Econometrics I Lecture 3: Linear Regression

Paul T. Scott NYU Stern

Fall 2021

Paul T. Scott NYU Stern

L3 - Linear Regression

Fall 2021 1 / 64

글 🖌 🖌 글

Linear Regression: Introduction I

• A typical example of a linear regression model:

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

where

• i = 1, ..., n indexes observations

< □ > < 同 > < 回 > < 回 >

Linear Regression: Introduction I

• A typical example of a linear regression model:

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

where

- i = 1, ..., n indexes observations
- y_i , a scalar, is often referred to as the **dependent variable**

(4) (日本)

Linear Regression: Introduction I

• A typical example of a linear regression model:

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

where

- i = 1, ..., n indexes observations
- y_i , a scalar, is often referred to as the **dependent variable**
- x_{k,i}, is the *i*th observation of the *k*th **explanatory variable**, or **independent variable** (independent of what?), or **regressor**

< □ > < □ > < □ > < □ > < □ > < □ >

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

where

- i = 1, ..., n indexes observations
- y_i , a scalar, is often referred to as the **dependent variable**
- x_{k,i}, is the *i*th observation of the *k*th **explanatory variable**, or **independent variable** (independent of what?), or **regressor**
- The β_k terms represent the parameters.
 There are K parameters, one for each regressor.

< □ > < □ > < □ > < □ > < □ > < □ >

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

where

- i = 1, ..., n indexes observations
- y_i , a scalar, is often referred to as the **dependent variable**
- x_{k,i}, is the *i*th observation of the *k*th **explanatory variable**, or **independent variable** (independent of what?), or **regressor**
- The β_k terms represent the parameters.
 There are K parameters, one for each regressor.
- ε_i is the **disturbance** or **error term**

・ロト ・四ト ・ヨト ・ヨト

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

where

- i = 1, ..., n indexes observations
- y_i , a scalar, is often referred to as the **dependent variable**
- x_{k,i}, is the *i*th observation of the *k*th **explanatory variable**, or **independent variable** (independent of what?), or **regressor**
- The β_k terms represent the parameters.
 There are K parameters, one for each regressor.
- ε_i is the **disturbance** or **error term**
- Only the y_i and $x_{k,i}$ terms are observed by the econometrician.

イロト 不得 トイラト イラト 一日

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_K x_{K,i} + \varepsilon_i,$$

Today's questions:

- Where does it come from?
- What assumptions do we need to estimate β?
- How do we estimate β?
- How to interpret estimates?
- What are the estimator's (finite sample) properties?

Linear Regression: Matrix Notation

• We can express the linear regression model in vector notation

$$\mathbf{y} = \mathbf{X}oldsymbol{eta} + oldsymbol{arepsilon}$$
 ,

where

- **y** is a $n \times 1$ vector of observations of the dependent variable
- **X** is a $n \times K$ vector of observations of the dependent variable
- β is a $K \times 1$ vector of parameters
- ε is a $n \times 1$ vector of error terms
- Note that each row of this equation corresponds to the previous equation for a single observation *i*:

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i,$$

 Conventions: Roman symbols are observed, Greek are not, bold means vector notation, bold capitals means matrices.

Paul T. Scott NYU Stern

- Let's suppose we are interested in how worker's wages depend on education and experience (known as Mincerian regression or Mincer earnings function for Jacob Mincer).
- y_i could be worker *i*'s wages, and $\mathbf{x}'_i = (1, edu_i, exp_i)$, where
 - edu_i is worker i's education (in years)
 - exp_i is worker i's work experience (in years)
 - Note that the regressors include a constant

- y_i is *i*'s wages, and $\mathbf{x}'_i = (1, edu_i, exp_i)$
- The data matrices for the Mincerian regression might look like this:

$$\mathbf{y} = \begin{pmatrix} 12\\ 35\\ 20\\ \vdots \end{pmatrix} \qquad \qquad \mathbf{X} = \begin{pmatrix} 1 & 12 & 2\\ 1 & 16 & 5\\ 1 & 12 & 21\\ \vdots & \vdots & \vdots \end{pmatrix}$$

< □ > < 同 > < 回 > < Ξ > < Ξ

$$\mathbf{y} = \mathbf{X} \boldsymbol{eta} + \boldsymbol{arepsilon}$$

- This equation doesn't mean much until we say something about the error term ε .
- The first (and strongest) assumption on error terms we will consider is **strict exogeneity**:

$$E\left[arepsilon|\mathbf{X}
ight]=\mathbf{0}$$

• Note that the law of iterated expectations implies

$$E\left[\varepsilon_{i}\right]=E_{\mathbf{X}}\left[E\left[\varepsilon_{i}|\mathbf{X}\right]\right]=0.$$

It also implies that the error terms are uncorrelated with the regressors: $Cov[\varepsilon_i, \mathbf{X}] = 0$ and $Cov[\varepsilon_i, \mathbf{x}_i] = 0$.

< 回 > < 三 > < 三

Strict Exogeneity and Conditional Means

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{1}$$
$$\boldsymbol{E} \left[\boldsymbol{\varepsilon} | \mathbf{X} \right] = 0 \tag{2}$$

• Now, we have a meaningful model.

