# POOLED REGRESSION

FELIPE BUCHBINDER

# A POOLED REGRESSION TREATS PANEL DATA *AS IF IT WERE* POOLED DATA

# SERIOUSLY, WHAT DOES THIS MEAN?

# POOLED DATA :
# A DIFFERENT SAMPLE OF ENTITIES AT EACH TIME PERIOD
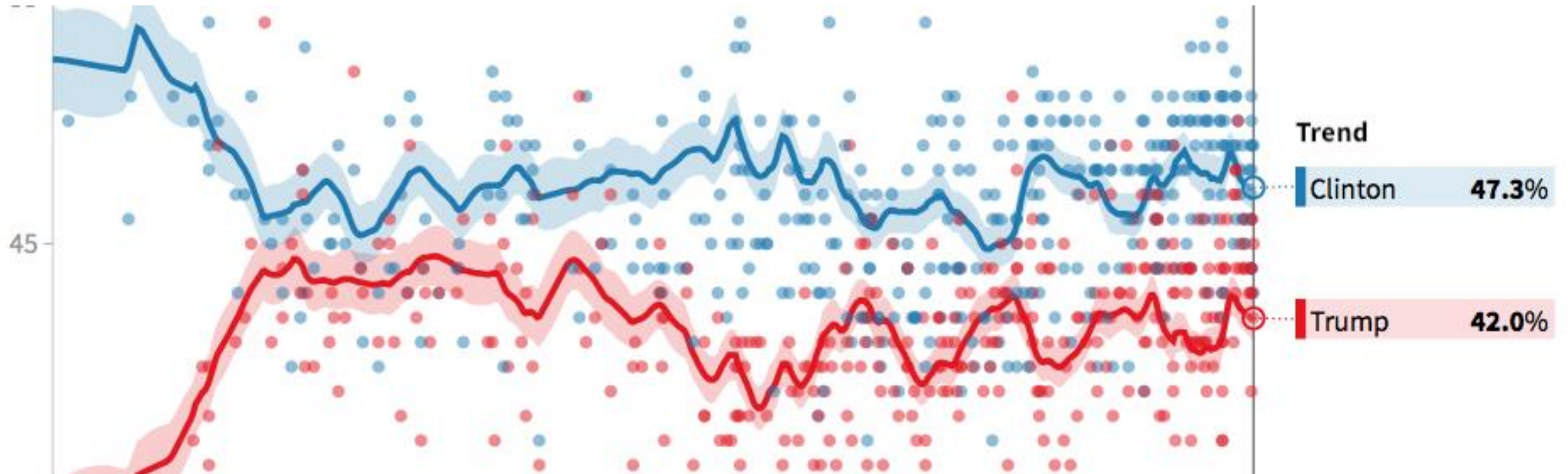
Sample 3

Sample 2

Sample 1

# POOLED DATA IS LIKE A TIME SERIES OF CROSS-SECTIONS

# ELECTION POLLS ARE AN EXAMPLE OF POOLED DATA

## A DIFFERENT GROUP OF PEOPLE IS SAMPLED AT EACH TIME

IN POOLED REGRESSION, JUST FORGET YOU HAVE PANEL DATA AND RUN AN ORDINARY LINEAR REGRESSION INSTEAD

THIS MEANS WE STACK ALL OBSERVATIONS AND TREAT THEM AS INDEPENDENT OBSERVATIONS

AWESOME! CAN I DO THAT?

# IGNORING A RELEVANT VARIABLE TYPICALLY LEADS TO BIASED ESTIMATES OF THE COEFFICIENTS…

$$\mathbf{b} = (\mathbf{X^T X})^{-1} \mathbf{X^T} y$$

$$= (\mathbf{X^T X})^{-1} \mathbf{X^T} (\mathbf{X\beta} + \mathbf{U} + \boldsymbol{\epsilon})$$

$$= \underbrace{(\mathbf{X^T X})^{-1} \mathbf{X^T X}}_{\mathbf{I}} \boldsymbol{\beta} + (\mathbf{X^T X})^{-1} \mathbf{X^T U} + (\mathbf{X^T X})^{-1} \mathbf{X^T} \boldsymbol{\epsilon}$$

$$= \boldsymbol{\beta} + (\mathbf{X^T X})^{-1} \mathbf{X^T U} + (\mathbf{X^T X})^{-1} \mathbf{X^T} \boldsymbol{\epsilon}$$

