# Support The Guardian
Available for everyone, funded by readers

**Contribute →**   **Subscribe →**

Search jobs     Sign in     Search

International edition

# The Guardian

News     Opinion     Sport     Culture     Lifestyle     More ⌄

Fashion   Food   Recipes   Love & sex   Health & fitness   Home & garden   Women   Men   Family   Travel   Money

**Running**

# Any amount of running reduces risk of early death, study finds

**Previous research suggested health benefits increased with greater volume of running**

## Nicola Davis
🐦 @NicolaKSDavis

Mon 4 Nov 2019
23.30 GMT

3149

**most viewed**

Business

# Want to live longer? Try getting a dog.

Dog ownership significantly lowers mortality risk; now researchers are trying to find out just how they keep people alive.

Researchers have attached a laundry list of health benefits to dog ownership. Dogs not only "offer companionship, reduce anxiety and loneliness, increase self-esteem, and improve overall mood," but also force their humans to exercise and spend more time outdoors. (iStock)

By **Christopher Ingraham**

# Japanese passenger cars sold in the US
## correlates with
# Suicides by crashing of motor vehicle

Japanese cars sold

1200 thousand cars
1000 thousand cars
800 thousand cars
600 thousand cars

Suicides by crashing

140 suicides
120 suicides
100 suicides
80 suicides

1999  2000  2001  2002  2003  2004  2005  2006  2007  2008  2009

● Suicides by crashing     ◆ Japanese cars sold

tylervigen.com

http://tylervigen.com/spurious-correlations

# Number of people who drowned by falling into a pool

correlates with

## Films Nicolas Cage appeared in



Nicholas Cage    ◆ Swimming pool drownings
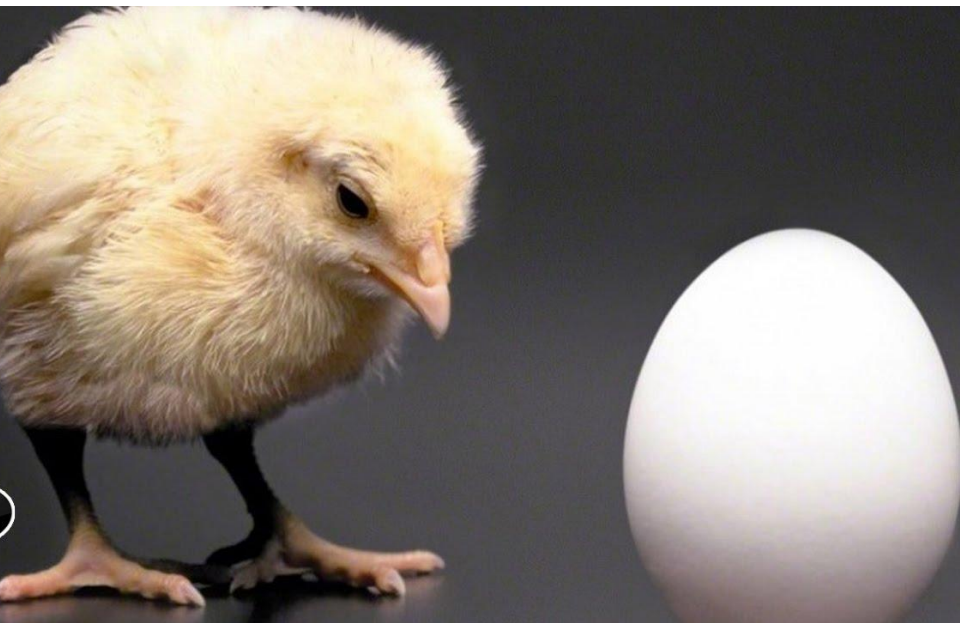
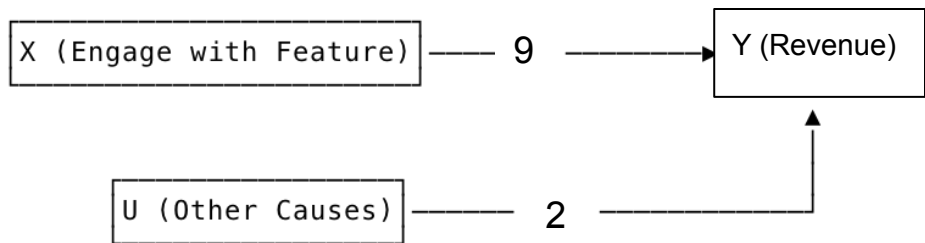# WHY CAUSALITY MATTERS

---

- Actions are taken based on their potential outcomes
- Can still make accurate predictions based on correlations

zalando

# WHY CAUSALITY MATTERS

- A machine learning model may say that feature X has a statistically significant positive impact on Revenue (~9 euros more per hour) and the model has a very high R^2: ~0.96.

