

# Quantitative Analysis and Empirical Methods

## 2) Research Design

Jan Rovny

Sciences Po, Paris, CEE / LIEPP  
University of Gothenburg, CERGU / Political Science

1 Explaining Relationships

2 Research Design

# Introduction

- We have learned about converting concepts into variables
- Our key interest is to understand relationships between variables
- And even more importantly, understand how changes in one variable affect changes in another variable
- To do this, we must first consider the design of our research and consider the threats to our research validity

# Explaining Relationships

# Dependent and Independent Variables

- Dependent variable: values you are trying to explain
- Independent variable (predictor, regressor, explanatory variable):
  - what we think causes or is causally related to the dependent variable

Hypothesis { Independent variable  $\rightarrow$  Relationship  $\rightarrow$  Dependent variable  
or  
Dependent variable =  $f$ (Independent variable)

# Hypothesis

- Hypothesis is an explicit testable statement
- Usually associates the dependent and an independent variable
- It can be about a single variable



# Characteristics of Hypotheses 1

- 1) Hypotheses should be as specific as possible and provide a clear tendency between the variables
  - H1.1 Economic development is related to democracy.
  - H2.1 Democracies are more peaceful.
  - H3.1 Income varies across countries and over time.

# Characteristics of Hypotheses 1

- 1) Hypotheses should be as specific as possible and provide a clear tendency between the variables
  - H1.1 Economic development is related to democracy.
  - H1.2 The greater its economic development, the more democratic a society is.
  - H2.1 Democracies are more peaceful.
  - H3.1 Income varies across countries and over time.



# Characteristics of Hypotheses 1

- 1) Hypotheses should be as specific as possible and provide a clear tendency between the variables
  - H1.1 Economic development is related to democracy.
  - H1.2 The greater its economic development, the more democratic a society is.
  - H2.1 Democracies are more peaceful.
  - H2.2 Democracies are less likely to go to war with each other, but go to war with non-democracies as often as non-democracies fight each other.
  - H3.1 Income varies across countries and over time.

# Characteristics of Hypotheses 1

- 1) Hypotheses should be as specific as possible and provide a clear tendency between the variables
  - H1.1 Economic development is related to democracy.
  - H1.2 The greater its economic development, the more democratic a society is.
  - H2.1 Democracies are more peaceful.
  - H2.2 Democracies are less likely to go to war with each other, but go to war with non-democracies as often as non-democracies fight each other.
  - H3.1 Income varies across countries and over time.
  - H3.2 Income is distributed more equally in Sweden than in the United States.

# Characteristics of Hypotheses 1

- 1) Hypotheses should be as specific as possible and provide a clear tendency between the variables
  - H1.1 Economic development is related to democracy.
  - H1.2 The greater its economic development, the more democratic a society is.
  - H2.1 Democracies are more peaceful.
  - H2.2 Democracies are less likely to go to war with each other, but go to war with non-democracies as often as non-democracies fight each other.
  - H3.1 Income varies across countries and over time.
  - H3.2 Income is distributed more equally in Sweden than in the United States.
  - H3.3 Income inequality has grown since 1990 in eastern Europe.

## Characteristics of Hypotheses 2

- 2) Hypotheses must be falsifiable (and not tautological)
  - H1.1 Individuals either vote or abstain.
  - H2.1 Strict judges impose strict punishment on convicted defendants.

## Characteristics of Hypotheses 2

- 2) Hypotheses must be falsifiable (and not tautological)
  - H1.1 Individuals either vote or abstain.
  - H1.2 More informed voters are more likely to vote.
  - H2.1 Strict judges impose strict punishment on convicted defendants.

## Characteristics of Hypotheses 2

- 2) Hypotheses must be falsifiable (and not tautological)
  - H1.1 Individuals either vote or abstain.
  - H1.2 More informed voters are more likely to vote.
  - H2.1 Strict judges impose strict punishment on convicted defendants.
  - H2.2 More junior judges impose stricter punishment on convicted defendants.

## Characteristics of Hypotheses 3

- 3) Hypotheses must be empirically testable  
We must be able to collect information (data) and examine it to see if the hypothesized relationship holds.
  - H1.1 It was God's will for Portugal to win the Euro Cup.

# Characteristics of Hypotheses 3

- 3) Hypotheses must be empirically testable  
We must be able to collect information (data) and examine it to see if the hypothesized relationship holds.
  - H1.1 It was God's will for Portugal to win the Euro Cup.
  - H1.2 It was better passing that led to Portugal winning the Euro Cup.



# The Bigger Picture

- A set of propositions – some of which are testable hypotheses – is a **theory**
- Not all propositions are testable, we call these **assumptions**
- A simplified form of a theory is called a **model**, which can be expressed verbally, graphically or mathematically.
- We will learn to model hypotheses and theories using Stata.