Taking the expected value…

$$\mathbb{E}(\mathbf{b}) = \boldsymbol{\beta} + (\mathbf{X^T X})^{-1} \mathbf{X^T U} + (\mathbf{X^T X})^{-1} \underbrace{\mathbb{E}(\mathbf{X^T} \boldsymbol{\epsilon})}_{0}$$

$$\therefore$$

$$\mathbb{E}(\mathbf{b}) = \boldsymbol{\beta} + (\mathbf{X^T X})^{-1} \mathbf{X^T U}$$

Thus, in general, $\mathbb{E}(\mathbf{b}) \neq \boldsymbol{\beta}$:

<span style="color:red">**b** is a biased estimate of **β**</span>

# IGNORING A RELEVANT VARIABLE TYPICALLY LEADS TO BIASED ESTIMATES OF THE COEFFICIENTS...

$$b = (X^TX)^{-1}X^Ty$$

$$= (X^TX)^{-1}X^T(X\beta + U + \epsilon)$$

$$= \underbrace{(X^TX)^{-1}X^TX}_{I}\beta + (X^TX)^{-1}X^TU + (X^TX)^{-1}X^T\epsilon$$

$$= \beta + (X^TX)^{-1}X^TU + (X^TX)^{-1}X^T\epsilon$$

Taking the expected value...

$$\mathbb{E}(b) = \beta + (X^TX)^{-}$$

...unless $X^TU = 0$,
In which case Pooled Regression yields unbiased estimates of the coefficients!

$$\mathbb{E}(b) = \beta + (X^TX)^{-1}X^TU$$

Thus, in general, $\mathbb{E}(b) \neq \beta$:

$b$ is a biased estimate of $\beta$

WHAT DOES $X^T U = 0$ MEAN IN PRACTICE?

IT MEANS THAT THE UNOBSERVED VARIABLES ($U$) ARE UNRELATED TO THE OBSERVED CHARACTERISTICS OF THE ENTITIES ($X$)

# EXERCISE

You will be given some examples of regression. For each regression…

- Think of possible sources of unobserved heterogeneities (U).

- What theoretical argument would you need to make in order to convince a reviewer that pooled data is a valid approach to tackle this regression problem? (i.e. what would $X^T U = 0$ mean in this case?)

# EXERCISE

1. Child Mortality ~ Democracy

2. Spending ~ Income

3. Wages ~ Education level

4. Crime Rate ~ Unemployment

$X^T U = 0$ MEANS COEFFICIENTS ARE UNBIASED, BUT THIS DOES NOT MEAN POOLED REGRESSION WORKS PERFECTLY ...

