

ECON-C4100 - Capstone: Econometrics I

Lecture 7: Multiple regression #2: estimation

Otto Toivanen

- At the end of this lecture, you
 - 1 understand what how multivariate regression differs from univariate regression.
 - 2 understand how and why to carry out a multivariate regression analysis.
 - 3 appreciate the assumptions made in multivariate regression analysis.
 - 4 are aware of the most common pitfalls in regression analysis.

- ① How do the individual coefficients compare to univariate results?
- ② What explains the difference(s)?
- ③ What about statistical significance of individual coefficients?
- ④ What about several / all coefficients?
- ⑤ What about R^2 ?

- 6 What is the interpretation of individual coefficients?
- 7 (under what assumptions) does OLS work?
- 8 How to choose which explanatory variables to include / exclude?
- 9 What if the world is more complicated than linear?
- 10 What all can go wrong, and how would I know / find out?

Q6 What is the interpretation of individual coefficients?

- Our estimation equation is:

$$Income = \beta_0 + \beta_{AgeMV} Age + \beta_{GMV} G + u_{MV} \quad (1)$$

Interpretation of individual coefficients

- Regression yields the **conditional expectation** of the dependent variable Y :

$$\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad (2)$$

$$\mathbb{E}[Income_i|\mathbf{X} = \mathbf{x}] = \beta_0 + \beta_{AgeMV} Age_i + \beta_{GMV} G_i \quad (3)$$

- By plugging into the regression those values of \mathbf{X} that we are interested we get the conditional expectation of Y .

Interpretation of individual coefficients

- Example: the expected income of a woman of 35 years of age is:

$$\begin{aligned}\mathbb{E}[Income_i | \mathbf{X} = 35, 1] = \\ 12819.2 + 298.55 \times 35 - 4545.02 \times 1 = 18723.4\end{aligned}$$

- Example: the expected income of a man of 55 years of age is:

$$\begin{aligned}\mathbb{E}[Income_i | \mathbf{X} = 55, 0] = \\ 12819.2 + 298.55 \times 55 - 4545.02 \times 0 = 29239.45\end{aligned}$$

Interpretation of individual coefficients

- Coefficients as partial derivatives:

$$\frac{\partial \mathbb{E}[Income_i | \mathbf{X}]}{\partial X_k} = \beta_k \quad (4)$$

- With discrete explanatory variables cannot take derivatives, so a coefficient measures the change in Y from a one unit change in X_k :

$$\beta_k = \mathbb{E}[Income_i | \mathbf{X}, X_k = m] - \mathbb{E}[Income_i | \mathbf{X}, X_k = m - 1] \quad (5)$$

- Notice that in both, we fix all other variables (their effect on Y).
- So β_G is the effect of gender on income, **keeping the effect of Age constant**.

Q7 What are the OLS assumptions?

- 1 Strict exogeneity: $\mathbb{E}(u|\mathbf{X}) = 0$.
- 2 (\mathbf{X}_i, Y) , $i = 1, \dots, n$ are independent and identically distributed across observations.
- 3 \mathbf{X}_i and Y_i (u_i) have finite *fourth* moments.
- 4 No **perfect multicollinearity** (\mathbf{X} has **full column rank**).
- 5 Auxiliary: u_i is homoscedastic.

Perfect multicollinearity

- Analog: To solve a system of equations, you need as many equations as you have unknowns.
- Two variables are perfectly (multi)collinear if one is a perfect linear function of the other.
- Example: Think of a phenomenon with two **mutually exclusive and exhaustive** outcomes, A and B.
- A dummy taking value 1 if A is true and 0 otherwise: $D_A = 1 - D_B$, where D_B is the dummy taking value 1 if B is true and 0 otherwise ("**dummy variable trap**").
- Perfect collinearity = correlation $+/- 1$.

- Recall from previous lecture the two-variable model (App 6.2. in S&W).
- 2 explanatory variables and homosc. errors, $\rho_{X_1, X_2} \neq 0$. Then

$$\sigma_{\beta_1}^2 = \frac{1}{n} \frac{1}{1 - \rho_{X_1, X_2}^2} \frac{\sigma_u^2}{\sigma_{X_1}^2}$$

Collinearity of two variables

- Notice what happens when $\rho_{X_1, X_2} \rightarrow 1$.
- Collinearity (= "high" ρ_{X_1, X_2}) increases the standard error(s) of the other coefficient(s).
- Correlation between the \mathbf{X} s is a two-edged sword:
 - ① It removes omitted variable bias.
 - ② It reduces the efficiency gains from introducing a further explanatory variable.

