# 1. The Experimental Ideal Hospitals' influence on health ...

 $\Rightarrow$  Do hospitals improve people's health?

 National Health Interview Survey about how healthy people feel (1 excellent health; 5 poor health) 2

Result suggests that hospitals make people sicker

| Group       | Sample Size | Mean health status | Std. Error |
|-------------|-------------|--------------------|------------|
| Hospital    | 7774        | 2.79               | 0.014      |
| No Hospital | 90049       | 2.07               | 0.003      |

 $\Rightarrow$  Explanation: People who go to the hospital are sicker and people who never went think they are healthier in the first place

# Describe problem more analytically

Dummy variable:

$$D_i = \begin{cases} 0 & \text{if not treated in hospital} \\ 1 & \text{if treated in hospital} \end{cases}$$

► Health status: *y<sub>i</sub>* 

Observed Outcome:

$$y_i = \begin{cases} y_{1i} & \text{if } D_i = 1\\ y_{0i} & \text{if } D_i = 0\\ = y_{0i} + (y_{1i} - y_{0i})D_i \end{cases}$$

### Potential outcomes

 $y_{i} = \begin{cases} y_{1i} & \text{if individual has been to the hospital} \\ y_{0i} & \text{if individual has not been to the hospital} \\ = y_{0i} + (y_{1i} - y_{0i})D_{i} \end{cases}$ (1)

► The causal effect of hospitalization for an individual is then:

$$y_{1i} - y_{0i}$$

Problem: We never see both potential outcomes for any one person!

# 1. The Experimental Ideal **Comparing conditional means**

► Difference in average health:

$$\mathbb{E}[y_i|D_i = 1] - \mathbb{E}[y_i|D_i = 0] = \mathbb{E}[y_{1i}|D_i = 1] - \mathbb{E}[y_{0i}|D_i = 1]$$

Observed difference in average health average treatment effect on the tr

$$+\underbrace{\mathbb{E}[y_{0i}|D_i=1]-\mathbb{E}[y_{0i}|D_i=0]}_{\mathbf{E}[y_{0i}|D_i=0]}$$

selection bias

 Average causal effect of hospitalization on those who were hospitalized:

$$\mathbb{E}[y_{1i}|D_i = 1] - \mathbb{E}[y_{i0}|D_i = 1] = \mathbb{E}[y_{1i} - y_{0i}|D_i],$$

where  $\mathbb{E}[y_{1i}|D_i = 1]$  is the health of the hospitalized, and  $\mathbb{E}[y_{0i}|D_i = 1]$  is the average health of those had they not been hospitalized (unobservable)

► Sick people are more likely to get treated; causes *selection bias* 

6

▶ <u>Problem</u>: Selection bias prevents us from studying *observable* quantitity  $\mathbb{E}[y_i|D_i = 1] - \mathbb{E}[y_i|D_i = 0]$ 

► Goal of empirical economic research: Overcome selection bias and say something about the causal effect of a variable

 Random assignment of D<sub>i</sub> solves the selection problem because it makes D<sub>i</sub> independent

Note that

$$\begin{split} \mathbb{E}[y_i|D_i = 1] - \mathbb{E}[y_i|D_i = 0] &= \mathbb{E}[y_{1i}|D_i = 1] - \mathbb{E}[y_{0i}|D_i = 0] \\ &= \mathbb{E}[y_{1i}|D_i = 1] - \mathbb{E}[y_{0i}|D_i = 1] \\ &= \mathbb{E}[y_{1i} - y_{0i}|D_i = 1] = \mathbb{E}[y_{1i} - y_{0i}] \end{split}$$

since under independence of  $y_{0i}$  and  $D_i$ , we have  $\mathbb{E}[y_{0i}|D_i = 0] = \mathbb{E}[y_{0i}|D_i = 1].$ 

 $\Rightarrow$  Selection bias is eliminated!

# Education Research: Are smaller classes better?

- Many non-experimental studies find no/little link between class size and children learning
- Typical error: just comparing test outcomes without accounting for selection bias
- ► Tennessee STAR experiment: Randomized trials; implemented in 1985/86; ran for 4 years; invovled > 11,000 children
- Key question: Did the randomization successfully balance subjects' characteristics across the different treatment groups?

9

#### Description

a cross-section from 1985-89

number of observations : 5748

observation : individuals

country : United States

#### Usage

data(Star)

#### Format

A dataframe containing :

tmathssk total math scaled score treadssk total reading scaled score classk type of class, a factor with levels (regular,small.class,regular.with.aide) totexpk years of total teaching experience sex a factor with levels (boy,girl) freelunk qualified for free lunch ? race a factor with levels (white,black,other) schidkn school indicator variable

#### 1. The Experimental Ideal —

# Does randomization work?

|                      |       | Class si | ze           |
|----------------------|-------|----------|--------------|
|                      | Small | Regular  | Regular/Aide |
| 1. Free lunch        | 0.473 | 0.475    | 0.499        |
| 2. White/Asian       | 0.683 | 0.677    | 0.661        |
| 3. Score percentiles | 0.744 | 0.733    | 0.733        |
| 4. Girls             | 0.485 | 0.486    | 0.487        |
|                      |       |          |              |

• Differences are small  $\Rightarrow$  Randomization seems to work!

# **Regression Analysis of Experiments**

► Assume that treatment effect is the same for everyone:

$$y_{1i} - y_{0i} = \rho$$

▶ Then, (1) becomes

$$y_i = \underbrace{\alpha}_{\mathbb{E}[y_{0i}]} + \underbrace{\rho}_{(y_{1i} - y_{0i})} D_i + \underbrace{\eta_i}_{y_{0i} - \mathbb{E}[y_{0i}]},$$

where  $\eta_i$  is the random part of  $y_{0i}$ .

