X^\prime X$ is singular; ie, the determinant $|X^\prime X|=0$

Happens when two or more columns of **X** are linearly dependent on each other

Common causes: including a variable twice, or a variable and itself times a constant

Very rare—except in panel data, as we will see

Matrix inversion—and thus LS regression—is impossible here

What if our covariates are correlated but not perfectly so? Then they are <u>not</u> linearly dependent

Then they are <u>not</u> linearly dependent

The regression coefficients are identified (a unique estimate exists for each β)

Then they are <u>not</u> linearly dependent

The regression coefficients are identified (a unique estimate exists for each β)

Regression with partial collinearity is unbiased & efficient.

But if the correlation among the X's is high, there is little to distinguish them

This leads to noisy estimates and large standard errors

Then they are <u>not</u> linearly dependent

The regression coefficients are identified (a unique estimate exists for each β)

Regression with partial collinearity is unbiased & efficient.

But if the correlation among the X's is high, there is little to distinguish them

This leads to noisy estimates and large standard errors

Those large se's are *correct*

Inappropriately, this situation is sometimes called "multicollinearity"

Inappropriately, this situation is sometimes called "multicollinearity"

Technically, multicollinearity describes only perfect linear dependence

Linear regression does not"fail" when correlation among X is "high".

Inappropriately, this situation is sometimes called "multicollinearity"

Technically, multicollinearity describes only perfect linear dependence

Linear regression does not"fail" when correlation among **X** is "high".

There is no "fix" for high correlation: it is not a statistical problem.

Have highly correlated **X** and large se's?

Inappropriately, this situation is sometimes called "multicollinearity"

Technically, multicollinearity describes only perfect linear dependence

Linear regression does not"fail" when correlation among X is "high".

There is no "fix" for high correlation: it is not a statistical problem.

Have highly correlated **X** and large se's? Then you lack sufficient data to precisely answer your research question

- So far, we have (implicitly) taken our regressors, X, as fixed X is not dependent on Y
- Fixed = pre-determined = exogenous

- So far, we have (implicitly) taken our regressors, X, as fixed X is not dependent on Y Fixed = pre-determined = exogenous
- Y consists of a function of X plus an error
- Y is thus endogenous to X

endogenous = "determined within the system"

What if **Y** helps determine **X** in the first place? That is, what if there is reciprocal causation?

What if **Y** helps determine **X** in the first place?

That is, what if there is reciprocal causation?

Very common in political science:

- campaign spending and share of the popular vote.
- policy attitudes and party identification
- arms races and war, etc.
- exchange rate policy and inflation

What if **Y** helps determine **X** in the first place?

That is, what if there is reciprocal causation?

Very common in political science:

- campaign spending and share of the popular vote.
- policy attitudes and party identification
- arms races and war, etc.
- exchange rate policy and inflation

In these cases, Y and X are both endogenous

Least squares is biased in this case

It will remain biased even as you add more data

In other words, it is inconsistent, or biased even as $N
ightarrow \infty$

Linear regression allows us to model the mean of a variable well

Y could be any linear function of $oldsymbol{eta}$ and ${f X}$

But LS always assumes the variance of that variable is the same:

 $\sigma^{\rm 2}$, a constant

We don't think **Y** has constant mean. Why expect constant variance?

In fact, heteroskedasticity—non-constant error variance—is very common



A common pattern of heteroskedasticity: Variance and mean increase together

Here, they are both correlated with the covariate X



In a fuzzy sense, X is a necessary but not sufficient condition for Y

This is usually an important point substantively. Heteroskedasticity is <u>interesting</u>, not just a nuisance



We can usually find heteroskedasticity by plotting the residuals against each covariate

Look for a pattern. Often a megaphone



But other patterns are possible.

Above, there is a dramatic difference in variance in different parts of the dataset



The same diagnostic reveals this problem.