$$Y_{it} = X_{it}\beta + \underbrace{U_i}_{} + \epsilon_{it}$$

Ignored in pooled regression

$$Y_{it} = X_{it}\beta + \boxed{U_i + \epsilon_{it}}$$
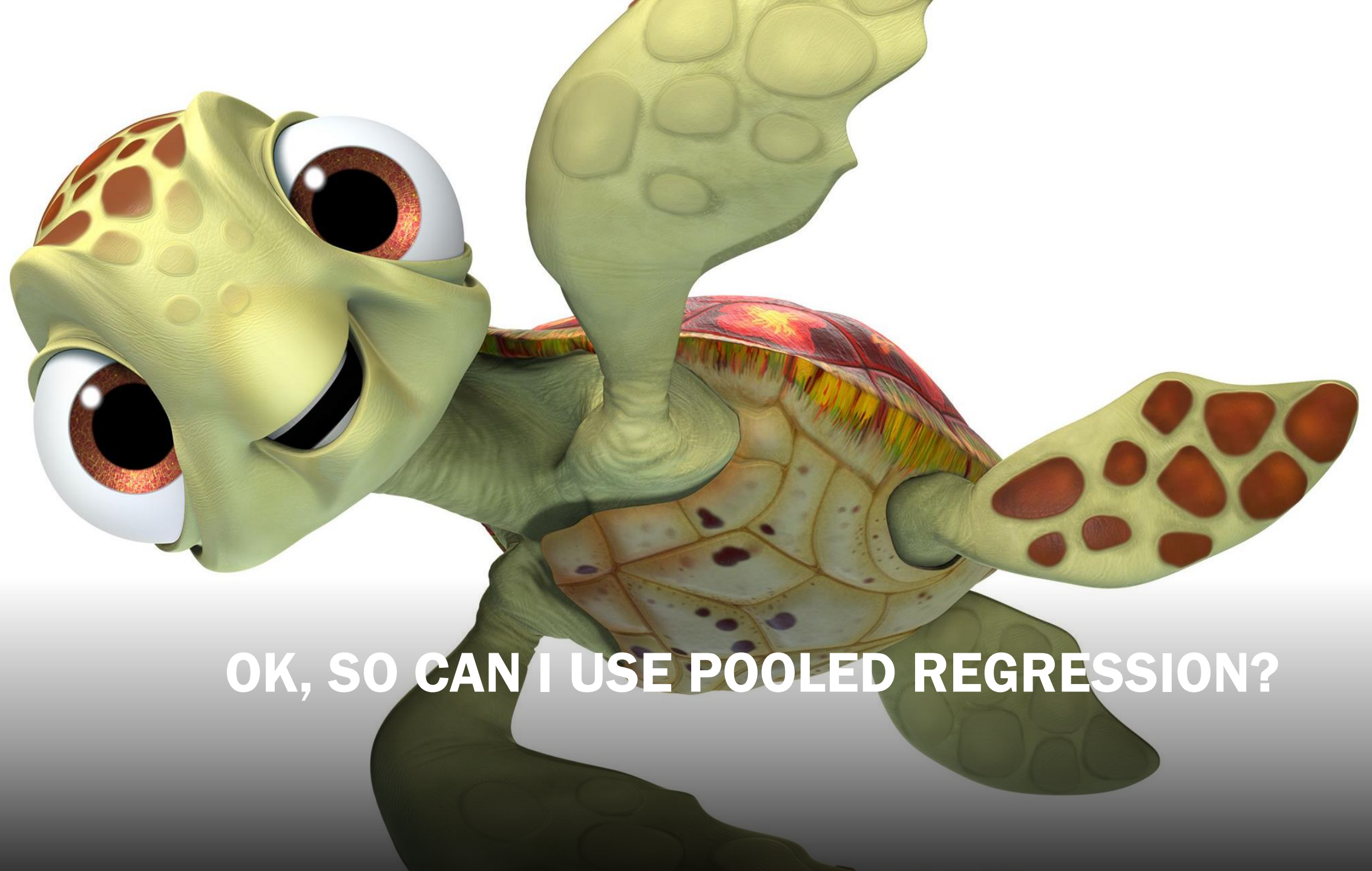
Pooled regression's residual