- 

zalando

## WHY CAUSALITY MATTERS

- Interest in Fashion is actually the cause of Engagement and Revenue
- The positive effect of Engagement with Revenue is actually just the positive effect of Interest in Fashion passing through (and actually decreased) by Time on Site.

# ESTABLISHING CAUSALITY:

# A/B TESTING

- Gold Standard: Randomised control trial
- Not always possible/ethical/affordable/in line with strategy

zalando

# ESTABLISHING CAUSALITY:

# OBSERVATIONAL METHODS

- Data is obtained passively, without designing an experiment
- Care must be taken to control for Confounding variables and selection bias
- 

zalando

# ESTABLISHING CAUSALITY:

# OBSERVATIONAL METHODS

- Simpson's Paradox
- https://www.youtube.com/watch?v=ebEkn-BiW5k

zalando

## METHODS

- Directed Acyclic Graphs

- Instrumental Variables Analysis

- Matching

zalando

# DIRECTED ACYCLIC GRAPHS - DAGS

**x** Too many steps in the checkout

**DIRECTED**

**y** Cart is abandoned

zalando

ACYCLIC

zalando

# Structural Causal Graphs



- Great tool for conceptualising causal relationships

- Mathematically grounded

- Implies Linear relationships

- R package **ggdag**

- Python package **causality**

zalando

# Structural Causal Graphs



- Two models specified here
- Left has genuine causal relationship between x4 and x5
- Right has a spurious correlation between x4 and x5
- They are in fact both caused by a common cause x6 - confounder
- graph inference

https://medium.com/@akelleh/causal-graph-inference-b3e3afd47110

zalando

# Structural Causal Graphs



- Two models specified here
- Left has genuine causal relationship between x4 and x5
- Right has a spurious correlation between x4 and x5
- They are in fact both caused by a common cause x6 - confounder
- graph inference

https://medium.com/@akelleh/causal-graph-inference-b3e3afd47110

zalando

**INSTRUMENTAL VARIABLE ANALYSIS**

Q: What to do when you can only A/B test something that is associated with the thing you want to test ?

A: Instrumental Variable Analysis  - Use a third variable that you can measure that is strongly associated with the one you actually want to measure

zalando

Instrumental Variables have three characteristics
1.   Associated/Correlated with variable whose impact we want to understand