Heteroskedasticity of this type often appears in panel datasets, where there are groups of observations from different units that each share a variance Every observation consists of a systematic component $(\mathbf{x}_i \boldsymbol{\beta})$ and a stochastic component (ε_i)

Generally, we can think of the stochastic component as an n-vector ϵ following a multivariate normal distribution:

 $\boldsymbol{\varepsilon} \sim \mathcal{MVN}(\boldsymbol{0},\boldsymbol{\Sigma})$

Aside: how the Multivariate Normal distribution works

Consider the simplest multivariate normal distribution, the joint distribution of two normal variables \mathbf{x}_1 and \mathbf{x}_2

As usual, let μ indicate a mean, and σ a variance or covariance

$$\begin{array}{ll} \mathsf{X} &=& \mathcal{MVN}(\boldsymbol{\mu},\boldsymbol{\Sigma}) \\ \begin{bmatrix} \mathsf{x}_1 \\ \mathsf{x}_2 \end{bmatrix} &=& \mathcal{MVN}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} \\ \sigma_{1,2} & \sigma_2^2 \end{bmatrix} \right) \end{array}$$

The MVN is more than the sum of its parts:

There is a mean and variance for each variable, <u>and</u> covariance between each pair



$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 0\\ 0 \end{array}\right], \left[\begin{array}{c} 1 & 0\\ 0 & 1 \end{array}\right]\right)$$



The standard MVN, with zero means, unit variances, and no covariance, looks like a higher dimension version of the normal: a symmetric mountain of probability



$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 0\\ 2 \end{array}\right], \left[\begin{array}{c} 1 & 0\\ 0 & 1 \end{array}\right]\right)$$



Shifting the mean of \mathbf{x}_2 moves the MVN in one dimension only Mean shifts affect only one dimension at a time



$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 2\\ 2 \end{array}\right], \left[\begin{array}{c} 1 & 0\\ 0 & 1 \end{array}\right]\right)$$



We could, of course, move the means of our variables at the same time.

This MVN says the most likely outcome is both x_1 and x_2 will be near 2.0


$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 0\\ 0 \end{array}\right], \left[\begin{array}{c} 0.33 & 0\\ 0 & 1 \end{array}\right]\right)$$



Shrinking the variance of \mathbf{x}_1 moves the mass of probability towards the mean of \mathbf{x}_1 , but leaves the distribution around \mathbf{x}_2 untouched



$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 0\\ 0 \end{array}\right], \left[\begin{array}{c} 0.33 & 0\\ 0 & 3 \end{array}\right]\right)$$



Increasing the variance of x_2 spreads the probability out, so we are less certain of x_2 , but just as certain of x_1 as before



$$\left[\begin{array}{c} \mathbf{x}_1 \\ \mathbf{x}_2 \end{array} \right] = \mathcal{MVN} \left(\left[\begin{array}{c} \mathbf{0} \\ \mathbf{0} \end{array} \right], \left[\begin{array}{c} \mathbf{0.33} & \mathbf{0} \\ \mathbf{0} & \mathbf{0.33} \end{array} \right] \right)$$



If the variance is small on all dimensions, the distribution collapses to a spike over the means of all variables



In this case, we are fairly certain of where all our variables tend to lie



$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 0\\ 0 \end{array}\right], \left[\begin{array}{c} 1 & 0.8\\ 0.8 & 1 \end{array}\right] \right)$$



In this special case, with unit variances, the covariance is also the correlation, so our distribution say x_1 and x_2 are correlated at r = 0.8



A positive correlation between our variables makes the MVN asymmetric, with greater mass on likely combinations



$$\left[\begin{array}{c} \mathbf{x}_1\\ \mathbf{x}_2 \end{array}\right] = \mathcal{MVN}\left(\left[\begin{array}{c} 0\\ 0 \end{array}\right], \left[\begin{array}{cc} 1 & -0.8\\ -0.8 & 1 \end{array}\right]\right)$$



A negative correlation makes <u>mismatched</u> values of our covariates more likely