Note that

$$E[\mathbf{y}|\mathbf{X}] = E[\mathbf{X}\beta + \varepsilon|\mathbf{X}]$$
$$= \mathbf{X}\beta + E[\varepsilon|\mathbf{X}]$$
$$= \mathbf{X}\beta,$$

so equations (1) and (2) already imply that the conditional mean is a linear function of X.

Paul T. Scott NYU Stern

Strict Exogeneity Interpretation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{1}$$
$$\boldsymbol{E}\left[\boldsymbol{\varepsilon}|\mathbf{X}\right] = 0 \tag{2}$$

- Strict exogeneity captures the idea that x is varied in the data without changing the mean of the unobservable factors affecting y.
- Strict exogeneity is very plausible in the context of experimental variation (especially in double blind studies)
- Unfortunately, social scientists are sometimes unable to rely on experimental variation, and strict exogeneity is rarely plausible in the context of naturally occurring data. For this reason, we will consider different assumptions on the error terms, but this is the starting point.

Omitted Variables and Endogeneity

• Suppose we estimate the following Mincerian regression:

$$\ln(wage_i) = \beta_0 + \beta_1 edu_i + \beta_2 exp_i + \beta_3 exp_i^2 + \varepsilon_i.$$

• Furthermore suppose the true model is:

$$\ln(wage_i) = \beta_0 + \beta_1 edu_i + \beta_2 exp_i + \beta_3 exp_i^2 + \beta_4 ability_i + \varepsilon_i$$

• Then, when estimating the first equation (because ability is unobserved), the error term is effectively:

$$\widetilde{\varepsilon}_i = \beta_4 ability_i + \varepsilon_i$$

Note that even if ε_i satisfies strict exogeneity, ε̃_i will not if, for instance, ability is correlated with education. We say x is endogenous if E [ε|x] ≠ 0.

・ 何 ト ・ ヨ ト ・ ヨ ト

Wherefore Linearity?

• When would
$$E[y_i | \mathbf{x_i}] = \mathbf{x}'_i \boldsymbol{\beta}$$
 be true?

• We might start with the conditional mean as a general function of x

$$y_i = f(\mathbf{x}_i) + \varepsilon_i,$$

and then what we've done is impose that f is a linear function.

- We might also think of the model as a linear approximation (using Taylor's theorem).
 - Actually, it can be a polynomial approximation, not just a linear approximation . . .

・ 何 ト ・ ヨ ト ・ ヨ ト

- The linear regression framework does not impose that y is a linear function of any particular variable x
- The regressors in **x**_i can include squared terms, higher powers, and other functions of a variable x
- For example, $\mathbf{x}'_i = (1, edu_i, exp_i, exp_i^2)$ in the Mincerian regression would allow declining (or increasing) returns to work experience.
- The "Linear" part of "Linear Regression" really means linear-in-parameters, which is much less restrictive than being linear with respect to a particular variable.

< □ > < 同 > < 回 > < 回 > < 回 >

Linear-in-Parameters II

- The dependent variable can also involve nonlinear transformations.
- For example, y_i in the Mincerian regression is typically the natural log of worker *i*'s wage:

$$\ln(wage_i) = \beta_0 + \beta_1 edu_i + \beta_2 exp_i + \beta_3 exp_i^2 + \varepsilon_i.$$

• Models of demand (or supply) sometimes have the form

$$\ln(q_i) = \beta_0 + \beta_1 \ln(p_i) + \varepsilon_i,$$

in which case β_1 represents the price elasticity of demand (supply) and does not depend on the units that prices and quantities are measured in (Check this). Economists love logs!

• What would it take for a model to not be linear in parameters?

Paul T. Scott NYU Stern

< □ > < □ > < □ > < □ > < □ > < □ >

- Another way to motivate the linear regression function is from the multivariate normal distribution.
- Suppose $(y, x)' \sim \mathcal{N}(\mu, \Sigma)$, where

$$\mu = \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix} \qquad \Sigma = \begin{pmatrix} \sigma_y^2 & \rho \sigma_y \sigma_x \\ \rho \sigma_y \sigma_x & \sigma_x^2 \end{pmatrix}$$

(日) (四) (日) (日) (日)

Multivariate Normal Data



Paul T. Scott NYU Stern

L3 - Linear Regression

Fall 2021 15 / 64

Multivariate Normal Marginal Densities



 Here, I approximate the conditional density Pr (y | x = 0) by plotting the density of y when we focus on observations with x ∈ (-.04, .04)

Paul T. Scott NYU Stern

Fall 2021 16 / 64

R code

```
# install.packages('MASS')
# install.packages('ggplot2')
# install.packages('reshape')
library(MASS)
library(ggplot2)
library(reshape)
mu < - c(0, 0)
Sigma <- matrix(c(1, -1, -1, 4), 2, 2)
xy <- mvrnorm(n = 5000, mu, Sigma)</pre>
plot(xy, xlab="x",ylab="y")
xy <- mvrnorm(n = 100000, mu, Sigma)
xy.df <- as.data.frame(xy)</pre>
names(xy.df) <- c("x", "y")
xy.stacked <- melt(xy.df)</pre>
ggplot(xy.stacked, aes(value, fill = variable)) + geom_density(
    alpha = 0.2)
xy.selected <- xy.df[xy.df$x>-.04 & xy.df$x<.04, ]</pre>
ggplot(xy.selected, aes(y)) + geom_density(alpha = 0.2)
                                               Paul T. Scott NYU Stern
                              L3 - Linear Regression
                                                             Fall 2021
```

17 / 64

Multivariate Normal and Regression Equation

• Suppose $(y, x)' \sim \mathcal{N}(\mu, \Sigma)$, where

$$\mu = \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix} \qquad \qquad \Sigma = \begin{pmatrix} \sigma_y^2 & \rho \sigma_y \sigma_x \\ \rho \sigma_y \sigma_x & \sigma_x^2 \end{pmatrix}$$

• It's also the case that

$$E(y|x) = \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x)$$

and

$$y = E(y|x) + \varepsilon$$

with ε normally distributed.