Collinearity vs. multicollinearity

- Collinearity refers to the (high) correlation between two variables.
- **Multicollinearity** is a characteristic of a matrix (vector) \mathbf{X} .
- while the pair-wise correlations between elements of \mathbf{X} may be "not so high", the aggregate effect of them may lead to inflated standard errors.

Q8 How to choose the explanatory variables?

- Too few explanatory variables \rightarrow possible omitted variable bias.
- Too many variables may lead to multicollinearity and inflated se's.
- Note: "too many" requires correlation among explanatory variables.
- Can one test one's way out of this?
- No, but tests do help.

How to choose the explanatory variables?

- There are tests of individual and of joint significance. Why cannot I run these on autopilot?
- Case #1: start from a small model, add variables according to some (statistical) criterion.
- Case #2: start from a large model, drop variables according to some (statistical) criterion.
- Case #3: use machine learning methods (for later). Designed especially for the case where number of variables $>$ number of observations.

How to choose the explanatory variables?

- What goes wrong?
 - ① Statistical significance \neq economic significance.
 - ② Statistical significance \neq economic relevance.
 - ③ You may end up with variables that are highly correlated with Y, but have no real connection to it.
 - ④ Multiple testing leads to wrong (too good) test results.

How to choose the explanatory variables?

- The principled approach:
 - ① Before touching your data, write down a protocol.
 - ② Base explanatory (control) variables on theory and existing knowledge.
 - ③ Specify a testing protocol.
 - ④ Execute.

How to choose the explanatory variables?

- The practical approach:
 - ① Try to be as close to the principled approach as possible.
 - ② Learning allowed and encouraged → new/respecification.
 - ③ Robustness testing.

How to choose the explanatory variables?

- Robustness testing:
 - ① It is rarely the case that there is a "right model" that you can (re)cover.
 - ② Ask: are your results sensitive to small, well-justified changes to your model?
 - ① Adding (meaningful) variables.
 - ② Deleting variables.
 - ③ Changing functional form.
 - ④ Changing assumptions about the error term.

Q8 What if the world is not linear?

- Well, nothing prevents us from making our model nonlinear.
 - ① Keep the \mathbf{X} base-variables the same, but make the function $f(\mathbf{X})$ more complicated.
 - ② Transform the variables.
- Let's start by making $f(\mathbf{X})$ more complicated.
- Let's remind ourselves of what the income - age graph looks like.
- But before that let's remind ourselves of what polynomial functions are.

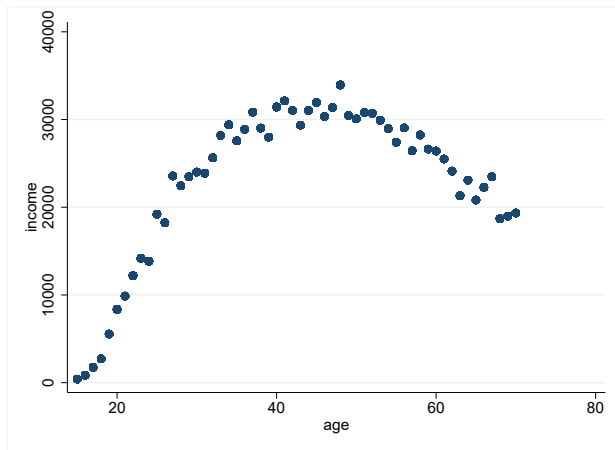
What if the world is not linear?

- Polynomial of order k (of one variable):

$$\begin{aligned} p(X) &= \sum_{i=0}^k \alpha_i X^i \\ &= \alpha_0 + \alpha_1 X + \alpha_2 X^2 + \dots \alpha_k X^k \end{aligned}$$

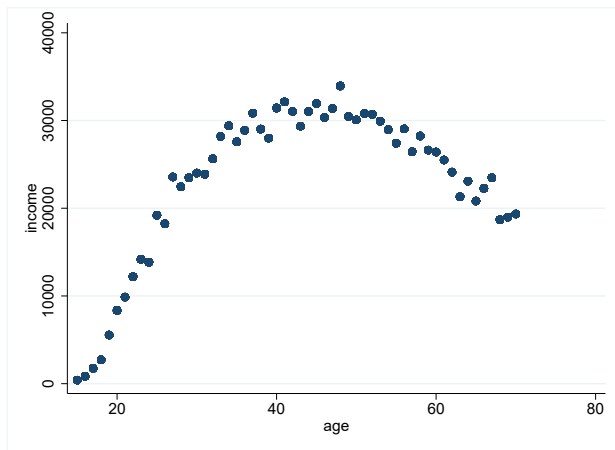
- Polynomials may consist of several variables.