In our current example, we have a huge multivariate normal distribuion:

each observation has its own mean and variance, and a covariance with every other observation

Suppose we have four observations. The Var-cov matrix of the disturbances is then

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \sigma_{14} \\ \sigma_{21} & \sigma_2^2 & \sigma_{23} & \sigma_{24} \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 & \sigma_{34} \\ \sigma_{41} & \sigma_{42} & \sigma_{43} & \sigma_4^2 \end{bmatrix}$$

In its most "ordinary" form, linear regression puts strict conditions on the variance-variance matrix, $\pmb{\Sigma}$

Again, assuming we have only four observations, the Var-cov matrix is

$$\Sigma = \sigma^2 \mathbf{I} = \begin{bmatrix} \sigma^2 & 0 & 0 & 0 \\ 0 & \sigma^2 & 0 & 0 \\ 0 & 0 & \sigma^2 & 0 \\ 0 & 0 & 0 & \sigma^2 \end{bmatrix}$$

Could treat each observation as consisting of $\mathbf{x}_i \boldsymbol{\beta}$ and a separate, univariate normal disturbance, each with the same variance, σ^2 .

This is the usual linear regression set up

Suppose the distrurbances are heteroskedastic.

Now each observation has an error term drawn from a Normal with its own variance

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & 0 & 0 \\ 0 & \sigma_2^2 & 0 & 0 \\ 0 & 0 & \sigma_3^2 & 0 \\ 0 & 0 & 0 & \sigma_4^2 \end{bmatrix}$$

Suppose the distrurbances are heteroskedastic.

Now each observation has an error term drawn from a Normal with its own variance

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 & 0 \\ 0 & \sigma_2^2 & 0 & 0 \\ 0 & 0 & \sigma_3^2 & 0 \\ 0 & 0 & 0 & \sigma_4^2 \end{bmatrix}$$

Still no covariance across disturbances.

Even so, we now have more parameters than we can estimate.

If every observation has its own unknown variance, we cannot estimate them

This MVN looks like the first example of a ridge: steeper in some directions than others, but not "tilted"

Heteroskedasticity does <u>not</u> bias least squares

- But LS is inefficient in the presence of heteroskedasticity
- More efficient estimators give greater weight to observations with low variance
- They pay more attention to the signal, and less attention to the noise

Heteroskedasticity does <u>not</u> bias least squares

But LS is inefficient in the presence of heteroskedasticity

More efficient estimators give greater weight to observations with low variance

They pay more attention to the signal, and less attention to the noise

Heteroskedasticity tends to make se's incorrect, because they depend on the estimate of σ^2

Researchers often try to "fix" standard errors to deal with this

(more on this later)

Suppose each disturbance has its own variance, and may be correlated with other disturbances

The most general case allows for both <u>heteroskedasticity</u> & <u>autocorrelation</u>

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \sigma_{14} \\ \sigma_{21} & \sigma_2^2 & \sigma_{23} & \sigma_{24} \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 & \sigma_{34} \\ \sigma_{41} & \sigma_{42} & \sigma_{43} & \sigma_4^2 \end{bmatrix}$$

LS is unbiased but inefficient in this case

The standard errors will be wrong, however

Key application: time series.

Current period is usually a function of the past

So when is least squares unbiased?

When is it efficient?

When are the standard errors correct?

#	Assumption	Formal statement	Consequence of violation
---	------------	------------------	--------------------------

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$\mathbf{E}(\mathbf{X}_{\mathbf{\mathcal{E}}}) = 0$	

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$E(X\varepsilon) = 0$	Biased, even as ${\it N} ightarrow\infty$

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$E(X \epsilon) = 0$	Biased, even as $N o \infty$

3 Disturbances have mean 0 $E(\varepsilon) = 0$

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$\mathrm{E}(X\varepsilon)=0$	Biased, even as $N ightarrow \infty$
3	Disturbances have mean 0	$\mathrm{E}(\boldsymbol{\varepsilon}) = 0$	Biased, even as $N o \infty$