• Normally distributed error terms: the first case we will consider (and easiest to analyze)

Paul T. Scott NYU Stern

→ Ξ →

Our next assumption is that the matrix of regressors has full rank
X is a n × K matrix with rank K.

(日) (四) (日) (日) (日)

- Our next assumption is that the matrix of regressors has full rank
 X is a n × K matrix with rank K.
- A model of consumption that violates full rank:

 $C = \beta_1 + \beta_2$ Salary + β_3 Nonsalary income + β_4 Total income + ε

(日) (四) (日) (日) (日)

- Our next assumption is that the matrix of regressors has full rank
 X is a n × K matrix with rank K.
- A model of consumption that violates full rank:

 $C = \beta_1 + \beta_2$ Salary + β_3 Nonsalary income + β_4 Total income + ε

• Conditional on total income, any increase in salary must be met by a proportional decrease in nonsalary income

< □ > < □ > < □ > < □ > < □ > < □ >

Full Rank II

• Consider the following model of consumption:

 $C = \beta_1 + \beta_2 \text{Salary} + \beta_3 \text{Nonsalary income} + \beta_4 \text{Total income} + \varepsilon$

Total income = Salary + Nonsalary income,

SO

 $\begin{array}{lll} \mathcal{C} &=& \beta_1 + (\beta_2 + \beta_4) \operatorname{Salary} + (\beta_3 + \beta_4) \operatorname{Nonsalary\,Income} + \varepsilon \\ \mathcal{C} &=& \beta_1 + \widetilde{\beta}_2 \operatorname{Salary} + \widetilde{\beta}_3 \operatorname{Nonsalary\,Income} + 0 \cdot \operatorname{Total\,Income} + \varepsilon \end{array}$

• Assuming $\beta_4 \neq 0$ in the original equation, we have constructed an empirically equivalent equation with different parameters. That is, for this model, we have different values for β that are *observationally* equivalent.

イロト 不得 トイラト イラト 一日

- Exercise: suppose I wanted to build a model of the approval ratings of major party nominees for US president.
- I want to include the following regressors:
 - Years holding elected office
 - Age
 - Gender
 - Indicator variable for being married to a former president
- What's the problem? Compare to the previous case with salaries any difference?

Linear Independence

- Another term for the full rank assumption (*Rank* (X) = K) is linear independence. Multicollinearity refers to a lack of linear independence. When a model has multicollinearity, we say it is not identified.
- Note that this is distinct from statistical independence.
- Linear independence means that one variable cannot be algebraically predicted by others.
- Statistical independence means that the variation in two random variables is unrelated. Linear independence does not imply statistical independence. Why not? Does statistical independence imply linear independence?

< ロ > < 同 > < 回 > < 回 > < 回 > <

Assumptions so far

• The assumptions we have introduced so far:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$
(1)

$$E\left[\boldsymbol{\varepsilon}|\mathbf{X}\right] = 0$$
(2)

$$Rank(\mathbf{X}) = K$$
(3)

- Note that we have not made any assumptions on the statistical properties of X. There is no need to do so; X can be fixed or random *for now*. What matters is the distribution of the error terms conditional on X.
- Not all assumptions will be maintained throughout the course (or even this lecture). When we consider formal results, I will be explicit about which assumptions are needed.

Homoscedasticity

• Another important assumption is homoscedasticity:

$$E\left[\varepsilon\varepsilon'|\mathbf{X}\right] = \sigma^2\mathbf{I}$$

where **I** is a $n \times n$ identity matrix.

• Given that $E[\varepsilon|\mathbf{X}] = 0$, this means that

$$Var\left[\varepsilon\right] = \sigma^2 \mathbf{I}$$

- In words, homoscedasticity says that each error term has the same variance; i.e., the variance of ε_i is not related to x_i.
 Heteroscedastcity is the alternative. Can you think of some cases of heteroscedasticity?
- The assumption also rules out correlation between the error terms for different observations.

• A convenient assumption is that error terms are normally distributed:

$$oldsymbol{arepsilon} | oldsymbol{X} \sim \mathcal{N}\left(oldsymbol{0}, \sigma^2 oldsymbol{\mathsf{I}}
ight)$$

• Where are we going?

- This assumption will make it easy to derive normally distributed estimators
- Ultimately, the central limit theorem can be used to derive asymptotic normality of estimators (that is, estimators that will be normally distributed as the sample gets large), so in practice it's rarely necessary to assume normal error terms.