What if the world is not linear?



- The figure suggests that an inverted - U shaped function could be a good fit.

What if the world is not linear?



- The figure suggests that an inverted - U shaped function could be a good fit.
- Let's try a quadratic function of age.

What if the world is not linear?

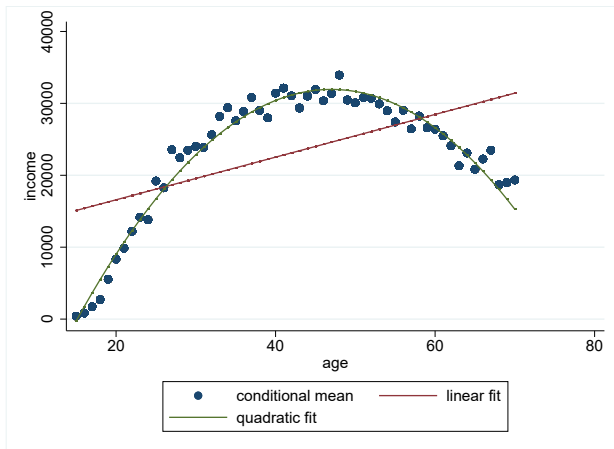
Table: Polynomial income regressions

	(1)	(2)
	income	income
Age	296.8*** (13.35)	2958.2*** (74.84)
Age2		-31.48*** (0.874)
Constant	10654.7*** (607.6)	-37549.0*** (1446.6)
Observations	5973	5973
r2	0.0764	0.241
F	493.9	949.9

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

What if the world is not linear?



What if the world is not linear?

- How to test for the order of the polynomial?
 - ① Start from an order that is a reasonable (high) one, such as 3 or 4.
 - ② Test down, i.e., whether the high(er) order term(s) are (jointly) significant.
- Notice: here you have a prior plan.
- Notice #2: a more modern version of this would involve a **semi-** or **non-parametric** approach (for later courses).
- Let's test whether a second order polynomial is sufficient by adding a third order term.
- Notice that for pedagogical purposes I am doing things in the **wrong** order.

What if the world is not linear?

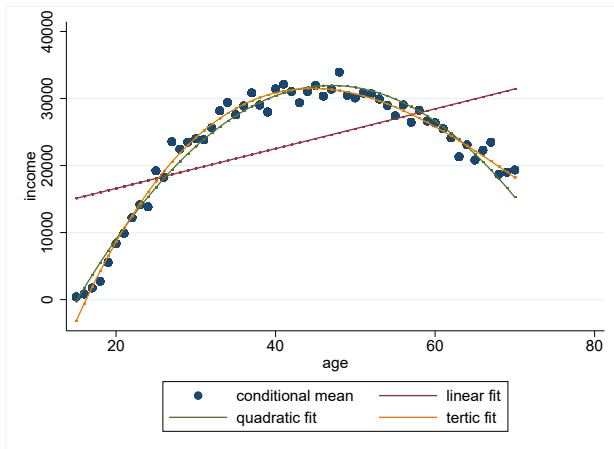
Table: Polynomial income regressions

	(1)	(2)	(3)
	income	income	income
Age	296.8*** (13.35)	2958.2*** (74.84)	4685.0*** (317.0)
Age2		-31.48*** (0.874)	-75.69*** (7.934)
Age3			0.346*** (0.0618)
Constant	10654.7*** (607.6)	-37549.0*** (1446.6)	-57597.8*** (3856.7)
Observations	5973	5973	5973
r2	0.0764	0.241	0.245
F	493.9	949.9	646.9

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

What if the world is not linear?



What if the world is not linear?

- Can you do more than use polynomials?
- Yes... though polynomials give very good approximations.

$$Y = f(X) + u$$

- Give $f(X)$ any shape you like.
- We will skip this for now (semi- and nonparametric estimation).

- What if there is reason to believe that the effect of X_1 depends on the value of X_2 ?
- Examples:
 - ① Returns to experience (=age) and/or education different by gender.
 - ② Effect of R&D subsidies different by firm size.