- 1. The Experimental Ideal
  - Conditional expectations with treatment switched on and off are

$$\mathbb{E}[y_i|D_i = 1] = \alpha + \rho + \mathbb{E}[\eta_i|D_i = 1]$$
$$\mathbb{E}[y_i|D_i = 0] = \alpha + \mathbb{E}[\eta_i|D_i = 0]$$

#### so that

$$\mathbb{E}[y_i|D_1 = 1] - \mathbb{E}[y_i|D_i = 0] = \rho + \underbrace{\mathbb{E}[\eta_i|D_i = 1] - \mathbb{E}[\eta_i|D_i = 0]}_{\text{selection bias}},$$

where  $\rho$  is the treatment effect.

- Selection bias corresponds to correlation between error η<sub>i</sub> and regressor D<sub>i</sub>.
- Note that

 $\mathbb{E}[\eta_i | D_i = 1] - \mathbb{E}[\eta_i | D_i = 0] = \mathbb{E}[y_{0i} | D_i = 1] - \mathbb{E}[y_{0i} | D_i = 0].$ 

- Hence, correlation reflects the difference in no-treatment potential outcomes between those who get treated and those who don't.
- Under random assignment, the selection bias is not present. Regressing y<sub>i</sub> on D<sub>i</sub> gives the causal effect of interest ρ.

#### 1. The Experimental Ideal —

|                | Depende  | nt Variable: Scor | e percentile |
|----------------|----------|-------------------|--------------|
|                | (1)      | (2)               | (3)          |
| small          | 0.011*** | 0.011***          | 0.011***     |
|                | (0.002)  | (0.002)           | (0.002)      |
| black          |          | -0.024***         | -0.009***    |
|                |          | (0.002)           | (0.002)      |
| girl           |          | 0.012***          | 0.011***     |
| -              |          | (0.002)           | (0.001)      |
| freelunch      |          |                   | -0.027***    |
|                |          |                   | (0.002)      |
| totexpk        |          |                   | 0.001***     |
| ·              |          |                   | (0.0001)     |
| Constant       | 0.733*** | 0.735***          | 0.735***     |
|                | (0.001)  | (0.001)           | (0.002)      |
| Observations   | 5,748    | 5,748             | 5,748        |
| R <sup>2</sup> | 0.007    | 0.051             | 0.102        |

### 1. The Experimental Ideal -Controls

Covariates play two roles here:

 Assignment to classes of different sizes was random within schools, but not across schools (different school types may have different average class sizes)

 $\Rightarrow$  Need to include school-specific fixed effects

Control for student characteristics

► Controls  $x_i$  uncorrelated with treatment  $D_i \Rightarrow$  no affect on  $\rho$ , but may generate more precise estimates of causal effect

Yields the model

$$y_i = \alpha + \rho D_i + \mathbf{x}'_i \gamma + \eta_i$$

# 2. Regression Fundamentals | 2.1 Conditional Expectations - Outline I

#### 1. The Experimental Ideal

#### 2. Regression Fundamentals 2.1 Conditional Expectations

- 2.2 Ordinary Least Squares
- 2.3 Asymptotic OLS Inference
- 2.4 Hypothesis Testing

#### 3. Applying Regression

- 3.1 Marginal Effects
- 3.2 Dummies & Saturated Models
- 3.3 Causality and CIA
- 3.4 Difference-in-Differences
- 3.5 Clustered Standard Errors

- 2. Regression Fundamentals | 2.1 Conditional Expectations ------
  - ► We define the (n × 1) vector of dependent y and the (n × k) matrix of independent variables X as follows:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1},$$
$$\mathbf{X} = \begin{pmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_n \end{pmatrix}_{n \times k} = \begin{pmatrix} x_{11} & x_{12} & \vdots & x_{1k} \\ x_{21} & x_{22} & \vdots & x_{2k} \\ \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \vdots & x_{nk} \end{pmatrix}_{n \times k}.$$

### 2. Regression Fundamentals | 2.1 Conditional Expectations — Conditional Expectation Function

• Definition of the CEF for continuous  $y_i$  and conditional density  $f_{y_i}(\cdot | \mathbf{x}_i = x)$ 

$$\mathbb{E}[y_i|\boldsymbol{x}_i=x] = \int tf_{y_i}(t|\boldsymbol{x}_i=x)dt$$

► Definition of the CEF for discrete  $y_i$  and conditional density  $f_{y_i}(\cdot | \mathbf{x}_i = x)$ 

$$\mathbb{E}[y_i|\mathbf{x}_i=x] = \sum_t tf_{y_i}(t|\mathbf{x}_i=x)$$

CEF is random because x<sub>i</sub> is random and the CEF is a function of x<sub>i</sub>!

# 2. Regression Fundamentals | 2.1 Conditional Expectations — CEF-Decomposition Property

An implication of the law of iterated expectations is the possibility to split up a random variable in two pieces:

$$y_i = \mathbb{E}[y_i | \mathbf{x}_i] + \varepsilon_i$$

 $\triangleright$  (*i*)  $\varepsilon_i$  is "mean-independent" of  $\mathbf{x}_i$ , i.e.  $\mathbb{E}[\varepsilon_i | \mathbf{x}_i] = 0$ 

 $\triangleright$  (*ii*)  $\varepsilon_i$  is uncorrelated with any function of  $x_i$ 

▶ y<sub>i</sub> can be decomposed into a piece that is perfectly correlated with x<sub>i</sub> (the CEF) and a piece, which is orthogonal to any function of x<sub>i</sub> CEF-Prediction Property: Let m(x<sub>i</sub>) be any function of x<sub>i</sub>.
 The CEF solves

$$\mathbb{E}[y_i | \mathbf{x}_i] = \underset{m(\mathbf{x}_i)}{\operatorname{argmin}} \mathbb{E}[(y_i - m(\mathbf{x}_i))^2],$$

so it is the minimum mean squared error (MMSE) predictor of  $y_i$  given  $x_i$ .