- We used β to denote the true vector of parameters; let's use b to denote estimates of or candidates for β.
- A residual is the fitted value given a particular potential parameter vector:

$$e_i\left(\mathbf{b}\right) = y_i - \mathbf{x}'_i \mathbf{b}$$

- A residual captures the degree to which the linear prediction x'_ib explains the dependent variable y_i.
- Formally, we should think of residuals as being a function of candidate parameter vectors, but we will often just write *e_i*.

A B b A B b

Least Squares

- It makes intuitive sense to want to find a value of **b** that makes residuals small, so that the estimated model explains the data well.
- There are many possible ways to think about making the residuals small, but by far the most popular criterion is the **sum of squared residuals**:

SSR (**b**) =
$$\sum_{i=1}^{n} e_i (b)^2 = \sum_{i=1}^{n} (y_i - \mathbf{x}'_i \mathbf{b})^2 = \mathbf{e} (\mathbf{b})' \mathbf{e} (\mathbf{b})$$

• The least squares estimator minimizes the sum of squared residuals:

$$\hat{\mathbf{b}}_{LS} \equiv rg\min_{\mathbf{b}} SSR\left(\mathbf{b}
ight)$$

• For linear models with homoscedastic errors, the least squares estimate is typically called the **Ordinary Least Squares** estimator.

Paul T. Scott NYU Stern

L3 - Linear Regression

Fall 2021 27 / 64

• Let's consider a model with only a single variable regressor (and a constant):

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

- Similar to the example in the introduction, we might have
 - ▶ y_i: Dependent Variable (e.g., test score of student i)
 - x_i: Independent Variable (e.g., class size of student i)
 - ε_i : regression error (e.g., noise in the model)

・ 何 ト ・ ヨ ト ・ ヨ ト

Derivation of Bivariate Linear Regression Estimator

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$SSR(b_0, b_1) = \sum_{i=1}^{n} e_i (b_0, b_1)^2 = \sum_{i=1}^{n} (y_i - b_0 - b_1 x_i)^2$$

• (On board) What are the OLS estimates of (β_0, β_1) ?

Image: A matched black
OLS Estimator

• The OLS estimates are:

$$\hat{b}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{s_{x,y}}{s_x^2}$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

where the s_x^2 refers to the sample variance of x and $s_{x,y}$ refers to the sample covariance of x and y.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- Now let's go ahead and actually use the OLS estimator on Alice and Bob (recall the intro lecture)
- Suppose small classes have 15 students in them, large classes 20 students
- Our data is:
 - Alice: $(y_1, x_1) = (6, 15)$
 - Bob: $(y_2, x_2) = (3, 20)$
- Note that this is a riduclously small sample size (the smallest we could have and still solve for the OLS estimator).

・ 何 ト ・ ヨ ト ・ ヨ ト

Implementing the OLS Estimator

- (On board) Find the OLS estimator for β_0 and β_1 for the data $(y_1, x_1) = (6, 15)$ and $(y_2, x_2) = (3, 20)$
- Formulas:

$$\hat{b}_{1} = \frac{\sum_{i=1}^{n} (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} = \frac{cov(x, y)}{var(x)}$$
$$\hat{b}_{0} = \bar{y} - \hat{b}_{1}\bar{x}$$

• You will never have to do this algebra yourself. It becomes very cumbersome with lots of observations and with lots of regressors, which we now turn to.

< □ > < □ > < □ > < □ > < □ > < □ >

Multiple Linear Regression

• Let's return to the multiple linear regression setting (multiple regressors):

$$\mathbf{y}=\mathbf{X}m{eta}+m{arepsilon}$$
 ,

• We can write the sum of squared residuals as follows:

$$SSR(\mathbf{b}) = \mathbf{e}_{i}(\mathbf{b})' \mathbf{e}_{i}(\mathbf{b}) = \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\mathbf{b} + \mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b}$$

• The necessary condition for a minimum:

$$\frac{\partial SSR(\mathbf{b})}{\partial \mathbf{b}} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\mathbf{b} = 0$$

• Recalling the rules:

$$\frac{\partial \mathbf{u}' \mathbf{v}}{\partial \mathbf{v}} = \mathbf{u} \qquad \qquad \frac{\partial \mathbf{v}' \mathbf{A} \mathbf{v}}{\partial \mathbf{v}} = 2 \mathbf{A} \mathbf{v}$$

Paul T. Scott NYU Stern

Fall 2021 33 / 64

The OLS Formula

• The necessary condition for a minimum:

$$\frac{\partial SSR\left(\mathbf{b}\right)}{\partial \mathbf{b}} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\mathbf{b} = 0$$

• Thus, **b** must satisfy

$$\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{y}$$

• Assuming X'X is invertible (be careful!), we have

$$\widehat{\mathbf{b}}_{OLS} \equiv \left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{y}$$

• Note the similarity to the bivariate least squares estimator.

Full Rank and Identification

- The invertibility of X'X is not guaranteed in general, but it is implied by the full rank condition: Rank (X) = K.
- When ${\bf X}$ does not have full rank, neither does ${\bf X}'{\bf X},$ and

$$\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{y}$$

can be solved by multiple values of **b**.