Example: Effect of age on income depends on gender

$$\text{Income} = f(\text{Age}, G, u) = \beta_0 + \beta_{\text{Age}} \times \text{Age} + \beta_G \times G + u$$

$$\begin{aligned} \text{Income} = f(\text{Age}, G, u) = & \beta_0 + \beta_{\text{Age}} \times \text{Age} + \beta_g \times G \\ & + \beta_{\text{Age}G} \times \text{Age} \times G + u \end{aligned}$$

- What is now the expected income | gender?
- What is now the expected income | age?

Q9 Example: Effect of age on income depends on gender

- How to calculate the effect of age on income?
- Now depends on the value of G directly.
- Notice
 - ① without the interaction $Age \times G$ the effect of age on income independent of G (= the same no matter what value G takes).
 - ② not true any more with the interaction.
- **Note:** if you add an interaction, make sure to have the original variables in the specification as well.

What if the world is not linear?

Table: Polynomial income regressions

	(1)	(2)
	income	income
Age	298.5*** (13.23)	333.5*** (18.80)
Gender	-4545.0*** (422.9)	-1598.7 (1203.3)
Age_G		-69.16** (26.44)
Constant	12819.2*** (634.6)	11336.3*** (850.8)
Observations	5973	5973
r2	0.0939	0.0950
F	309.4	208.8

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Q9 Transformations of variables

- What is a transformation of a variable?
- Use some $g(X)$ instead of X .
- Most often use (natural) log of X .
- Sometimes $\frac{1}{X}$.
- Always use a *monotonic* transformation!

Which variable to transform?

- Y , X , or both (all)?
- Using logs *smooths* the data, i.e., decreases the differences across different values that the variable takes.
- Taking logs allows negative values for a non-negative variable (if value < 1)
- On the other hand, cannot take logs if < 0 .

Log approximation to % change

$$\ln(Y + \Delta Y) - \ln(Y) \cong \frac{\Delta Y}{Y}$$

Which variable to transform?

1. Only Y

$$\ln \text{Income} = \beta_0 + \beta_{\text{Age}} \times \text{Age} + \beta_G \times G + u$$

$$\text{Income} = e^{\beta_0 + \beta_{\text{Age}} \times \text{Age} + \beta_G \times G + u}$$

$$= e^{\beta_0} e^{\beta_{\text{Age}} \times \text{Age}} e^{\beta_G \times G} e^u$$

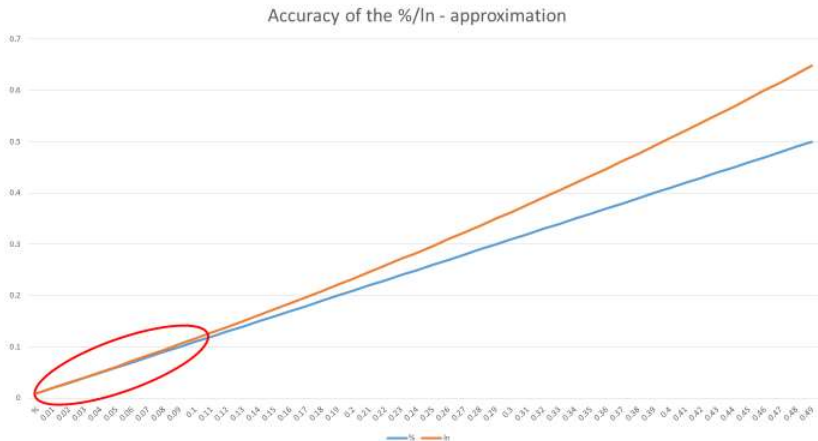
Which variable to transform?

2. Only X

$$\text{Income} = \beta_0 + \beta_{\ln \text{Age}} \times \ln \text{Age} + \beta_G \times G + u$$

- Interpretation of $\beta_{\ln \text{Age}}$?
- A 1% increase in Age is associated with at $0.01 \times \beta_{\ln \text{Age}}$ change in income.

Log approximation to % change



Which variable to transform?

3. Both Y and X

$$\ln \text{Income} = \beta_0 + \beta_{\ln \text{Age}} \times \ln \text{Age} + \beta_G \times G + u$$

- Interpretation of $\beta_{\ln \text{Age}}$?
- $\beta_{\ln \text{Age}} = \%$ -change in income due to a 1% increase in Age .
- In other words, $\beta_{\ln \text{Age}}$ is the **age elasticity of income**.

Which variable to transform?