$$Y_{it} = X_{it}\beta + U_i + \epsilon_{it}$$

$$\epsilon_{it}^{\textbf{Pooled}} = U_i + \epsilon_{it}$$

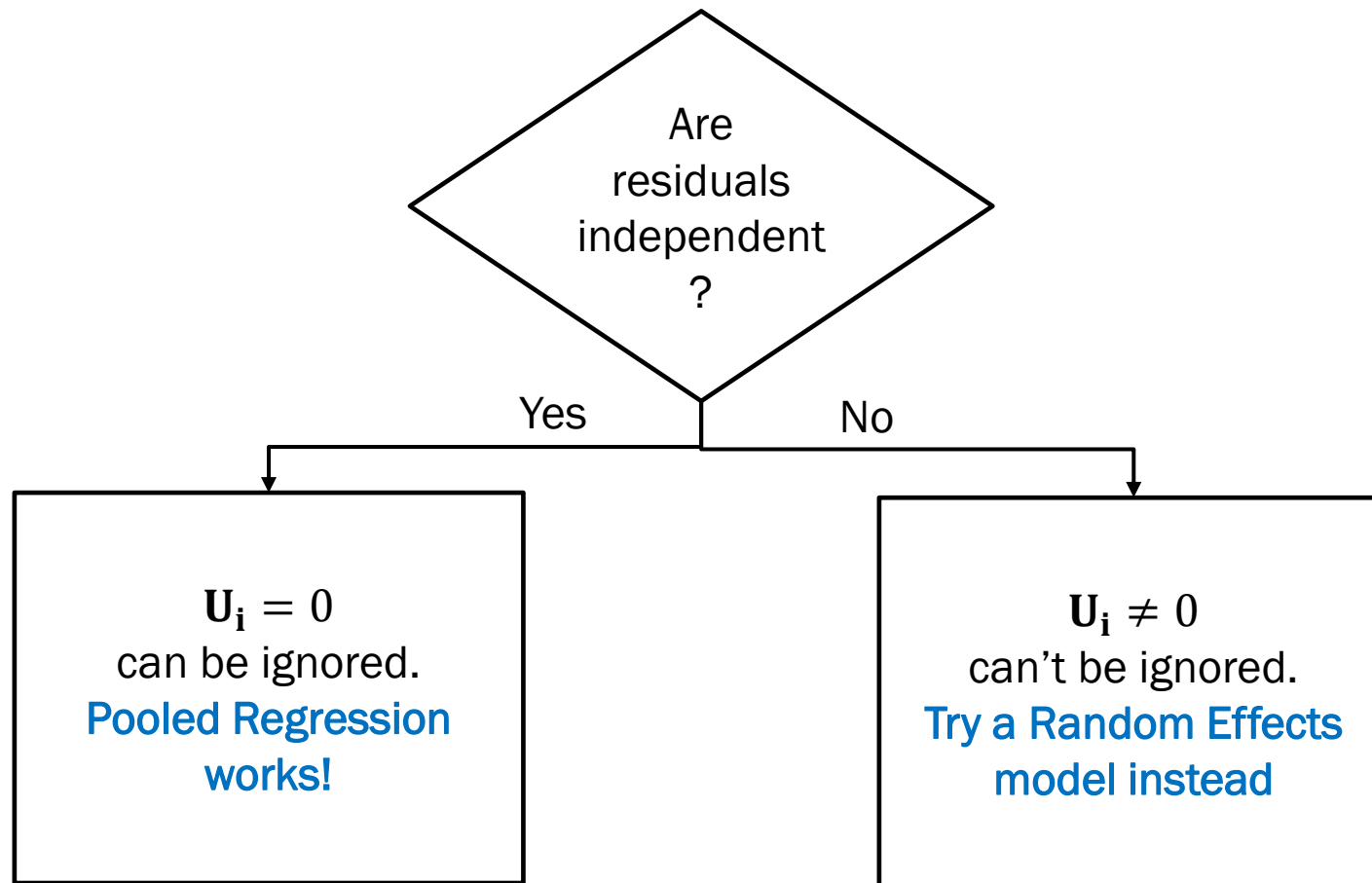$$Y_{it} = X_{it}\beta + U_i + \epsilon_{it}$$