• Linear algebra review: when a square matrix **A** is not invertible, it has a non-trivial nullspace. This means that

$\mathbf{A}\mathbf{b}=\mathbf{0}$

can be solved by multiple vectors \mathbf{b} . This implies that, for any \mathbf{c} ,

$\mathbf{A}\mathbf{b} = \mathbf{c}$

also has multiple solutions if it has at least one solution.

Paul T. Scott NYU Stern

Properties of Residuals

 $\bullet\,$ Let e denote the vector of OLS residuals, and consider

$$\mathbf{X}'\mathbf{e} = \mathbf{X}' \left(\mathbf{X}\mathbf{b}_{OLS} - \mathbf{y}\right)$$

Substituting in the OLS formula,

$$\mathbf{X}'\mathbf{e} = \mathbf{X}'\left(\mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{y} - \mathbf{y}\right)$$

• This simplifies to

$$\mathbf{X}'\mathbf{e} = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{y} = 0$$

 This implies that (1) the OLS residuals sum to zero, given that one of the regressors is a constant, and (2) the OLS residuals are uncorrelated with the regressors. Intuitively: there is no variation left in e that can be explained by X. • We saw that for a bivariate regression, the slope on the regressor is given by

$$\hat{b}_1 = rac{s_{X,Y}}{s_X^2}$$

• Does it follow that, with multiple regressors, the coefficient on the *k*th regressor is

$$\hat{b}_k = \frac{s_{X_k, y}}{s_{X_k}^2}?$$

• If not, under what conditions?

• Consider the following model of rectangular painting area:

$$\ln A_i = \beta_0 + \beta_1 \ln W_i + \beta_2 \ln H_i + \varepsilon_i$$

where

- W_i is the painting i's width
- ► *H_i* is the painting *i*'s width
- A_i is the painting *i*'s area, $A_i = W_i \cdot H_i$
- ε_i is measurement error in painting *i*'s area
- What should the β 's be in theory (i.e., given what you know about geometry)?

(4) (日本)

Silly Example II

• Suppose $\ln W_i$ and $\ln H_i$ are correlated in the population (naturally)

$$\begin{pmatrix} \ln W_i \\ \ln H_i \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_w \\ \mu_h \end{pmatrix}, \begin{pmatrix} \sigma_w^2 & \sigma_{wh} \\ \sigma_{wh} & \sigma_h^2 \end{pmatrix} \right)$$

• Then, notice that

$$Cov (\ln A_i, \ln W_i) = Cov (\ln W_i + \ln H_i, \ln W_i)$$

= Var (ln W_i) + Cov (ln W_i, ln H_i)

• Finally,

$$\frac{Cov\left(\ln A,\ln W\right)}{Var\left(\ln W\right)} = \frac{Var\left(\ln W\right) + Cov\left(\ln H,\ln W\right)}{Var\left(\ln W\right)} = 1 + \frac{Cov\left(\ln H,\ln W\right)}{Var\left(\ln W\right)}$$

→ ∃ →

Silly Example III

• In a bivariate regression, recall the slope would be given by

$$\frac{\textit{Cov}(\ln\textit{A},\ln\textit{W})}{\textit{Var}(\ln\textit{W})} = 1 + \frac{\textit{Cov}(\ln\textit{H},\ln\textit{W})}{\textit{Var}(\ln\textit{W})},$$

so the term $\frac{Cov(\ln H, \ln W)}{Var(\ln W)}$ represents bias.

- Implication: if the OLS coefficients in a multiple regression had the same formula as the coefficients in a bivariate regression, there would be bias.
- Exception: if the regressors ln *H* and ln *W* are uncorrelated, there will be no bias above, and OLS with multiple regressors delivers the same coefficients as if we ran a bivariate regression with each of the regressors separately.
- This example also makes a point about omitted variables bias: if height is unobserved, the regression of ln *A* on only ln *W* will deliver the biased coefficient above.

Paul T. Scott NYU Stern

L3 - Linear Regression

Fall 2021 40 / 64

$$\widehat{\mathbf{b}}_{OLS} \equiv \left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{y}$$

- Indeed, OLS with multiple regressors does not deliver the bivariate regression coefficient $s_{x_ky}/s_{x_k}^2$ for each regressor x_k (except in the case of orthogonal regressors)
- To gain some intuition for what this formula is doing, let's consider the (X'X) and X'y pieces separately in the context of the silly model of painting area.

- Suppose that $w_i = \ln W_i \mu_w$ and $\ln h_i = \ln H_i \mu_h$
- Let $\mathbf{X} = \begin{bmatrix} \mathbf{w} & \mathbf{h} \end{bmatrix}$
- It follows that

$$\mathbf{X}'\mathbf{X} = \begin{pmatrix} \sum (\ln W_i - \mu_w)^2 & \sum (\ln W_i - \mu_w) (\ln H_i - \mu_h) \\ \sum (\ln W_i - \mu_w) (\ln H_i - \mu_h) & \sum (\ln H_i - \mu_h)^2 \end{pmatrix}$$

Note: if we multiply that by n^{-1} , we have the sample analog of the covariance matrix.

A D N A B N A B N A B N

• Similarly,

$$\mathbf{X'y} = \left(\begin{array}{c} \sum \left(\ln W_i - \mu_w \right) \left(\ln A_i - \mu_a \right) \\ \sum \left(\ln H_i - \mu_h \right) \left(\ln A_i - \mu_a \right) \end{array} \right),$$

is the sample covariance of \mathbf{x} and \mathbf{y} (times n).