ANOVA Theorem: The unconditional variance of y<sub>i</sub> can be split up as follows:

$$\mathbb{V}[y_i] = \mathbb{V}(\mathbb{E}[y_i|\boldsymbol{x}_i]) + \mathbb{E}[\mathbb{V}(y_i|\boldsymbol{x}_i)].$$

# 2. Regression Fundamentals | 2.1 Conditional Expectations – **CEF and linear regression**

- Population regression problem: Least squares problem formulated in terms of non-random features of the joint distribution of dependent and independent variables.
- A population regression coefficient is the solution to a population least squares problem.
- Let the  $(k \times 1)$  regression coefficient vector  $\beta$  be defined as  $\beta = \underset{\tilde{\beta}}{\operatorname{argmin}} \mathbb{E}[(y_i - \mathbf{x}'_i \tilde{\beta})^2]$
- First order condition:

$$\mathbb{E}[oldsymbol{x}_i(y_i - oldsymbol{x}_i' ilde{oldsymbol{eta}})] = 0$$

yielding

$$\boldsymbol{\beta} = \mathbb{E}[\boldsymbol{x}_i \boldsymbol{x}_i']^{-1} \mathbb{E}[\boldsymbol{x}_i' \boldsymbol{y}_i]. \tag{2}$$

2. Regression Fundamentals | 2.1 Conditional Expectations

#### ► By construction:

$$\mathbb{E}[\mathbf{x}_i(y_i - \mathbf{x}'_i\beta)] = \mathbf{0},$$

where

$$y_i - \mathbf{x}'_i \boldsymbol{\beta} = \varepsilon_i.$$

defines the population residual, which is uncorrelated with the regressors  $x_i$ .

## 2. Regression Fundamentals | 2.1 Conditional Expectations — Anatomy of multivariate regression

In a multivariate regression, i.e., with more than one (non-constant) regressor, the slope coefficient for the k-th regressor is given by coefficient

$$\beta_k = \frac{\mathbb{C}\mathsf{ov}\left[y_i, \tilde{x}_{ki}\right]}{\mathbb{V}[\tilde{x}_{ki}]}$$

where  $\tilde{x}_{ki}$  is the residual from a population regression of  $x_{ki}$  on all other covariates.

▶ Hence,  $\mathbb{E}[\mathbf{x}_i \mathbf{x}'_i]^{-1} \mathbb{E}[\mathbf{x}'_i y_i]$  is the  $(k \times 1)$  vector with k-th element given by  $\mathbb{C}$ ov  $[y_i, \tilde{x}_{ki}]/\mathbb{V}[\tilde{x}_{ki}]$ .

# 2. Regression Fundamentals | 2.1 Conditional Expectations — Justification for regression

**Theorem** (Linear CEF). Suppose the CEF is linear. Then, the population regression function is it, i.e.,  $\mathbb{E}[y_i|\mathbf{x}_i] = \mathbf{x}'_i \boldsymbol{\beta}$  with  $\boldsymbol{\beta}$  given by (2).

24

But when is the CEF linear?

- ► If the vector (y, x'<sub>i</sub>) is multivariate normally distributed (problem: variables are often discrete)
- If the regression models are saturated, i.e., there is a separate paramter for every possible combination of values the regressors can take

**Theorem** (Best Linear Predictor). The function  $\mathbf{x}'_i \boldsymbol{\beta}$  is the best linear predictor of  $\mathbf{y}$  given  $\mathbf{x}_i$  in a MSE sense. Hence,  $\mathbb{E}[\mathbf{y}_i | \mathbf{x}_i]$  is the best MMSE predictor of  $y_i$  given  $\mathbf{x}_i$  in the class of all functions of  $\mathbf{x}_i$ .

**Theorem** (Regression-CEF Theorem). The function  $\mathbf{x}'_i \boldsymbol{\beta}$  provides the minimal MSE linear approximation to  $\mathbb{E}[y_i | \mathbf{x}_i]$ , that is

$$eta = \mathop{\mathrm{argmin}}_{ ilde{eta}} \mathbb{E}[(\mathbb{E}[y_i | oldsymbol{x}_i] - oldsymbol{x}_i' ilde{eta})^2]$$

- While the CEF is the best unrestricted predictor (in a MSE sense) of the dependent variable, the regression provides the best <u>linear</u> predictor for it.
- ► If we want to approximate E[y<sub>i</sub>|x<sub>i</sub>] (as opposed to predicting y<sub>i</sub>), regression provides the best (in a MSE sense) <u>linear</u> approximation to it.

# 2. Regression Fundamentals | 2.1 Conditional Expectations — How to estimate $\beta$ ?

▶ The population regression coefficient

$$\boldsymbol{\beta} = \mathbb{E}[\boldsymbol{x}_i \boldsymbol{x}_i']^{-1} \mathbb{E}[\boldsymbol{x}_i' y_i].$$

is estimated by replacing the population moments  $\mathbb{E}[.]$  by their corresponding sample means:

$$\hat{\boldsymbol{\beta}} = \hat{\mathbb{E}}[\boldsymbol{x}_i \boldsymbol{x}_i']^{-1} \hat{\mathbb{E}}[\boldsymbol{x}_i' y_i] \\ = \left(\frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i \boldsymbol{x}_i'\right)^{-1} \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i' y_i \\ = (\boldsymbol{X}' \boldsymbol{X})^{-1} \boldsymbol{X}' \boldsymbol{y}$$

▶ Ordinary Least Squares (OLS) estimator



Duncan (carData)

R Documentation

28

#### **Duncan's Occupational Prestige Data**

Description

The Duncan data frame has 45 rows and 4 columns. Data on the prestige and other characteristics of 45 U. S. occupations in 1950.