# Model Example 1

Theory: Economic liberalization causes democratization

Model:

- liberalization  $\rightarrow$  investment
- investment  $\rightarrow$  profit
- profit  $\rightarrow$  bourgeoisie
- bourgeoisie  $\rightarrow$  political interests
- political interests  $\rightarrow$  democratization



# Model Example 1

Theory: Economic liberalization causes democratization

Model:

- liberalization → investment
- investment → profit
- profit → bourgeoisie
- bourgeoisie → political interests
- political interests → democratization



- sheep → commercial mindedness of aristocracy → economic liberalization (with apologies to Barrington Moore)

## Model Example 2

Theory: Vote choice is determined an individual's socio-economic profile

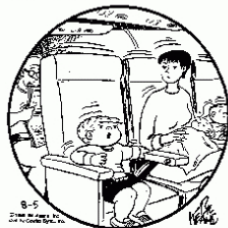
Model:

$$\textit{vote} = \beta_0 + \beta_1 \textit{gender} + \beta_2 \textit{age} + \beta_3 \textit{education} + \beta_4 \textit{income} + \beta_5 \textit{religiosity} + \beta_6 \textit{community}$$

# Causality

- The goal of social sciences is to provide causal explanations
- These are explanations that state that A causes B (rather than the opposite)
- Causality is impossible to prove (statistically), it is only possible to argue
- Correlation  $\neq$  Causation!!!

THE FAMILY CIRCUS



"I wish they didn't turn on that seatbelt sign so much! Every time they do, it gets bumpy."

# Causal Explanation

- 1) Time order
  - Cause must precede effect
  - Granger Causality
- 2) Covariation
  - $\Delta Y$  must be associated with  $\Delta X$
  - Cannot explain change with a constant
- 3) Non-spuriousness
  - The relationship between X and Y should not be driven by a third variable
  - Problem of proving a negative
  - Problem of multiple causation
- 4) Theoretical consistency
  - It needs to make sense...
  - Problem of data-mining

# Research Design

# Research Design

- Research design is a set of procedures for testing a hypothesis about a relationship.
- It is the method of assessing the effect(s) of the independent variable(s) or treatments on the dependent variable or outcome



# Counterfactual

- A way of assessing causality
- Idea that each individual can be exposed to two alternative states of a cause
- Assumption is that each individual has a potential outcome under each treatment
- That is, individual  $i$  can have an outcome  $y_i^1$  and  $y_i^0$
- Individual level causal effect can be defined as  $y_i^1 - y_i^0$
- However, we cannot observe both  $y_i^1$  and  $y_i^0$  at individual level
- We can only measure  $y_i$ ,
- $y_i = y_i^1$  in treatment state,  $y_i = y_i^0$  in control state

# Example

What is the effect of taking the math camp (C) on your math skills (S)?

- Take two groups of students: 1 and 0 (ideally identical)
- Measure initial skill:  $S_{0,t0}, S_{1,t0}$
- Group 1 takes C (treatment)
- Measure skills again:  $S_{0,t1}, S_{1,t1}$
- The effect of C on S is the difference between the differences in skills =  $(S_{1,t1} - S_{1,t0}) - (S_{0,t1} - S_{0,t0})$

# Forms of Validity

- Internal Validity
  - Considers whether our claims are valid (correct) in the particular experiment at hand
  - Experiment specific logic of design and validity of claims
  - Assessment is based on the 4 criteria of assessing causal explanations
- External Validity
  - Considers the extent to which the results obtained in a given study are generalizable to different populations, contexts, times...

# Forms of Research Design

- True Experiments
  - Researcher has control over sample selection, treatment, measurement and analysis
  - Typically division into experimental and control group with randomized or matched assignment before experimental treatment.
- Quasi-Experiments
  - No control over sample selection or treatment
  - Natural experiments (before / after some shock); panel data; cross-sectional data; controlled comparisons

# True Experiments

Group	Assignment	Obs 1	Treatment	Obs 2	Comparison
Experimental	Random	$O_{e1}$	Yes	$O_{e2}$	$O_{e2} - O_{e1}$
Control	Random	$O_{c1}$	No	$O_{c2}$	$O_{c2} - O_{c1}$

- Excel at internal validity but have limited external validity
  - Context
  - Non-representative sample

# Quasi-Experiments

Design							
Cross-section	X	O					
Panel study	$O_1$	X	$O_2$				
Trend study	$O_1$	X	$O_2$	$O_3$	X	$O_4$	$O_5$

- Stronger external validity, weaker internal validity
- Use of statistics to establish association between cause and effect
- With longitudinal studies may be able to establish time-order, but not always (cointegration)
- Problem of possible spuriousness threatens internal validity
- To deal with spuriousness we employ statistical control techniques

# Common threats to Internal Validity

- Omitted Variables (spuriousness)
- Regression to the mean
  - Groups selected on the basis of extreme scores will tend to display lower scores next time
- Non-random sample selection
  - Difference in comparison groups is not due to treatment but because groups were different to start with
  - Selecting on the dependent variable

# Example 1

At the age of 5 your parents give you an algebra exam and you fail. Rather than being punished, you argue with your parents that you need greater incentives to study. They agree to give you an ice cream every time you do your homework for the next 10 years. At the age of 15 you take another algebra exam and you do well. You claim that it is thanks to the ice cream.