Unobserved
Confounders
Eg Income

**Please** Sign Up
for Zalando
Service

Sign Up for
Zalando Service

Revenue

Encouragement

Action You Care
About

Outcome You
Care About

zalando

21

## INSTRUMENTAL VARIABLE ANALYSIS

Instrumental Variables have three characteristics

1. Associated/Correlated with variable whose impact we want to understand

2. Does not impact outcome, except via potential effect on the variable we want to understand (exclusion restriction)



**Please** Sign Up for Zalando Service

Encouragement

Sign Up for Zalando Service

Action You Care About

Revenue

Outcome You Care About

zalando

# INSTRUMENTAL VARIABLE ANALYSIS

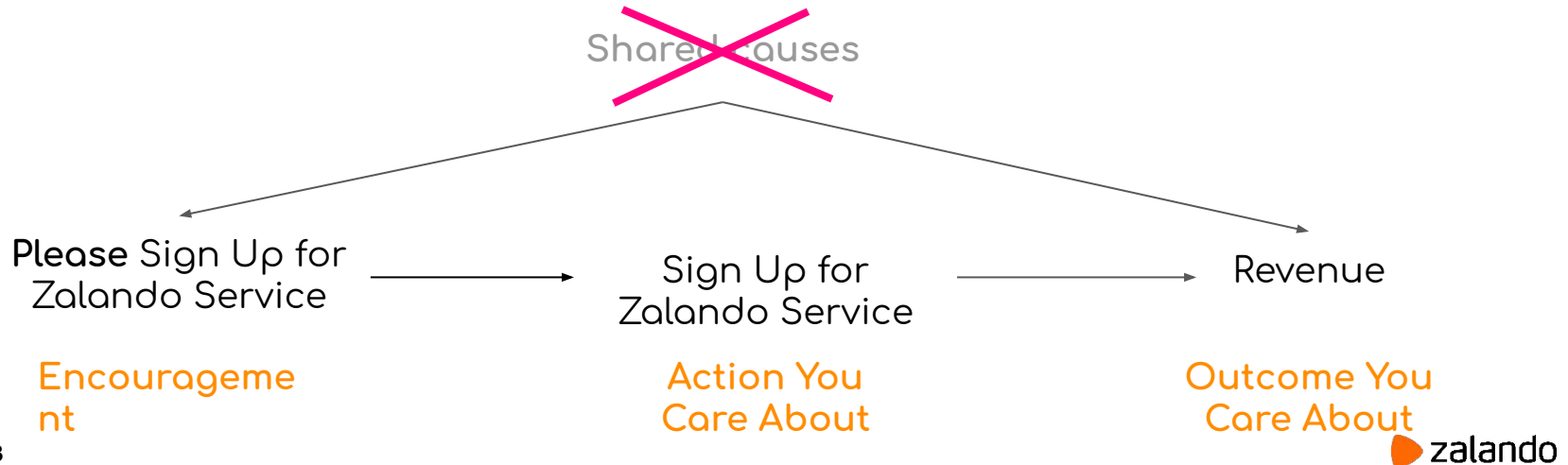Instrumental Variables have three characteristics
1. Associated/Correlated with variable whose impact we want to understand

2. Does not impact outcome, except via potential effect on the variable we want to understand (exclusion restriction)
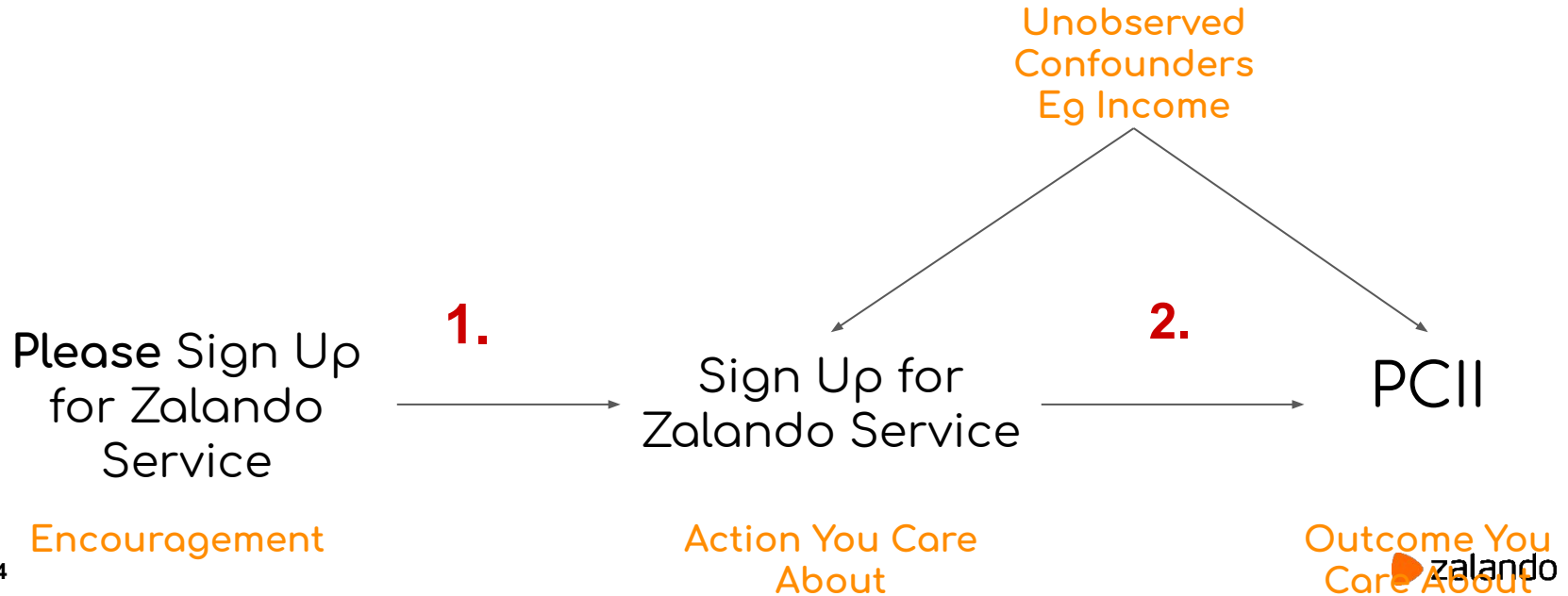
3. Outcome and instrument do not share causes