Usage

Duncan

Format

This data frame contains the following columns:

type

Type of occupation. A factor with the following levels: prof, professional and managerial; we, white-collar; be, blue-collar.

#### income

Percentage of occupational incumbents in the 1950 US Census who earned \$3,500 or more per year (about \$36,000 in 2017 US dollars).

education

Percentage of occupational incumbents in 1950 who were high school graduates (which, were we cynical, we would say is roughly equivalent to a PhD in 2017)

prestige

Percentage of respondents in a social survey who rated the occupation as "good" or better in prestige

#### Source

Duncan, O. D. (1961) A socioeconomic index for all occupations. In Reiss, A. J., Jr. (Ed.) Occupations and Social Status. Free Press [Table VI-1].

#### References

Fox, J. (2016) Applied Regression Analysis and Generalized Linear Models, Third Edition. Sage.

Fox, J. and Weisberg, S. (2019) An R Companion to Applied Regression, Third Edition, Sage.

[Package carData version 3.0-3 Index]

#### 2. Regression Fundamentals | 2.1 Conditional Expectations



| Statistic | N  | Mean   | St. Dev. | Min     | Pctl(25) | Pctl(75) | Max    |
|-----------|----|--------|----------|---------|----------|----------|--------|
| income    | 45 | 41.867 | 24.435   | 7       | 21       | 64       | 81     |
| education | 45 | 52.556 | 29.761   | 7       | 26       | 84       | 100    |
| prestige  | 45 | 47.689 | 31.510   | 3       | 16       | 81       | 97     |
| resmod2   | 45 | -0.000 | 20.513   | -59.088 | -12.557  | 13.795   | 49.858 |

#### 2. Regression Fundamentals | 2.1 Conditional Expectations



|                         | Prestige                     |
|-------------------------|------------------------------|
|                         | prestige                     |
| education               | 0.546***                     |
|                         | (0.098)                      |
| income                  | 0.599***                     |
|                         | (0.120)                      |
| Constant                | -6.065                       |
|                         | (4.272)                      |
| Observations            | 45                           |
| R <sup>2</sup>          | 0.828                        |
| Adjusted R <sup>2</sup> | 0.820                        |
| Residual Std. Error     | 13.369 (df = 42)             |
| F Statistic             | $101.216^{***}$ (df = 2; 42) |
| Note:                   | *p<0.1; **p<0.05; ***p<0.01  |

## 2. Regression Fundamentals | 2.1 Conditional Expectations — **How to estimate** $\beta_k$ **?**

▶ The population regression coefficient

$$\beta_k = \frac{\mathbb{C}\mathsf{ov}\left[y_i, \tilde{x}_{ki}\right]}{\mathbb{V}[\tilde{x}_{ki}]}$$

is estimated by

$$\hat{\beta}_k = \frac{\widehat{\mathbb{Cov}}[y_i, \tilde{x}_{ki}]}{\widehat{\mathbb{V}}[\tilde{x}_{ki}]} = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(\tilde{x}_{ki} - \bar{x}_k)}{\frac{1}{n} \sum_{i=1}^n (\tilde{x}_{ki} - \bar{x}_k)^2},$$

where  $\tilde{x}_{ki}$  is the residual from a regression of  $x_{ki}$  on all other covariates.

# R Partialling out income

|   | prestige                         |
|---|----------------------------------|
|   | prestige                         |
| residuals from regressing education on income | 0.546 <sup>**</sup><br>(0.219)   |
| Constant                                      | 47.689 <sup>***</sup><br>(4.441) |
| Observations                                  | 45                               |
| R <sup>2</sup>                                | 0.126                            |
| Adjusted R <sup>2</sup>                       | 0.106                            |
| Residual Std. Error                           | 29.794 (df = 43)                 |
| F Statistic                                   | $6.214^{**}$ (df = 1; 43)        |
| Note:   | *p<0.1; **p<0.05; ***p<0.01      |

# 





| Wage  | es (Ecdat)   | R Documentation |
|-------|--|-----------------|
| Pa    | nel Data of Individual Wages                         |                 |
| Des   | cription   |                 |
| a pan | el of 595 observations from 1976 to 1982             |                 |
| numb  | per of observations : 4165                           |                 |
| obsei | rvation : individuals                                |                 |
| count | try : United States                                  |                 |
| Usa   | ge   |                 |
| data  | (Nages)  |                 |
| Forr  | nat  |                 |
| A dat | aframe containing :                                  |                 |
| exp   |  |                 |
|       | years of full-time work experience                   |                 |
| wks   |  |                 |
|       | weeks worked   |                 |
| bluec | ol   |                 |
|       | blue collar ?  |                 |
| ind   |  |                 |
|       | works in a manufacturing industry ?                  |                 |
| south |  |                 |
|       | resides in the south ?                               |                 |
| smsa  |  |                 |
|       | resides in a standard metropolitan statistical are ? |                 |
| marri | ed   |                 |
|       | married ?  |                 |