- What is the dependent variable?
- What is the independent variable?
- What is the underlying research design?
- What are the possible validity problems?



# Example 1

At the age of 5 your parents give you an algebra exam and you fail. Rather than being punished, you argue with your parents that you need greater incentives to study. They agree to give you an ice cream every time you do your homework for the next 10 years. At the age of 15 you take another algebra exam and you do well. You claim that it is thanks to the ice cream.

- What is the dependent variable?
- Performance on an algebra exam (math skills)
- What is the independent variable?
  
- What is the underlying research design?
  
- What are the possible validity problems?

# Example 1

At the age of 5 your parents give you an algebra exam and you fail. Rather than being punished, you argue with your parents that you need greater incentives to study. They agree to give you an ice cream every time you do your homework for the next 10 years. At the age of 15 you take another algebra exam and you do well. You claim that it is thanks to the ice cream.

- What is the dependent variable?
- Performance on an algebra exam (math skills)
- What is the independent variable?
- Incentives (getting ice cream)
- What is the underlying research design?
  
- What are the possible validity problems?

# Example 1

At the age of 5 your parents give you an algebra exam and you fail. Rather than being punished, you argue with your parents that you need greater incentives to study. They agree to give you an ice cream every time you do your homework for the next 10 years. At the age of 15 you take another algebra exam and you do well. You claim that it is thanks to the ice cream.

- What is the dependent variable?
- Performance on an algebra exam (math skills)
- What is the independent variable?
- Incentives (getting ice cream)
- What is the underlying research design?
- Panel study (test, treatment, re-test)
- What are the possible validity problems?

# Example 1

At the age of 5 your parents give you an algebra exam and you fail. Rather than being punished, you argue with your parents that you need greater incentives to study. They agree to give you an ice cream every time you do your homework for the next 10 years. At the age of 15 you take another algebra exam and you do well. You claim that it is thanks to the ice cream.

- What is the dependent variable?
- Performance on an algebra exam (math skills)
- What is the independent variable?
- Incentives (getting ice cream)
- What is the underlying research design?
- Panel study (test, treatment, re-test)
- What are the possible validity problems?
- Omitted variables (age, school attendance)

## Example 2

A group of elite pilots train together for a month. At the end of each day the best performers are praised. The next day, those who were praised the day before rarely perform well.

- Why?

## Example 2

A group of elite pilots train together for a month. At the end of each day the best performers are praised. The next day, those who were praised the day before rarely perform well.

- Why?
- Regression to the mean

## Example 3

Does use of soft drugs lead to use of hard drugs? Some (imaginary) data: 97% of heroine addicts have smoked marijuana.

- What is the problem of this design?

## Example 3

Does use of soft drugs lead to use of hard drugs? Some (imaginary) data: 97% of heroine addicts have smoked marijuana.

- What is the problem of this design?
- Selection of the dependent variable (how many people who smoked marijuana did not go on to use heroine?)



## Example 3

Does use of soft drugs lead to use of hard drugs? Some (imaginary) data: 97% of heroine addicts have smoked marijuana.

- What is the problem of this design?
- Selection of the dependent variable (how many people who smoked marijuana did not go on to use heroine?)
- 100% of heroine addicts have drunk milk

# Common threats to External Validity

- Selection
  - Groups are not representative of the larger population (not randomized)
  - Example: A study on our class
- Out of sample extrapolation
  - Setting is out of bounds of treatment
  - Example: a study that administers dosages from 0 to 40ml when doctors administer 100.

# Example 1

You notice that students who spend 2 hours on the homework get a grade of 80%, while those who spend 4 hours on it get a grade of 95%. You conclude that spending 6 hours on homework will earn you a 110%.

- What is the problem here?

# Example 1

You notice that students who spend 2 hours on the homework get a grade of 80%, while those who spend 4 hours on it get a grade of 95%. You conclude that spending 6 hours on homework will earn you a 110%.

- What is the problem here?
- Out of sample extrapolation

## Example 2

Studying Eurosceptic parties, you notice that they are both on the left and right side of the political spectrum and conclude that there is no effect of left-right placement on attitudes towards the European Union.

- What is the problem here?

## Example 2

Studying Eurosceptic parties, you notice that they are both on the left and right side of the political spectrum and conclude that there is no effect of left-right placement on attitudes towards the European Union.

- What is the problem here?
- Selecting on the dependent variable