• Therefore, we have

$$\widehat{\mathbf{b}}_{OLS} = \left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{y} = \left(n^{-1}\mathbf{X}'\mathbf{X}\right)^{-1}\left(n^{-1}\mathbf{X}'\mathbf{y}\right) = S_{xx}^{-1}s_{xy}$$

where S_{xx} is the sample covariance matrix for **x** and s_{xy} is the sample covariance of **x** and *y*.

< □ > < 同 > < 回 > < 回 > < 回 >

Silly Example VII

$$\widehat{\mathbf{b}}_{OLS} = S_{xx}^{-1} s_{xy}$$

• The population moments for the model of painting areas are

$$Var\left(\mathbf{x}\right) = \begin{pmatrix} \sigma_{w}^{2} & \sigma_{wh} \\ \sigma_{wh} & \sigma_{h}^{2} \end{pmatrix} \qquad Cov\left(\mathbf{x}, y\right) = \begin{pmatrix} \sigma_{w}^{2} + \sigma_{wh} \\ \sigma_{h}^{2} + \sigma_{wh} \end{pmatrix}$$

Note that

$$\left(\operatorname{Var}\left(\mathbf{x}\right)\right)^{-1} = \frac{1}{\sigma_{w}^{2}\sigma_{h}^{2} - \sigma_{wh}^{2}} \left(\begin{array}{cc} \sigma_{h}^{2} & -\sigma_{wh} \\ -\sigma_{wh} & \sigma_{w}^{2} \end{array}\right)$$

 With some algebra, the population-moments version of OLS gives us the right coefficients:

$$(Var(\mathbf{x}))^{-1} Cov(\mathbf{x}, y) = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Paul T. Scott NYU Stern

The Frisch-Waugh Theorem I

• Think about separating **X** into two sub-matrices:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix}$$
 ,

with

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

Frisch-Waugh Theorem

The OLS regression of **y** on $[X_1, X_2]$ yields a subvector **b**₂ of coefficient estimates that is the same as the result from a regression of the residuals from a regression of **y** on X_1 are regressed on the residuals from a regression of X_2 on X_1 .

- In other words, let's start by regressing \mathbf{y} on \mathbf{X}_1 . Let's label the residuals from this regression \mathbf{y}^* .
- Let's also regress start by regressing X_2 on X_1 (think about regressing each column of X_2 on X_1). Let's label the residuals from this regression X_2^* .
- If we regress y^{*} on X^{*}₂, we get the same coefficient on X^{*}₂ that we would have had in the full regression of y on [X₁, X₂].
- An implication is that the coefficient on each variable can be thought of as the effect of that variable after controlling for all the other variables. Thus, OLS coefficients are sometimes called **partial regression coefficients**.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- Bivariate linear regression is easy to visualize (scatter plot with a line running through). Frisch-Waugh tells us how we can visualize the effect of a single variable from a multiple regression.
- One implication of Frisch-Waugh is that, if we de-mean all the variables and then run a regression with the de-meaned variables (but leaving out the constant term in **X**), we will get the same coefficients on all the variables.

Units and Coefficients

• In simple linear regression, it's easy to see that re-scaling (changing the units of x) will rescale the parameter estimates in the opposite way:

$$\hat{b}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

- If you multiply each value of x by λ , the numerator gets multiplied by λ , and the denominator by λ^2 , meaning \hat{b}_1 gets divided by λ . The same is true in the multiple regression framework (Frisch-Waugh makes this easier to see).
- Exercise: if the regressor x in a bivariate regression is the log of a variable, show that rescaling the original variable does not affect \hat{b}_1 . What about \hat{b}_0 ?

Paul T. Scott NYU Stern

• The **total sum of squares**, or the **total variation** in the dependent variable:

$$SST = (\mathbf{y} - \mathbf{i}\overline{y})'(\mathbf{y} - \mathbf{i}\overline{y})$$

• The dependent variable decomposes into a prediction and a residual:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e} = \widehat{\mathbf{y}} + \mathbf{e}$$

where $\mathbf{X}\mathbf{b} = \widehat{\mathbf{y}}$ is the prediction or **fitted value** of *y*.

• We can think about decomposing the variation in y into variation in \hat{y} and **e**. Intuitively, we want the variation in \hat{y} to account for as much as possible

< □ > < □ > < □ > < □ > < □ > < □ >

Define

$$\mathbf{M}_0 = I - n^{-1} \mathbf{i} \mathbf{i}'.$$

note that M_0 is symmetric and **idempotent** (on board).