#### 2. Regression Fundamentals | 2.1 Conditional Expectations



| Statistic    | N     | Mean   | St. Dev. | Min   | Pctl(25) | Pctl(75) | Max   |
|--------------|-------|--------|----------|-------|----------|----------|-------|
| experience   | 4,165 | 19.854 | 10.966   | 1     | 11       | 29       | 51    |
| weeks worked | 4,165 | 46.812 | 5.129    | 5     | 46       | 50       | 52    |
| industry     | 4,165 | 0.395  | 0.489    | 0     | 0        | 1        | 1     |
| education    | 4,165 | 12.845 | 2.788    | 4     | 12       | 16       | 17    |
| log wage     | 4,165 | 6.676  | 0.462    | 4.605 | 6.395    | 6.953    | 8.537 |



|                         | Dependent variable:            |
|-------------------------|--------------------------------|
|                         | logarithmic Wage               |
| years of education      | 0.065***                       |
|                         | (0.002)                        |
| Constant                | 5.839***                       |
|                         | (0.031)                        |
| Observations            | 4,165                          |
| R <sup>2</sup>          | 0.155                          |
| Adjusted R <sup>2</sup> | 0.155                          |
| Residual Std. Error     | 0.424 (df = 4163)              |
| F Statistic             | $764.526^{***}$ (df = 1; 4163) |
| Note:                   | *p<0.1; **p<0.05; ***p<0.01    |

### Regression Fundamentals | 2.1 Conditional Expectations Regression fit vs. conditional means



- 2. Regression Fundamentals | 2.1 Conditional Expectations -
  - Iterating expectations in the formula for  $\beta$  yields

$$\beta = \mathbb{E}[\mathbf{x}_i \mathbf{x}_i']^{-1} \mathbb{E}[\mathbf{x}_i y_i] = \mathbb{E}[\mathbf{x}_i \mathbf{x}_i']^{-1} \mathbb{E}[\mathbf{x}_i \mathbb{E}[y_i | \mathbf{x}_i]].$$

- ► Hence, an implication of the regression-CEF theorem is that regression coefficients can be obtained by using E[y<sub>i</sub>]x<sub>i</sub>] as dependent variable instead of y<sub>i</sub>.
- Note that

$$\mathbb{E}[(\mathbb{E}[y_i|\boldsymbol{x}_i] - \boldsymbol{x}_i'\tilde{\boldsymbol{\beta}})^2] = \sum_{u} (\mathbb{E}[y_i|\boldsymbol{x}_i = \boldsymbol{u}] - \boldsymbol{u}'\tilde{\boldsymbol{\beta}})^2 g_{\boldsymbol{x}}(\boldsymbol{u}),$$

where  $g_x(u)$  is the probability mass function.

 $\Rightarrow \beta$  can be constructed from a weighted least squares regression of  $\mathbb{E}[y_i|\mathbf{x}_i = \mathbf{u}]$  on  $\mathbf{u}$ , where  $\mathbf{u}$  runs over the values of  $\mathbf{x}_i$  and each observation is weighted based on the distribution of  $\mathbf{x}_i$ .

|                         | Dependent variable:         |
|-------------------------|-----------------------------|
|                         | Means by years of education |
| years of education      | 0.057***                    |
| -                       | (0.006)                     |
| Constant                | 5.953***                    |
|                         | (0.071)                     |
| Observations            | 14                          |
| R <sup>2</sup>          | 0.872                       |
| Adjusted R <sup>2</sup> | 0.862                       |
| Residual Std. Error     | 0.095 (df = 12)             |
| F Statistic             | $81.912^{***}$ (df = 1; 12) |
| Note:                   | *p<0.1; **p<0.05; ***p<0.01 |

|                         | Dependent variable:         |
|-------------------------|-----------------------------|
|                         | Means by years of education |
| years of education      | 0.065***                    |
|                         | (0.004)                     |
| Constant                | 5.839***                    |
|                         | (0.056)                     |
| Observations            | 14                          |
| R <sup>2</sup>          | 0.951                       |
| Adjusted R <sup>2</sup> | 0.947                       |
| Residual Std. Error     | 0.012 (df = 12)             |
| F Statistic             | $235^{***}$ (df = 1; 12)    |
| Note:                   | *p<0.1; **p<0.05; ***p<0.01 |

## 2. Regression Fundamentals | 2.2 Ordinary Least Squares -Outline I

#### 1. The Experimental Ideal

#### 2. Regression Fundamentals

2.1 Conditional Expectations

#### 2.2 Ordinary Least Squares

2.3 Asymptotic OLS Inference2.4 Hypothesis Testing

#### 3. Applying Regression

- 3.1 Marginal Effects
- 3.2 Dummies & Saturated Models
- 3.3 Causality and CIA
- 3.4 Difference-in-Differences
- 3.5 Clustered Standard Errors

# 2. Regression Fundamentals | 2.2 Ordinary Least Squares – Assumptions

(A1) The model is linear in the parameters and with the conditional mean specified by  $\mathbb{E}[\boldsymbol{y}|\boldsymbol{X}] = \boldsymbol{X}\boldsymbol{\beta}$ .

- (A2) The regressor matrix  $\boldsymbol{X}$  is of full column rank k, i.e.,  $rk(\boldsymbol{X}) = k < n$  (with probability one).
- (A3) The regressors  $x_i^{(j)}$ , j = 1, ..., k are *strictly exogenous* with  $\mathbb{E}[\varepsilon | \mathbf{X}] = \mathbf{0}$  implying  $\mathbb{E}[\varepsilon_i | \mathbf{x}_1, ..., \mathbf{x}_n] = \mathbf{0}$  for all i = 1, ..., n.

#### 2. Regression Fundamentals | 2.2 Ordinary Least Squares -

Define the conditional covariance matrix of  $\varepsilon$  as

$$\mathbb{V}[arepsilon|oldsymbol{X}] = \mathbb{E}[arepsilonarepsilon'|oldsymbol{X}] =: oldsymbol{\Psi}.$$

(A4) The error terms are assumed to be uncorrelated, i.e.,  $\mathbb{E}[\varepsilon_i \varepsilon_j | \mathbf{X}] = 0$  for  $i \neq j$ . Then,  $\Psi = \mathbf{D}$  with

$$\boldsymbol{D} := diag[\sigma_1^2, \ldots, \sigma_n^2].$$

(A5) The error terms are assumed to be (conditionally) homoscedastic, i.e., they have equal variances across the sample, i.e.,  $\mathbb{V}[\varepsilon_i|\mathbf{X}] = \mathbb{E}[\varepsilon_i^2|\mathbf{X}] = \sigma^2$  for all i = 1, ..., n. Then,  $\mathbf{D} = \sigma^2 \mathbf{I}_n$ , where  $\mathbf{I}_n$  denoting an *n*-dimensional identity matrix.