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$E(X\varepsilon) = 0$	Biased, even as $N ightarrow \infty$
3	Disturbances have mean 0	$E(\varepsilon) = 0$	Biased, even as $N ightarrow \infty$
4	No serial correlation	$\mathrm{E}(\varepsilon_i\varepsilon_j)=0, i\neq j$	

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$\mathrm{E}(X\boldsymbol{\varepsilon})=0$	Biased, even as $N ightarrow \infty$
3	Disturbances have mean 0	$\mathrm{E}(oldsymbol{arepsilon})=0$	Biased, even as $N ightarrow \infty$
4	No serial correlation	$\mathbf{E}(\varepsilon_i\varepsilon_j)=0, i\neq j$	Unbiased but ineff. se's wrong

5 Homoskedastic errors $E(\varepsilon'\varepsilon) = \sigma^2 I$

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$\mathrm{E}(X\boldsymbol{\varepsilon})=0$	Biased, even as $N ightarrow \infty$
3	Disturbances have mean 0	$E(\boldsymbol{\varepsilon}) = 0$	Biased, even as $N ightarrow \infty$
4	No serial correlation	$\mathrm{E}(\varepsilon_i\varepsilon_j)=0, i\neq j$	Unbiased but ineff. se's wrong
5	Homoskedastic errors	$\mathrm{E}(\varepsilon'\varepsilon) = \sigma^2 I$	Unbiased but ineff. se's wrong

6 Gaussian error distrib $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

#	Assumption	Formal statement	Consequence of violation
1	No (perfect) collinearity	$\operatorname{rank}(X) = k, k < n$	Coefficients unidentified
2	X is exogenous	$\mathrm{E}(X\boldsymbol{arepsilon})=0$	Biased, even as $N ightarrow \infty$
3	Disturbances have mean 0	$\mathrm{E}(oldsymbol{arepsilon})=0$	Biased, even as $N ightarrow \infty$
4	No serial correlation	$\mathrm{E}(\varepsilon_i\varepsilon_j)=0, i\neq j$	Unbiased but ineff. se's wrong
5	Homoskedastic errors	$\mathrm{E}(\varepsilon'\varepsilon) = \sigma^2 \mathbf{I}$	Unbiased but ineff. se's wrong
6	Gaussian error distrib	$arepsilon \sim \mathcal{N}(0,\sigma^2)$	se's wrong unless $N ightarrow \infty$

(Assumptions get stronger from top to bottom, but 4 & 5 could be combined)

Gauss-Markov Theorem

It is easy to show $\beta_{\rm LS}$ is linear and unbiased, under Asps 1–3: If $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\mathrm{E}(\boldsymbol{\epsilon}) = 0$, then by substitution

$$\hat{\beta}_{\rm LS} = (X'X)^{-1}X'(X\beta + \varepsilon)$$

 $= \beta + (X'X)^{-1}X'\varepsilon$

So long as

- $(X'X)^{-1}$ is uniquely identified,
- \cdot X is exogenous or at least uncorrelated with arepsilon , and
- $E(\varepsilon) = 0$ (regardless of the distribution of ε)

Then $\mathsf{E}(\hat{oldsymbol{eta}}_{\mathrm{LS}})=oldsymbol{eta}$

 $\rightarrow \beta_{\rm LS}$ is unbiased and a linear function of ${\bf y}.$

If we make assumptions 1–5, we can make a stronger claim When there is no serial correlation, no heteroskedasticity, no endogeneity, and no perfect collinearity, then

Gauss-Markov holds that LS is the best linear unbiased estimator (BLUE)

BLUE means that among linear estimators that are unbiased, $\hat{\beta}_{\rm LS}$ has the least variance.

But, there might be a nonlinear estimator with lower MSE overall, unless ...

If in addition to Asp 1–5, the disturbances are normally distributed (6), then

Gauss-Markov holds LS is Minimum Variance Unbiased (MVU)