- Notice that $\mathbf{y}'\mathbf{M}_0\mathbf{y}$ is the SST.
- Also, we can show that

$$\mathbf{y}'\mathbf{M}_0\mathbf{y} = \mathbf{b}'\mathbf{X}'\mathbf{M}_0\mathbf{X}\mathbf{b} + \mathbf{e}'\mathbf{e}$$
,

where the first term on the RHS is known as the **regression sum of** squares (SSR), and the second term is the **error sum of squares** (SSE).

$$SST = SSR + SSE$$

Coefficient of Determination

• The coefficient of determination is defined as the proportion of the variation in the dependent variable explained by the model:

$$R^{2} = \frac{SSR}{SST} = \frac{\mathbf{b}' \mathbf{X}' \mathbf{M}_{0} \mathbf{X} \mathbf{b}}{\mathbf{y}' \mathbf{M}_{0} \mathbf{y}} = 1 - \frac{\mathbf{e}' \mathbf{e}}{\mathbf{y}' \mathbf{M}_{0} \mathbf{y}}$$

- R² is a measure of goodness of fit that always goes up as we add more regressors. It doesn't tell us whether it's "worth it" to add a new regressor to the model. (Later we will talk about why overfitting can be bad in finite samples.)
- Adjusted R² is more useful for model selection because it incorporates a penalty for the number of regressors:

$$\overline{R}^{2} = 1 - \frac{\mathbf{e}'\mathbf{e}/(n-K)}{\mathbf{y}'\mathbf{M}_{0}\mathbf{y}/(n-1)}$$

• Given assumptions (1)-(3), homoscedasticity, and normal error terms,

$$\mathbf{b}_{OLS} | \mathbf{X} \sim \mathcal{N} \left(\boldsymbol{\beta}, \sigma^2 \left(\mathbf{X}' \mathbf{X} \right)^{-1} \right)$$

- Proof on board
- Note: proof that it's unbiased and expression for variance do not require normally error terms, but it would be hard to say what the finite sample distribution is, exactly, without normality.

Omitted Variables Bias

 Suppose the econometrician only observes regressors X, but the true model is

$$\mathbf{y}=\mathbf{X}oldsymbol{eta}+\mathbf{z}\gamma+arepsilon$$
 ,

• The OLS estimator will equal

$$\mathbf{b} = \left(\mathbf{X}'\mathbf{X}
ight)^{-1}\mathbf{X}'\mathbf{y} = oldsymbol{eta} + \left(\mathbf{X}'\mathbf{X}
ight)^{-1}\mathbf{X}'\mathbf{z}\gamma + \left(\mathbf{X}'\mathbf{X}
ight)^{-1}\mathbf{X}'arepsilon$$

- The last term is mean zero given the strict exogeneity assumption.
- Note that the second term will not be zero if **X** and **z** are correlated; i.e. if $\mathbf{X}'\mathbf{z} \neq \mathbf{0}$.
- Implication: correlation between omitted variables and the observed regressors makes OLS biased.

Bias with multiple variables

• Consider a model with two variables:

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \varepsilon_i,$$

and suppose that

$$x_{1,i} = u_{1,i} + u_{3,i} + \varepsilon_i x_{2,i} = u_{2,i} + u_{3,i}$$

where the *u*'s and ε are all independently distributed.

 We know that OLS will deliver a biased estimate of β₁, but will OLS still be consistent for β₂?

(4) (3) (4) (4) (4)

Bias with multiple variables

• Consider a model with two variables:

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \varepsilon_i,$$

and suppose that

$$x_{1,i} = u_{1,i} + u_{3,i} + \varepsilon_i x_{2,i} = u_{2,i} + u_{3,i}$$

where the *u*'s and ε are all independently distributed.

- We know that OLS will deliver a biased estimate of β₁, but will OLS still be consistent for β₂?
- No! The correlation between x_1 and x_2 leads to bias even in β_2 .
 - Thus, endogeneity problems are a big deal not only for our variables of interest, but also when it comes to control variables.

• Consider "controlling" for x_1 with the wrong estimate b_1 of β_1 :

$$y_i - b_1 x_{1,i} = \beta_0 + \beta_2 x_{2,i} + \varepsilon_i + (\beta_1 - b_1) x_{1,i}$$

• Given that $b_1 \neq \beta_1$, we get some stuff in the error term:

$$y_i - b_1 x_{1,i} = \beta_0 + \beta_2 x_{2,i} + \varepsilon_i + (\beta_1 - b_1)(u_{1,i} + u_{2,i} + \varepsilon_i),$$

and we now have an endogeneity problem because $x_{2,i} = u_{2,i} + u_{3,i}$ and $u_{2,i}$ is now in the error term.

• Note: we could formalize this intuition using the Frisch-Waugh Theorem.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- What if we include a variable in **X** that doesn't actually affect y?
- This can be accommodated in the linear regression framework as a regressor that has a coefficient $\beta_k = 0$.
- Thus, this does not create any bias in the other coefficients, and the expected coefficient on the irrelevant variable is zero.
- However, this can reduce the precision of the estimates of the other coefficients, and it's not generally a good idea to add as many variables to a regression as you can (**overfitting**).