$$\epsilon_{it}^{\text{Pooled}} = U_i + \epsilon_{it}$$
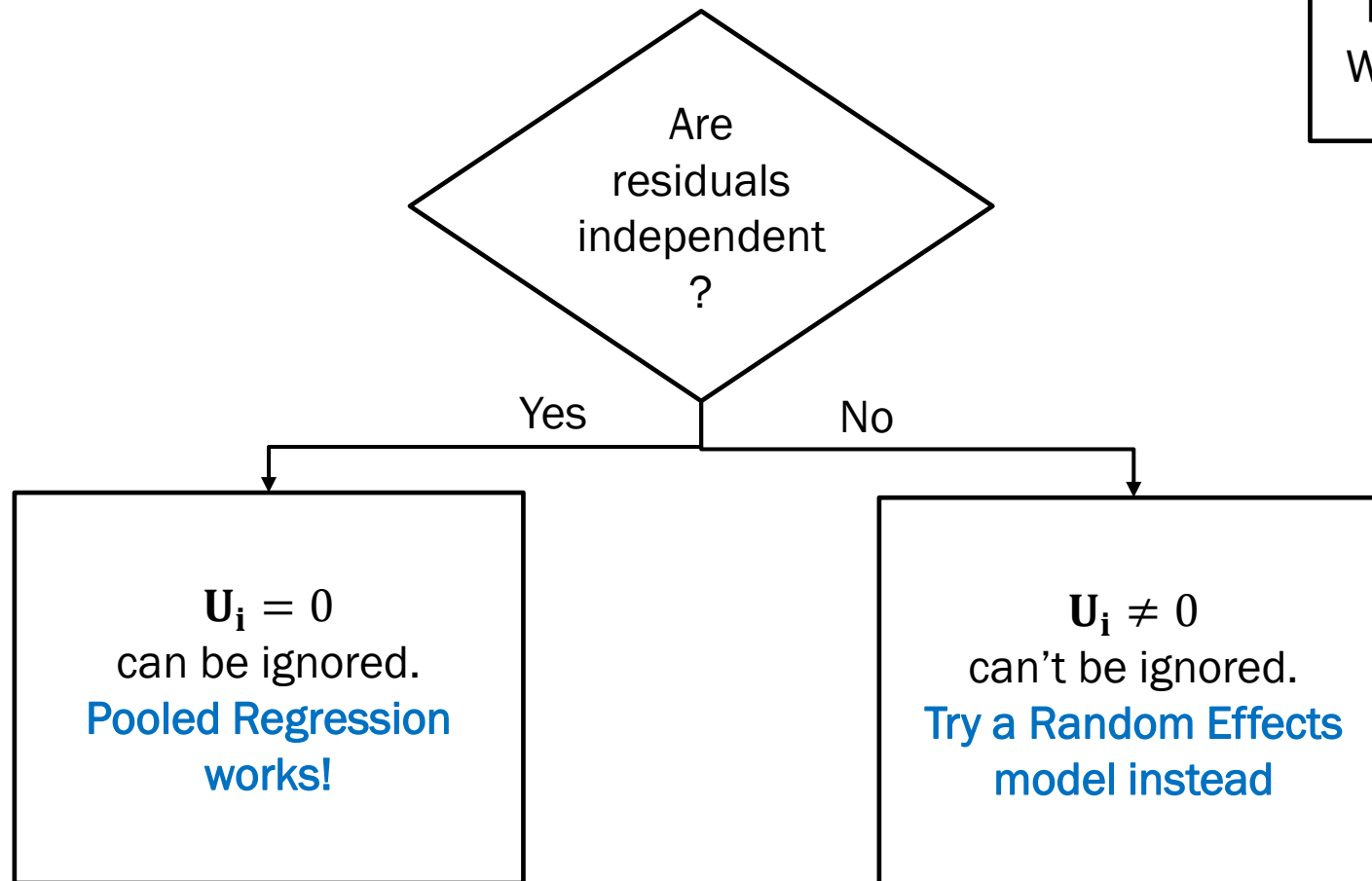
Residuals are serially correlated