Stata code

```
1 gen lnincome = ln(income)
2 gen lnage = ln(age)
3 regr income age gender if year == 15 & income != . & income_age_m != ., robust
4     eststo linear
5 regr lnincome age gender if year == 15 & income != . & income_age_m != ., robust
6     eststo loglin
7 regr income lnage gender if year == 15 & income != . & income_age_m != ., robust
8     eststo linlog
9 regr lnincome lnage gender if year == 15 & income != . & income_age_m != ., robust
10    eststo loglog
11 estout linear loglin linlog loglog, cells(b(star fmt(3)) se(par fmt(2))) ///
12 stats(r2 r2_a F N, fmt(%9.5f %9.5f %9.0g))
```

Table: Polynomial income regressions

	(1)	(2)	(3)	(4)
	income	lnincome	income	lnincome
Age	298.5*** (13.23)	0.0206*** (0.000723)		
Gender	-4545.0*** (422.9)	-0.143*** (0.0226)	-4503.1*** (412.7)	-0.140*** (0.0217)
lnAge			14059.7*** (486.4)	0.982*** (0.0269)
Constant	12819.2*** (634.6)	8.970*** (0.0350)	-26075.9*** (1806.8)	6.236*** (0.100)
Observations	5973	5695	5973	5695
r ²	0.0939	0.130	0.137	0.194
F	309.4	425.3	475.3	686.8

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Interpretations of β_{Age} , $\beta_{\ln Age}$ (mean age = 43)

- Linear: a 1 unit = 1 year ($\approx 2.3\%$) increase in age increases income by
300€
- Log-linear: a 1 unit = 1 year increase in age increases income by
 $2.1\% \rightarrow 0.021 \times 23296.67\text{€} \approx 490\text{€}$
- Linear-log: a 1% (note: $\ln(1.01) \approx 0.01$) increase in age (≈ 0.43 years) increases income by
 $0.01 \times 14059\text{€} \approx 140\text{€} \rightarrow$ effect of 1 year increase $\approx 2.5 \times 140\text{€} \approx 350\text{€}$
- Log-log: a 1% increase in age increases income by 0.982%. Effect of a 1 year increase
 $\approx 2.5 \times 0.00982 \times 23296.67\text{€} \approx 570\text{€}$

Interpretations of β_{Age} , $\beta_{\ln Age}$ (mean age = 43)

- Using the log of the dependent variable \rightarrow coefficient interpretation in **percent**.
- Typically, in economic data, using log of explanatory variable leads to higher R^2 .
- Many economic variables have a lower limit (income cannot be negative), but OLS assumes that the support is the real line.
 - \rightarrow log transformation allows coverage of the real line.
 - \rightarrow log transformation necessitates $Y(X) > 0$.

Q10 What can go wrong?

- 1 Internal validity.
- 2 External validity.

- 1 Omitted variable bias.
- 2 Functional form misspecification (mistake).
- 3 Measurement error in variable(s).
- 4 Sample selection (OVB).
- 5 Simultaneous (reverse) causality (OVB).
- 6 Non-homoskedastic errors.

Internal validity 1. - Omitted variable bias

- The relevant condition the one we have already discussed.
- "Judicious" choice of controls.
- Add variables.
- There are further solutions. We will get to these.

Internal validity 2. - Functional form

- How can you be sure?
 - ① Tests between the functional forms you try.
 - ② Note: can easily test only those functional forms that are "nested".

Example #1: 1st and 2nd order polynomial **nested** (= one is a restricted version of the other).

Example #2: log-log and linear are non-nested.

- Try out different ones and check robustness of your results (see earlier).

Internal validity 3. - Measurement error in variables

- Case #1: Y measured with error, error random.

$$Y_{observed} = Y + error$$

- Let's have a look at our regressions:

Regression we'd like to estimate:

$$Y = \beta_0 + \beta_1 X + u$$

Regression we can estimate:

$$\begin{aligned} Y_{observed} &= Y + error = \beta_0 + \beta_1 X + u + error \\ &= \beta_0 + \beta_1 X + v \end{aligned}$$

- Measurement error in Y not a big problem (as long as random).
- Leads to higher standard errors, but no bias.