< □ > < □ > < □ > < □ > < □ > < □ >

$$\mathbf{b}_{OLS} | \mathbf{X} \sim \mathcal{N} \left(\boldsymbol{eta}, \sigma^2 \left(\mathbf{X}' \mathbf{X}
ight)^{-1}
ight)$$

- $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ is the covariance matrix of the parameter estimates it tells us how precise \mathbf{b}_{OLS} is as an estimate of β .
- The diagonal elements of σ² (X'X)⁻¹ are the variances of each of the parameter estimates. The square roots of those are standard errors. Note that standard errors are standard deviations of the distribution of an estimator.
- The off-diagonal elements of $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ are also important (especially when we get to hypothesis testing). They indicate whether two parameter estimates are correlated.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Estimating Standard Errors

$$\mathbf{b}_{OLS} | \mathbf{X} \sim \mathcal{N} \left(\boldsymbol{\beta}, \sigma^2 \left(\mathbf{X}' \mathbf{X} \right)^{-1} \right)$$

- $(\mathbf{X}'\mathbf{X})^{-1}$ is just data, but σ^2 must be estimated (recall that it's the variance of the normally distributed error terms).
- Intuitively, the residuals are estimates of the error terms, so the sample variance of residuals is what we use to estimate σ²:

$$s^2 = rac{\mathbf{e}'\mathbf{e}}{n-K}$$

Similar to how the unbiased estimator of the variance of a random variable with unknown mean requires us to divide by n − 1, here we have to divide by n − K, where K is the number of parameters being estimated.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- Consider a parameter vector γ and consider $\mathbf{x}'\gamma$ as a predictor of y.
- The mean squared error of this predictor is

$$MSE = E\left[\left(y - \mathbf{x}' \boldsymbol{\gamma}\right)^2\right]$$

• It's important to notice that this can be written

$$MSE = E\left[\left(y - E\left[y|\mathbf{x}\right]\right)^{2}\right] + E\left[\left(E\left[y|\mathbf{x}\right] - \mathbf{x}'\gamma\right)^{2}\right]$$

 Result: OLS is the value of γ that minimizes mean square error. This does not require normality of the error terms. (Proof on board)

Gauss-Markov Theorem

In the linear regression model with homoscedastic errors, The OLS estimator is the **best linear unbiased estimator** (BLUE).

- "Best" means estimator with lowest variance most precise
- Normally distributed error not assumed here
- If error terms are not homoscedastic, OLS is still unbiased given strict exogeneity, but lower-variance estimators are possible (see: Generalized Least Squares)
- Proof on board

Finite Sample vs. Asymptotic Properties

- It's rare to have small-sample properties of an estimator
- Econometric studies typically do Monte Carlo (simulation) studies to learn about finite-sample performance of estimators.
- For most estimators, we derive asymptotic properties, i.e.,

$$\sqrt{n} \left(\mathbf{b}_{OLS} - \boldsymbol{\beta} \right) \Rightarrow_{d} \mathcal{N} \left(\mathbf{0}, \boldsymbol{\Sigma} \right)$$

- The above statement says that the OLS estimator **converges in distribution** to a normally distributed variable.
- Part of that result is that OLS is consistent. Formally, consistency means that an estimator converges in probability to the true value. Intuitively, this means that an estimator will be unbiased in very large samples (and also that the variance of the estimator becomes small in large samples).

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

Asymptotics: Setup

- Asymptotic analysis refers to analyzing statistics of the data as the number of observations *n* grows to infinity (e.g., the central limit theorem).
- To do asymptotic analysis, we specify a data generating process.
- A data generating process describes the distribution of a sequence of observations (x_i, ε_i) for i = 1, 2, ... n.
- The simplest case is assuming that observations are **independent and identically distributed** (i.i.d)
- In other settings, the observations are allowed to be correlated, but in a limited way that requires the dependence between "distant" observations to go to zero – see **stationarity** and **ergodicity** (e.g., time series, and more recently some spatial analysis)
- We also need the data to be "well-behaved", i.e. $E[\mathbf{xx}']$ must have full rank and the error terms must have finite variance.

• We can use the central limit theorem to show that OLS is asymptotically normally distributed:

$$\mathbf{b}_{OLS} \sim^{\mathbf{a}} \mathcal{N}\left(\boldsymbol{\beta}, \frac{\sigma^2}{n} E\left[\mathbf{x}_i \mathbf{x}_i'\right]^{-1}\right)$$

where \sim^a signifies convergence in distribution.

• Note that this is just like the finite-sample distribution we got with normally distributed error terms.

Mincerian Regression, Cornwell and Rupert (1988)

Residuals:

Min	1 Q	Median	3 <mark>Q</mark>	Max
-2.2034	-0.2379	-0.0071	0.2327	2.1380

Coefficients:

	Estimate	Std. Error	<mark>t</mark> value	Pr(> <mark>t</mark>)	
(Intercept)	5.245e+00	7.170e-02	73.153	< 2e-16	***
ED	5.654e-02	2.612e-03	21.644	< 2e-16	***
EXP	4.045e-02	2.174e-03	18.605	< 2e-16	***
EXP2	-6.811e-04	4.783e-05	-14.242	< 2e-16	***
WKS	4.485e-03	1.090e-03	4.115	3.94e-05	***
OCC	-1.405e-01	1.472e-02	-9.544	< 2e-16	***
SOUTH	-7.210e-02	1.249e-02	-5.773	8.37e-09	***
SMSA	1.390e-01	1.207e-02	11.513	< 2e-16	***
MS	6.736e-02	2.063e-02	3.265	0.0011	**
FEM	-3.892e-01	2.518e-02	-15.457	< 2e-16	***
UNION	9.015e-02	1.289e-02	6.993	3.13e-12	***

Residual standard error: 0.3524 on 4154 degrees of freedom Multiple R-squared: 0.4183, Adjusted R-squared: 0.4169F-statistic: 298.7 on 10 and 4154 DF, p-value: < 2.2e-16

イロト 不得下 イヨト イヨト 二日