OK, SO CAN I USE POOLED REGRESSION?
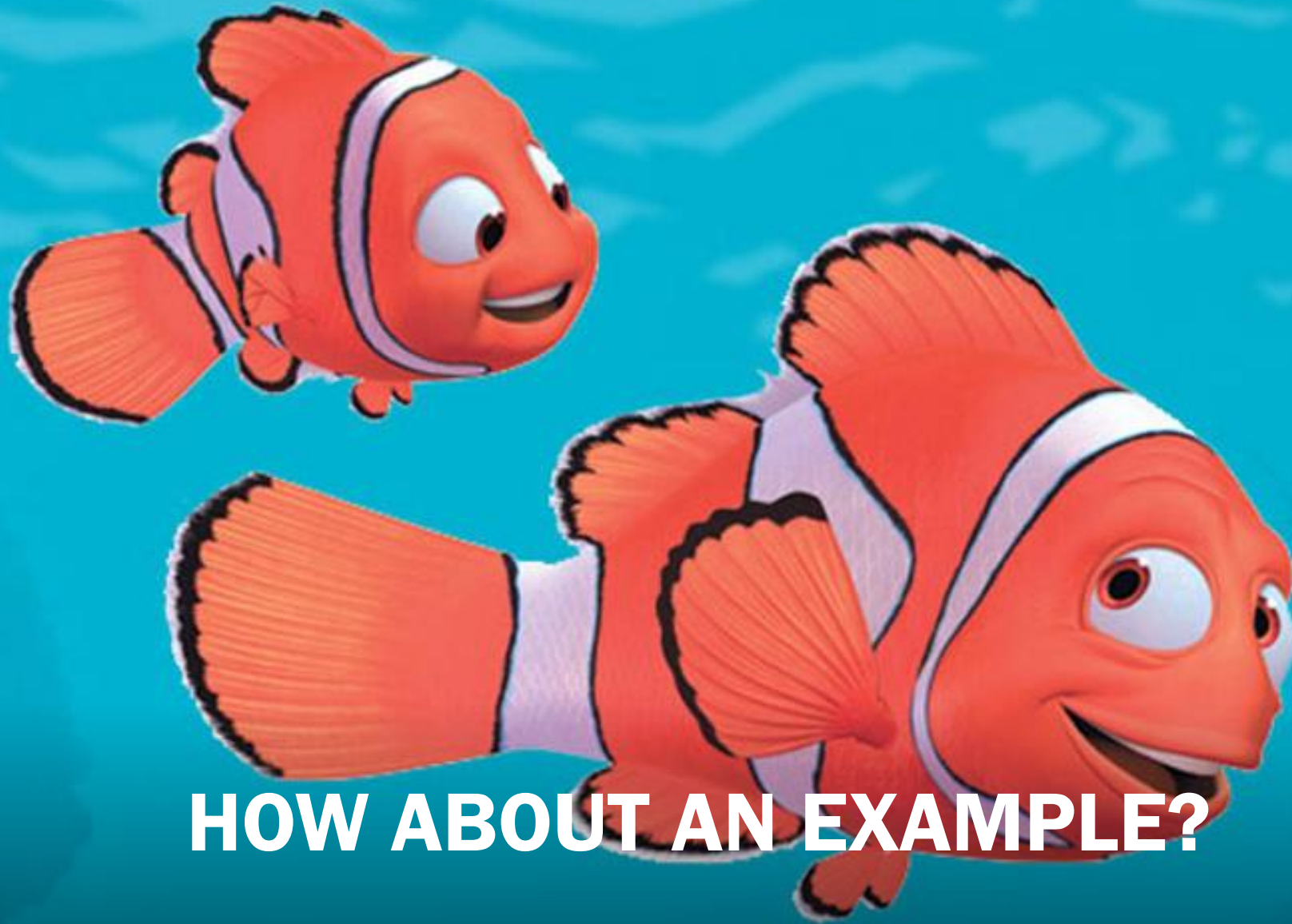
# BREUSCH-PAGAN LAGRANGE MULTIPLIERS TEST

# BREUSCH-PAGAN LAGRANGE MULTIPLIERS TEST

Recall that
$$\epsilon^{\mathbf{Pooled}} = \mathbf{U_i} + \epsilon$$
$$\mathbf{b} = \boldsymbol{\beta} + \left(\mathbf{X^T X}\right)^{-1}\mathbf{X^T U}$$
What happens if $\mathbf{U} = 0$?

Are residuals independent ?

Yes

No

$\mathbf{U_i} = 0$
can be ignored.
Pooled Regression works!

$\mathbf{U_i} \neq 0$
can't be ignored.
Try a Random Effects model instead

# BREUSCH-PAGAN LAGRANGE MULTIPLIERS TEST

$$\begin{cases} H_0: & \mathbf{U} = 0 \quad \text{Use Pooled Regression} \\ H_a: & \mathbf{U} \neq 0 \quad \text{Try Random Effects} \end{cases}$$

HOW ABOUT AN EXAMPLE?

