



### 6.3 – OLS in multiple regression

The population model:

$k \rightarrow$  counting # of  $X$ 's  $K < n$

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \epsilon_i \quad (6.1)$$

How to estimate the  $\beta$ s?

- Still want to minimize the sum of squared residuals (the sum of "vertical distances"):

$$\min_{b_0, b_1, \dots, b_k} \sum_{i=1}^n e_i^2 \quad \leftarrow \text{"prediction errors"}$$

$k \rightarrow$  # of slopes  
 $1 \rightarrow$  intercept

$$\frac{\partial \sum e_i^2}{\partial b_0} = 0 \quad \frac{\partial \sum e_i^2}{\partial b_1} = 0 \quad \frac{\partial \sum e_i^2}{\partial b_2} = 0$$

$\Rightarrow$  solve 3 simultaneous equations

before:

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

- Take  $(k + 1)$  derivatives, set them equal to zero, solve
- The new formula is too difficult to show (unless we use matrices, which we won't)

grad econometrics:  
 $b = (X'X)^{-1}X'Y$  ← don't study!

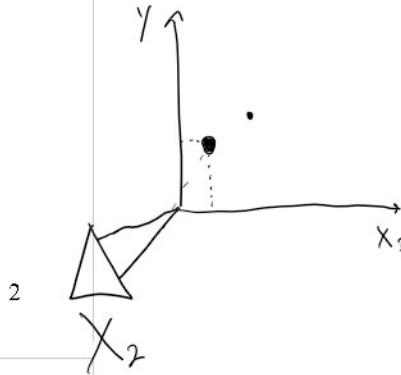
$b_0 =$   
 $b_1 = ?$  changes  
 $b_2 =$

The resulting estimated model:

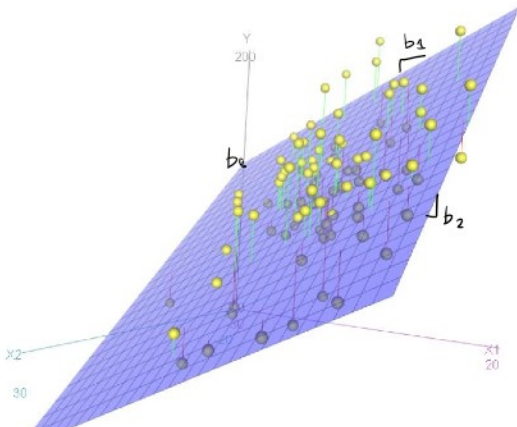
$$\hat{Y}_i = b_0 + b_1 X_{1i} + b_2 X_{2i} + \dots + b_k X_{ki} \quad (6.2)$$

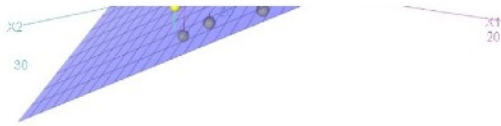
can't be interpreted as a line! (It's a  $k$ -dimensional hyperplane).

We can still try to visualize things if we have only 2  $X$  variables, however:

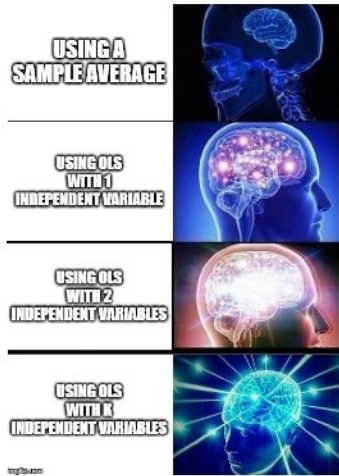


$H_0: \beta_2 = 0$   
 $t = \frac{b_2 - 0}{s.e.(b_2)}$





3



4

### 6.3.2 Interpretation

Let's look at a population model with two  $X$  variables:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \epsilon_i \quad (6.3)$$

- $Y$  is still the dependent variable
  - $X_1$  and  $X_2$  are the <sup>explanatory</sup> independent variables (the regressors)
  - $i$  still denotes an observation number
  - $\beta_0$  is the population intercept
  - $\beta_1$  is the effect of  $X_1$  on  $Y$ , holding all else constant ( $X_2$ )
  - $\beta_2$  is the effect of  $X_2$  on  $Y$ , holding all else constant ( $X_1$ )
  - $\epsilon$  is the regression error term (containing all the omitted factors that affect  $Y$ )
- X<sub>2</sub> has come out of  $\epsilon$*
- "ceteris paribus"*

5

No perfect correlation (-1 or 1)

## 6.4: A2 – No perfect multicollinearity

Now that we have multiple  $X$  variables in our model, we need to make an **additional assumption** in order for OLS to work:

**There is no perfect multicollinearity.** This means:

- No two variables (or combinations of variables) are exactly linearly related
- No two variables are perfectly correlated

6

For example, exact linear relationships between  $X$ s are:

- $X_1 = X_2$
- $X_1 = 100X_2$
- $X_1 = 1 + X_2 - 3X_3$

If you know  $X_1$ , you know  $X_2$  in first two examples).