Shared causes

Please Sign Up for Zalando Service → Sign Up for Zalando Service → Revenue

Encouragement

Action You Care About

Outcome You Care About

zalando

**INSTRUMENTAL VARIABLE ANALYSIS**

Instrumental Variables Analysis  via   Two Stage Least Squares



Unobserved Confounders Eg Income

**1.**

**2.**

**Please** Sign Up for Zalando Service

Sign Up for Zalando Service

PCII

Encouragement

Action You Care About

Outcome You Care About

## First Stage

$$\text{Sign up} = \alpha + \beta \text{Encouragement}$$

## Second Stage

$$\text{Revenue} = \alpha + \beta \hat{\text{Sign}} \text{ up}$$

Prediction from first stage

zalando

# INSTRUMENTAL VARIABLE ANALYSIS

Instrumental Variables Analysis  via   Two Stage Least Squares

- R package: AER package ivreg method

- Python package:  Linear Models package IV2SLS method

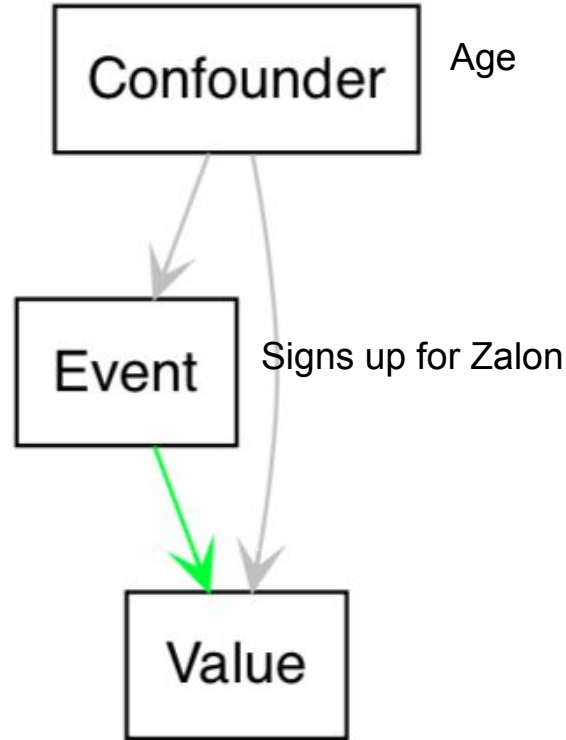- Book: Pearl & Mackenzie 2018, 249

Limitations:

- Instrument must be justified - have a strong association with the variable of interest

- Instrument must not have direct impact on the outcome

- If you misuse instrumental variables, you can get even more biased results than if you had not used them at all.

zalando

**MATCHING**

- A method to estimate causality without an experiment, or in the presence of a "Natural" experiment
- E.g. Launch of a new product

zalando
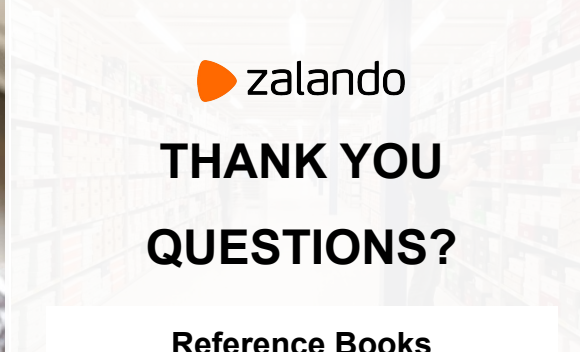
## MATTING



Confounder — Age

Event — Signs up for Zalon

Value

- Must be able to measure the confounder
- Naive estimate of the Event and the Value is biased due to the confounder
- Instead measure all relationships
- "Control" for Age
- "Matching Estimator"
- https://cran.r-project.org/web/packages/Matching/Matching.pdf

zalando

**SUMMARY**

- Correlation is not Causation

- Understanding cause is important for business decision making

- Can still make accurate predictions without knowing cause

- Controlled Experiments are the Gold Standard

- Can still estimate cause from Observational Data

- Methods:

  - Directed Acyclic Graphs

  - Instrumental Variables

  - Matching

zalando

## zalando

# THANK YOU
# QUESTIONS?

**Reference Books**

*Causality: Models, Reasoning and Inference* (2009) Judea Pearl

*The Book of Why: The New Science of Cause and Effect* (2018) Pearl & Mackenzie

**Reference Blogs**

https://medium.com/causal-data-science Adam Kelleher

**My Github & Slides**

https://github.com/alicelynch/meetup-talks