### 2. Regression Fundamentals | 2.2 Ordinary Least Squares – Hat Matrix and Residual Maker

▶ The OLS residual vector is given by

$$oldsymbol{e} = egin{pmatrix} e_1 \ e_2 \ dots \ e_n \end{pmatrix} = oldsymbol{y} - \hat{oldsymbol{y}} = oldsymbol{y} - oldsymbol{X} \hat{eta}.$$

► Vector of fitted dependent variables:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{eta}$$
  
=  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$   
=  $P\mathbf{y}$ ,

where  $\boldsymbol{P} := \boldsymbol{X} (\boldsymbol{X}' \boldsymbol{X})^{-1} \boldsymbol{X}'$  is a projection matrix.

2. Regression Fundamentals | 2.2 Ordinary Least Squares -----

Therefore:

$$\boldsymbol{e} = \boldsymbol{y} - \boldsymbol{P} \boldsymbol{y} = (\boldsymbol{I} - \boldsymbol{P}) \boldsymbol{y} = \boldsymbol{M} \boldsymbol{y},$$

where

$$\boldsymbol{M} := \boldsymbol{I} - \boldsymbol{P} = \boldsymbol{I} - \boldsymbol{X}(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'$$

is called "residual maker" or "annihilator".

- For the matrix M it can be shown that MX = 0, M' = M, and MM = M
- ▶ Then, the relationship between e and  $\varepsilon$  is

$$\boldsymbol{e} = \boldsymbol{M}(\boldsymbol{X}\boldsymbol{eta} + \boldsymbol{arepsilon}) = \boldsymbol{M}\boldsymbol{arepsilon},$$

2. Regression Fundamentals | 2.2 Ordinary Least Squares – Numerical Properties of OLS

In terms of the residual vector, the F.O.C. is given by

$$oldsymbol{X}'oldsymbol{e} = \sum_{i=1}^n oldsymbol{x}_i e_i = oldsymbol{0}.$$

▷ That is, each column of **X** is orthogonal to e. ⇒  $\hat{y}'e = 0$ 

▶ If the regression equation has an intercept:

$$\sum_{i=1}^{n} x_{i1} \cdot e_i = \sum_{i=1}^{n} e_i = 0.$$

#### 2. Regression Fundamentals | 2.2 Ordinary Least Squares -Frisch-Waugh Theorem

Assume the regression model can be written as

$$oldsymbol{y} = oldsymbol{X}_1oldsymbol{eta}_1 + oldsymbol{X}_2oldsymbol{eta}_2 + arepsilon$$

and the OLS estimator is given by  $\hat{\boldsymbol{\beta}} = (\hat{\beta}_1', \hat{\beta}_2')' = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}.$ 

► Theorem (Frisch-Waugh) Define  $M_1 = I_N - X_1(X'_1X_1)^{-1}X'_1, \quad M_2 = I_N - X_2(X'_2X_2)^{-1}X'_2$ as the residual makers based on regressions on  $X_1$  and  $X_2$ . Then,  $\hat{\beta}_1$  ( $\hat{\beta}_2$ ) is obtained as the OLS estimator in a regression of  $M_1y$  ( $M_2y$ ) on  $M_1X_2$  ( $M_2X_1$ ). Hence,  $\hat{\beta}_1 = (X'_1M_2X_1)^{-1}(X'_1M_2y)$  $\hat{\beta}_2 = (X'_2M_1X_2)^{-1}(X'_2M_1y).$  ► The Frisch-Waugh theorem proves the anatomy of multivariate regression discussed and illustrated in Section 2.1.

### 2. Regression Fundamentals | 2.2 Ordinary Least Squares -Mean and Covariance of OLS

▶ **Proposition** (Unbiasedness): Under (A1)-(A3) , the OLS estimator  $\hat{\beta}$  is unbiased, i.e.

$$\mathbb{E}[\hat{\boldsymbol{\beta}}|\boldsymbol{X}] = \boldsymbol{\beta} \qquad (\forall \boldsymbol{\beta}).$$

 ▶ Proposition (Variance of the OLSE): Under assumptions (A1)
 - (A4), the conditional covariance matrix of the OLS estimator *β* is given by

$$\boldsymbol{V}_{\hat{\boldsymbol{\beta}}} := \mathbb{V}[\hat{\boldsymbol{\beta}} | \boldsymbol{X}] = (\boldsymbol{X}' \boldsymbol{X})^{-1} \boldsymbol{X}' \boldsymbol{D} \boldsymbol{X} (\boldsymbol{X}' \boldsymbol{X})^{-1}.$$
(3)

If we additionally impose (A5), i.e., homoscedasticity, (3) collapses to

$$\boldsymbol{V}_{\hat{\beta}} = \sigma^2 (\boldsymbol{X}' \boldsymbol{X})^{-1}. \tag{4}$$

2. Regression Fundamentals | 2.2 Ordinary Least Squares **Unbiased Estimation of**  $\sigma^2$  and  $\mathbb{V}[\hat{\beta}]$ 

• Proposition (Mean of SSE=  $\sum_{i=1}^{n} e_i^2$ ):  $\mathbb{E}[SSE|\mathbf{X}] = \mathbb{E}[\mathbf{e}'\mathbf{e}|\mathbf{X}] = (n-k)\sigma^2.$ 

$$\hat{\sigma}^2 = \frac{SSE}{n-k} = \frac{e'e}{n-k}$$

50

is an unbiased estimator of  $\sigma^2$ , i.e.  $\mathbb{E}[\frac{SSE}{n-k}|\mathbf{X}] = \sigma^2$ .