- Case #2: X measured with error, error random.

$$X_{observed} = X + error, \rho_{X,error} = 0$$

- This is the case of so-called **classical errors-in-variables**. This case is "well-behaved".
- Let's have a look at our regression:

Internal validity 3. - Measurement error in variables

- We would like to estimate

$$Y = \beta_0 + \beta_1 X + u$$

- However we only observe $X_{observed} = X + error$. Hence we need to rewrite

$$\begin{aligned} Y &= \beta_0 + \beta_1 X + u \\ &= \beta_0 + \beta_1 X_{observed} - \beta_1 (X_{observed} - X) + u \\ &= \beta_0 + \beta_1 X_{observed} - \beta_1 error + u \\ &= \beta_0 + \beta_1 X_{observed} + v \end{aligned}$$

- One can show (see SW ch. 9.2):

$$\hat{\beta}_1 = \frac{\sigma_X^2}{\sigma_X^2 + \sigma_{error}^2} \beta_1$$

- $\frac{\sigma_X^2}{\sigma_X^2 + \sigma_{error}^2}$ = **signal to noise** - ratio.
- The larger is the role of the error, i.e., the larger is σ_{error}^2 relative to σ_X^2 , the more biased is $\hat{\beta}_1$.
- This is so-called **Attenuation bias**.

- Solution #1: Get better measures.
- Solution #2: Get a measure of σ_{error}^2 .
- Solution #3: A technical solution (**instrumental variables**) that we will get to later.

Internal validity 4. - Sample selection

- Your observations are not a random sample of the underlying population.
- Example #1: Estimate the returns to entrepreneurship using 5 year old firms.
- The non-profitable entrants exit.
- Example #2: Estimate the returns to graduating quickly.
- Those who graduate quickly have unobservable skills that make them (un)attractive to employers.

Internal validity 4. - Sample selection

- Example #3: estimate effects of R&D subsidies.
- Firms that get subsidies not avg. firms.
- Rule: think through and understand selection into your sample.
- Model selection into the sample.
- We will discuss this later, but in general is an advanced topic.

Internal validity 4. - Sample selection

- Sample selection can threaten internal validity - the parameters you obtain for the population of interest are biased.
- Sample selection can also threaten the external validity of your exercise, i.e., even if you get unbiased estimates for the population in question, your results do not generalize.

- Think of the determination of prices and quantities.
- Price affects how much is sold and produced.
- How much is bought and produced affects the price.
→ simultaneous causality.
- We will come back to this.

Internal validity 6. - Heteroskedasticity

- Deviations from homoskedasticity can take different forms depending on the data.
- With sequential observations, maybe also correlation over time (**autocorrelation**).
- With e.g. geographical data, correlation across observation units (**clustering**).
- Affects statistical precision (=standard errors) of individual coefficients, nothing else.
- Can be corrected by using **robust** standard errors (with potential loss of efficiency - but robust se's can be smaller than homosc. se's).
- In data with relevant other dimensions (e.g. geographical locations), clustered se's may be more appropriate than regular robust se's.

Internal validity 6. - heteroskedasticity

Stata code

```
1  regr income age gender if year == 15 & income != . & income_age_m != .  
2     eststo linear  
3  regr income age gender if year == 15 & income != . & income_age_m != ., robust  
4     eststo linear_het
```

Table: Polynomial income regressions

	(1)	(2)
	homosc.	heterosk. robust
Age	298.5*** (13.23)	298.5*** (11.82)
Gender	-4545.0*** (422.9)	-4545.0*** (421.6)
Constant	12819.2*** (634.6)	12819.2*** (566.5)
Observations	5973	5973
r ²	0.0939	0.0939
F	309.4	369.3

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

- Econometrics:

A branch of economics in which economic theory and statistical methods are fused in the analysis of numerical and institutional data

Hood, W. & Koopmans, T. (1953). Studies in econometric method. *Cowles Commission Monograph no. 14*, Wiley

- External validity = results generalize to other settings than the one studied.
 - Any (material) change to any of the components of your study jeopardizes external validity.
- ① Differences in (applicable) theory.
 - ② Differences in statistical method.
 - ③ Differences in data (including in populations).
 - ④ Differences in institutions.

- Example: Do our income - age results hold for (an)other year in the FLEED data?
- Let's compare results from current year 15 to year 10.

External validity: comparison of results from two very similar data

Table: Polynomial income regressions

	(1) year 15	(2) year 10
Age	298.5*** (13.23)	216.4*** (11.74)
Gender	-4545.0*** (422.9)	-4765.7*** (361.7)
Constant	12819.2*** (634.6)	12167.0*** (568.5)
Observations	5973	5779
r2	0.0939	0.0803
F	309.4	252.2

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

- Is any study externally valid?
- Yes and no.
- Best to ensure internal validity, and conduct many studies.