Including both variables would be redundant.

OLS can't handle it. (Like dividing by zero).

Some common examples of where the assumption of "no perfect multicollinearity" is violated in practice are when the same variable is measure in different units (such as square feet and square metres, or dollars and cents), and in the *dummy variable trap*.

7

The **Living.Area** variable measures the size of the house in square feet. Suppose that there was another variable in the data set that measured **house size in square metres** (1 square foot = 0.0929 square metre). We can create this variable in R using:

```
House.Size <- 0.0929 * Living.Area
```

and now let's include it in our OLS estimation:

```
summary(lm(Price ~ Fireplaces + Living.Area + House.Size))
```

```
Coefficients: (1 not defined because of singularities)
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  14.730146   5.007563   2.942  0.00331 **
Fireplaces    8.962440   3.389656   2.644  0.00827 **
Living.Area   0.109313   0.003041  35.951 < 2e-16 ***
House.Size    NA         NA         NA     NA
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 68980 on 1725 degrees of freedom
Multiple R-squared:  0.5095,    Adjusted R-squared:  0.5089
F-statistic: 895.9 on 2 and 1725 DF,  p-value: < 2.2e-16
```

8

### 6.4.1 The dummy variable trap

The dummy variable trap occurs when one too many dummy variables are included in the equation. For example, suppose that we have a dummy variable **female** that equals 1 if the worker is female. Suppose that we also have a variable **male** that equals 1 if the worker is male. There is an exact linear combination between the two variables:

$$\begin{array}{ccc} & \nearrow 0 & \nearrow 1 \\ \text{female} = 1 - \text{male} & & \\ & \searrow 1 & \searrow 0 \end{array}$$

OLS won't work for:

$$\text{wage} = \beta_0 + \beta_1 \times \text{male} + \beta_2 \times \text{female} + \epsilon$$

Much easier to fall into the trap for "categorical variables"

- Alberta = 1 if Location = AB; 0 otherwise
- British.Columbia = 1 if Location = BC; 0 otherwise
- Manitoba = 1 if Location = MB; 0 otherwise
- ⋮
- Yukon = 1 if Location = YT; 0 otherwise

Wage      Location → 13 dummy variables:  
 21.2      "AB"  
 14.1      "MB"  
 Alberta = 1 if Location = "AB"  
 Manitoba = 1 if " " "MB"

$$\text{wage} = \beta_0 + \beta_1 \text{Alberta} + \beta_2 \text{Manitoba} + \dots + \beta_{12}$$

- We would create 13 dummy variables using "location", but only include 12 of them in our equation
- The group that is left out becomes the "base group"
- We could also drop the intercept (but this isn't usually done)

$$\hat{\text{wage}} = b_0 + 2.11 \text{Alberta}$$

base group is ON:

$$\hat{\text{wage}} = b_0 + \textcircled{2.3} \text{Alberta} + \textcircled{0.2} \text{Manitoba}$$

Final note:

Non-linear transformations are ok! We will do this in chapter 8.

$$X_2 = X_1^2 \quad \text{OK} \checkmark$$

$\log(\lambda_2)$

NO perfect multicollinearity

→ very high correlation btw variables

### 6.4.2 Imperfect multicollinearity

Imperfect multicollinearity is when two (or more) variables are almost perfectly related (highly correlated).

#### Example

Pretend we know the pop. model:

$$Y = 2X_1 + 2X_2 + c$$

and that the correlation between  $X_1$  and  $X_2$  is 0.99.

```
summary(lm(Y ~ X1))
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -4.4165      3.8954  -1.134    0.263
X1           4.0762      0.4698   8.676 2.13e-11 ***
```

The estimated standard error is small, so that the  $t$ -statistic is large, and we are sure that  $X_1$  is statistically significant. However, the estimated  $\beta_1$  is twice as big as it should be. This is because of omitted variable bias.

```
summary(lm(Y ~ X1 + X2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-4.676	3.956	-1.182	0.243
X1	1.958	4.075	0.481	0.633
X2	2.128	4.066	0.523	0.603

Now, the estimated  $\beta$ s are closer to their true value of 2, but both appear to be statistically insignificant! (Note the large standard errors and small  $t$ -statistics.)

14

### The problem:

- Because  $X_1$  and  $X_2$  are correlated, difficult to attribute changes in  $X_1$  to changes in  $Y$  (same for  $X_2$ )
- $X_1$  and  $X_2$  are almost always changing together in a similar way
- *ceteris paribus* assumption is not feasible
- $\beta_1$  is the effect of  $X_1$  on  $Y$ , holding  $X_2$  constant

15

### How imperfect multicollinearity affects estimation

- large standard errors, wide confidence intervals
- adding and dropping variables results in large swings of the estimated values
- overall – makes us unsure about our results
- problem is difficult to address *increase  $n$  (can't do it usually)*
- can't drop variables (OVB)
- if you don't need to interpret the affected variables, it's not a problem

Imperfect Mult.

- What is it? (def<sup>n</sup>)
- The problem?

16