### 2. Regression Fundamentals | 2.2 Ordinary Least Squares -Estimating the OLS VC Matrix

► Under homoscedasticity, we have

$$\hat{\mathbb{V}}[\hat{\boldsymbol{\beta}}|\boldsymbol{X}] = \hat{\sigma}^2(\boldsymbol{X}'\boldsymbol{X})^{-1},$$

which is an unbiased estimator of  $\mathbb{V}[\hat{\beta}|\mathbf{X}]$ , i.e.  $\mathbb{E}[\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}|\mathbf{X}] = \mathbb{V}[\hat{\beta}|\mathbf{X}]$ .

▶ Under heteroscedasticity,  $V_{\hat{\beta}}$  can be estimated by

$$\hat{oldsymbol{V}}_{\hat{eta}}^{EW} = (oldsymbol{X}'oldsymbol{X})^{-1} \left(\sum_{i=1}^n oldsymbol{x}_i oldsymbol{x}_i' e_i^2
ight) (oldsymbol{X}'oldsymbol{X})^{-1}.$$

Developed by Eicker (1963) and introduced in econometrics by White (1980) and is refered to as *Eicker-White robust* or *heteroscedasticity-consistent* covariance estimator

# 2. Regression Fundamentals | 2.2 Ordinary Least Squares -**Computation of** $\hat{\mathbb{V}}[\hat{\beta}|X]$

```
# Residuals standard error hat(sigma)
sqrt(1/(n-k) * as.numeric(t(e)%*%e)) # by hand
summary (Im (lwage \sim ed + wks + exp, data = Wages)) $ sigma # by package
# VC matrix under homoscedasticity
# ____
VCV = (1/(n-k) * as.numeric(t(e) %*% e) * solve(t(X) %*% X))
vc < -vcov(lm(lwage ~ ed + wks + exp, data = Wages))
xtable :: xtable (vc, auto=TRUE)
# Standard errors of the estimated coefficients
(stder = sqrt(diag(VCV)))
# VC matrix under heteroskedasticity
model <- Im(lwage ~ wageedu + wks + exp, data = Wages)
vc < -vcovHC (model, type = 'HC0')
xtable::xtable(VCV. auto=TRUE)
```



#### Homoscedasticity:

| -           | (Intercept) | ed        | wks        | exp        |
|-------------|-------------|-----------|------------|------------|
| (Intercept) | 0.0044591   | -0.000074 | -0.0000692 | -0.0000115 |
| ed          | -0.0000740  | 0.000005  | 0.0000000  | 0.000003   |
| wks         | -0.0000692  | 0.000000  | 0.0000015  | 0.0000000  |
| exp         | -0.0000115  | 0.000000  | 0.0000000  | 0.0000003  |

Heteroscedasticity:

|             | (Intercept) | ed        | wks        | exp        |
|-------------|-------------|-----------|------------|------------|
| (Intercept) | 0.0044591   | -0.000074 | -0.0000692 | -0.0000115 |
| ed          | -0.0000740  | 0.000005  | 0.0000000  | 0.000003   |
| wks         | -0.0000692  | 0.000000  | 0.0000015  | 0.0000000  |
| exp         | -0.0000115  | 0.000000  | 0.0000000  | 0.0000003  |

2. Regression Fundamentals | 2.2 Ordinary Least Squares – Breusch/Pagan/Godfrey Test

• Linear model:  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 

► Heteroscedastic error terms:

$$\mathbb{V}[\varepsilon_i] := \sigma_i^2 = h(\gamma_0 + \mathbf{z}_i' \boldsymbol{\gamma}),$$

54

 $\begin{array}{l} \triangleright \ h(\cdot) > 0 \ \text{denotes an unknown function,} \\ \triangleright \ \textbf{z}_i \ \text{denotes a} \ (p \times 1) \ \text{vector of (known!) variance regressors.} \end{array}$ 

▶ Null hypothesis:  $H_0$  :  $\gamma = 0$ .

# 

(i) Estimate the model under the null hypothesis.

(ii) Compute the centered  $R^2$  of the auxiliary regression:

$$\mathbf{e}_i^2/\hat{\sigma}^2 = \delta_0 + \delta_1 z_{i1} + \ldots + \delta_p z_{ip} + \xi_i,$$

where  $e_i$  is the OLS residual of step (i).

(iii) Compute  $LM_{BPG} = nR^2$ .

# 

- ▶ Null hypothesis:  $\sigma_i^2 = \sigma^2$ .
- ► Idea: Comparing CV matrices  $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\Omega\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$  and  $(\mathbf{X}'\mathbf{X})^{-1}$ .
- (i) Estimate the model under the null hypothesis.
- (ii) Estimate the following auxiliary regression by OLS:

1

$$\mathbf{e}_i^2 = \delta_0 + \boldsymbol{\delta}' \boldsymbol{z}_i,$$

where  $z_i$  consists of the regressors, their squares and all cross terms.