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital, data=Grunfeld, model='pooling')

```
===============================================
                    Dependent variable:
                 ------------------------------
                              inv
                            Pooled
-----------------------------------------------
value                      0.116***
                           (0.006)

capital                    0.231***
                           (0.025)

Constant                  -42.714***
                           (9.512)

-----------------------------------------------
Observations                  200
R2                           0.812
Adjusted R2                  0.811
F Statistic        426.576*** (df = 2; 197)
===============================================
Note:            *p<0.1; **p<0.05; ***p<0.01
```

WHAT IF WE TRY TO CAPTURE $U_i$ BY USING A DUMMY FOR EACH $i$?

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

```
library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital,
data=Grunfeld, model='pooling')

ols <- lm(inv~value + capital +
factor(firm, data=Grunfeld)
```



|  | Dependent variable: | |
| --- | --- | --- |
|  | inv | |
|  | panel linear Pooled (1) | OLS OLS with Dummies (2) |
| value | 0.116*** (0.006) | 0.110*** (0.012) |
| capital | 0.231*** (0.025) | 0.310*** (0.017) |
| factor(firm)2 |  | 172.203*** (31.161) |
| factor(firm)3 |  | -165.275*** (31.776) |
| factor(firm)4 |  | 42.487 (43.910) |
| factor(firm)5 |  | -44.320 (50.492) |
| factor(firm)6 |  | 47.135 (46.811) |
| factor(firm)7 |  | 3.743 (50.565) |
| factor(firm)8 |  | 12.751 (44.053) |
| factor(firm)9 |  | -16.926 (48.453) |
| factor(firm)10 |  | 63.729 (50.330) |
| Constant | -42.714*** (9.512) | -70.297 (49.708) |
| Observations | 200 | 200 |
| R2 | 0.812 | 0.944 |
| Adjusted R2 | 0.811 | 0.941 |
| Residual Std. Error |  | 52.768 (df = 188) |
| F Statistic | 426.576*** (df = 2; 197) | 288.500*** (df = 11; 188) |
| Note: | | *p<0.1; **p<0.05; ***p<0.01 |

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

```
library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital,
data=Grunfeld, model='pooling')

ols <- lm(inv~value + capital +
factor(firm, data=Grunfeld)
```

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

```
library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital,
data=Grunfeld, model='pooling')

ols <- lm(inv~value + capital +
factor(firm, data=Grunfeld)
```

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital, data=Grunfeld, model='pooling')

ols <- lm(inv~value + capital + factor(firm, data=Grunfeld)

.

|  | Dependent variable: | |
|---|---|---|
|  | inv | |
|  | panel linear Pooled (1) | OLS OLS with Dummies (2) |
| value | 0.116*** (0.006) | 0.110*** (0.012) |
| capital | 0.231*** (0.025) | 0.310*** (0.017) |
| factor(firm)2 |  | 172.203*** (31.161) |
| factor(firm)3 |  | -165.275*** (31.776) |
| factor(firm)4 |  | 42.487 (43.910) |
| factor(firm)5 |  | -44.320 (50.492) |
| factor(firm)6 |  | 47.135 (46.811) |
| factor(firm)7 |  | 3.743 (50.565) |
| factor(firm)8 |  | 12.751 (44.053) |
| factor( |  | -16.926 (48.453) |
| factor( |  | 63.729 (50.330) |
| Constar |  | -70.297 (49.708) |
| Observations | 200 | 200 |
| R2 | 0.812 | 0.944 |
| Adjusted R2 | 0.811 | 0.941 |
| Residual Std. Error |  | 52.768 (df = 188) |
| F Statistic | 426.576*** (df = 2; 197) | 288.500*** (df = 11; 188) |

Note: *p<0.1; **p<0.05; ***p<0.01

Why did the $R^2$ increase?

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital, data=Grunfeld, model='pooling')

ols <- lm(inv~value + capital + factor(firm), data=Grunfeld)

# IS INVESTMENT DETERMINED BY COMPANY VALUE?
## THE GRUNFELD DATASET

library("plm")

library("stargazer")

pooled <- plm(inv ~ value + capital, data=Grunfeld, model='pooling')

ols <- lm(inv~value + capital + factor(firm, data=Grunfeld)

IT'S CALLED **FIXED EFFECTS**

STAY TUNED…