(iii) Then, the LM statistic is computed based on the  $R^2$  of the auxiliary regression

$$n \cdot R^2 \stackrel{a}{\sim}_{H_0} \chi^2_{(q)},$$

where the degrees of freedom q are the number of auxiliary regressors in (ii) excluding the intercept.



```
> bptest(model, data = Wages)
studentized Breusch-Pagan test
data: model
BP = 28.791, df = 3, p-value = 2.478e-06
> # White test
> bptest(model, ~wageedu*wks+wageedu*exp+wks*exp+1(wageedu)^2+1(exp)^2
+1(wks)^2, data = Wages)
studentized Breusch-Pagan test
data: model
BP = 57.347, df = 6, p-value = 1.554e-10
```

# 2. Regression Fundamentals | 2.2 Ordinary Least Squares ------Efficiency of OLS

Theorem (Gauss-Markov): Under the full ideal conditions (i) - (v), the least squares estimator β̂ is the best linear unbiased estimator (BLUE) of β, i.e.

$$\mathbb{V}[\tilde{\boldsymbol{\beta}}|\boldsymbol{X}] - \mathbb{V}[\hat{\boldsymbol{\beta}}_{OLS}|\boldsymbol{X}] = \sigma^2 \boldsymbol{D}' \boldsymbol{D} \quad \text{pos. semi-def.},$$

where  $\tilde{\boldsymbol{\beta}}$  is any linear unbiased estimator of  $\boldsymbol{\beta}$  and  $\boldsymbol{D}$  is given by

$$ilde{oldsymbol{eta}} - \hat{oldsymbol{eta}}_{OLS} = oldsymbol{D}'oldsymbol{y}.$$

Proposition (The LSE under Normality): Under the full ideal conditions (i) – (vi) the distribution of β̂ is given by:

$$\hat{\boldsymbol{eta}} | \boldsymbol{X} \sim \mathcal{N}(\boldsymbol{eta}, \sigma^2(\boldsymbol{X}'\boldsymbol{X})^{-1}).$$

► **Theorem**: For the normally distributed linear regression model, the distribution of the statistic  $e'e/\sigma^2$  is

$$rac{oldsymbol{e}'oldsymbol{e}}{\sigma^2}|oldsymbol{X}\sim\chi^2_{(n-k)}.$$

2. Regression Fundamentals | 2.2 Ordinary Least Squares -

Proposition (The variance of 
<sup>
<sup>ˆ</sup>σ<sup>2</sup></sup>): In a normally distributed regression model, the variance of 
<sup><sup>ˆ</sup>σ<sup>2</sup></sup> is given by:

$$\mathbb{V}[\hat{\sigma}^2|\boldsymbol{X}] = \frac{2\sigma^4}{n-k}.$$

▶ **Proposition** (Independence of  $\hat{\beta}$  and  $\hat{\sigma}^2$ ): In the normally distributed linear regression model,  $\hat{\beta}$  and  $\hat{\sigma}^2$  are independent random variables.

#### 

Coefficient of determination:

$$R^2 := \frac{\mathsf{SSR}}{\mathsf{SST}} = \frac{\hat{\mathbb{V}}[\hat{y}_i]}{\hat{\mathbb{V}}[y_i]}.$$

▶ If the model contains an intercept term, then

$$\sum_{i} (y_i - \bar{y})^2 = \sum_{i} (\hat{y}_i - \bar{y})^2 + \sum_{i} e_i^2$$
  
SST = SSR + SSE

► Then:

$$R^2 := 1 - \frac{\mathsf{SSE}}{\mathsf{SST}} = 1 - \frac{\hat{\mathrm{V}}[e_i]}{\hat{\mathrm{V}}[y_i]}$$

## 2. Regression Fundamentals | 2.2 Ordinary Least Squares — **Properties of** $R^2$

- ▶  $R^2 = \rho_{y,\hat{y}}^2$ , where  $\rho$  denotes the empirical correlation.
- If the model has an intercept,  $0 \le R^2 \le 1$ .
- The  $R^2$  follows an unknown distribution.
- Adding further variables leads to an increase in the  $R^2$ .
- ► An *R*<sup>2</sup> can have a reasonable size in spurious regressions if the regressors are non-stationary.
- ► Linear transformations of the regression model do not change the value of the *R*<sup>2</sup> coefficient.

### 2. Regression Fundamentals | 2.2 Ordinary Least Squares — **Alternative** $R^2$ **Concepts**

▶ Non-centered  $R^2$ :

$$\tilde{R}^2 = \frac{\sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2} = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n y_i^2}.$$

• General definition of an  $R^2$ :

$$R^2 = \mathbb{C}\mathrm{or}[y_i, \hat{y}_i]^2.$$

• Adjusted  $R^2$ :

$$\bar{R}^2 = 1 - rac{(n-1)\sum_{i=1}^n e_i^2}{(n-k)\sum_{i=1}(y_i - \bar{y})^2},$$

where k denotes the number of regressors in the model.

# 2. Regression Fundamentals | 2.2 Ordinary Least Squares — Multicollinearity

**Lemma**(Variance of  $\hat{\beta}_j$ ) Under assumptions (A1) to (A5),  $\mathbb{V}[\hat{\beta}_j | \mathbf{X}]$  can be written as

$$\mathbb{V}[\hat{\beta}_j|\boldsymbol{X}] = \frac{\sigma^2}{S_{x^{(j)}x^{(j)}}} \frac{1}{1 - R_j^2},$$

where  $S_{x^{(j)}x^{(j)}} := \sum_{i=1}^{n} (x_i^{(j)} - \bar{x}^{(j)})^2$  denotes the 'variation' of regressor  $x^{(j)}$  and  $R_j^2$  is the  $R^2$  of an (auxiliary) regression of  $x^{(j)}$  on all  $x^{(i)}$ ,  $(i \neq j)$ .

## 2. Regression Fundamentals | 2.2 Ordinary Least Squares – Indicators for Multicollinearity

Variance Inflation Factor:

$$\forall IF_j = \frac{1}{1 - R_j^2}$$

▶ Rule of thumb:  $VIF_j > 10$  implies serious multicollinearity

- Low values of det(X'X) indicate multicollinearity.
- ▶ If **X** if in correlation form:  $0 \le det(\mathbf{X}'\mathbf{X}) \le 1$ .