# What is Data Science?

Ernst C. Wit
Università della Svizzera italiana
Institute of Computing

wite@usi.ch

https://www.ci.usi.ch/

15 March 2022

# What do you think Data Science is about?



Data Science ...

# What do you think Data Science is about?
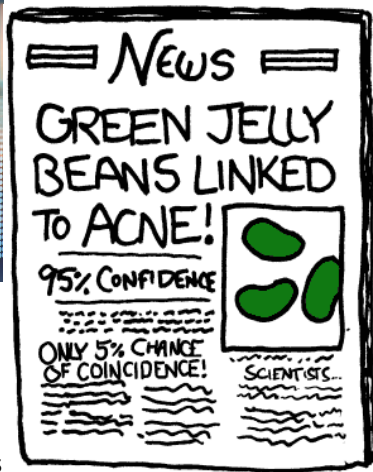


Data Science ...

- ... analyzes loads of data
- ... use artificial intelligence
- ... discovers secret patterns, such as

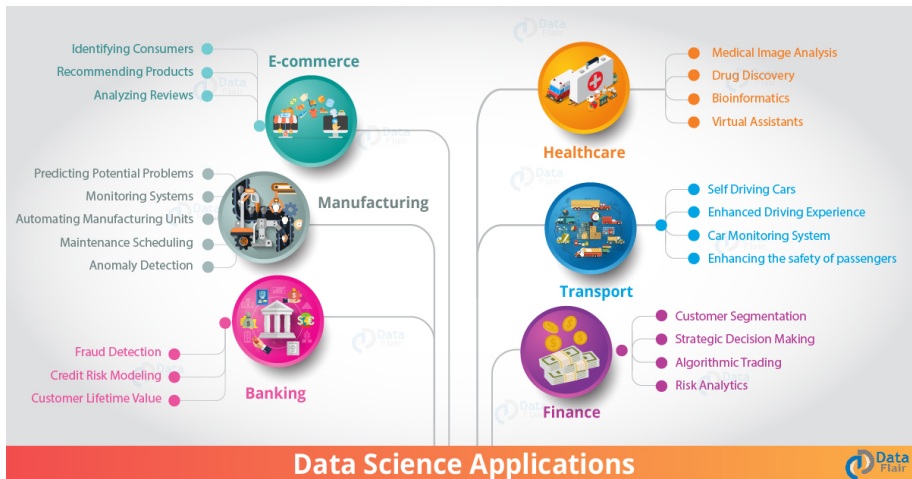# What do you think Data Science is about?



Data Science ...

- ... analyzes loads of data
- ... use artificial intelligence
- ... discovers secret patterns, such as

# Where is Data Science used?



Data Science Applications

# How does Data Science work?

It starts out with a **question**:

- What causes Y (e.g. fraud, Covid infection, engineering faults)?
- How to predict Y (e.g. consumers, new drugs, disease)?

It then gathers **data**...

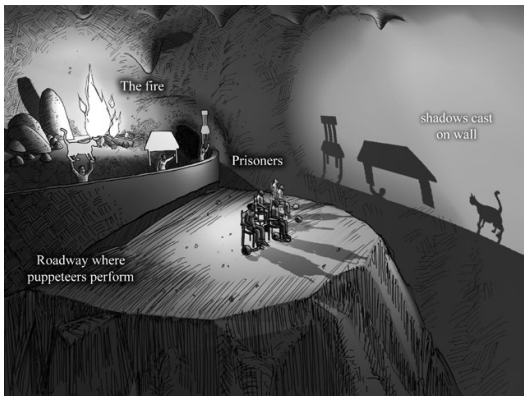- by designing a study
- by collecting what is available

Then use **statistical data science techniques** to analyze data:

- plotting data
- model relationships in data
- formulate conclusions

**shadows** = data, **cat/table/chair** = true answer
**flickering fire** = random noise, sampling, ascertainment bias, confounding



**Data Science** is controlling **fire** so **prisoners**
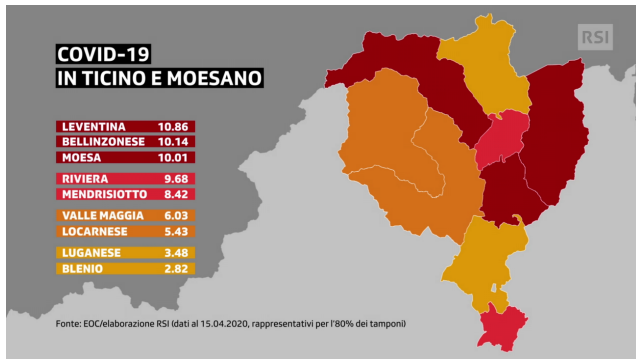learn **real** objects only from **shadows**.

# What **TRUTH** do we want to discover?

# Covid-19
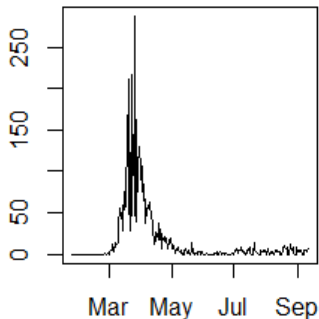
We want to study three questions in this class:

1. How did pandemic evolve in Ticino?

2. How was pandemic affected by inter-cantonal transmission?
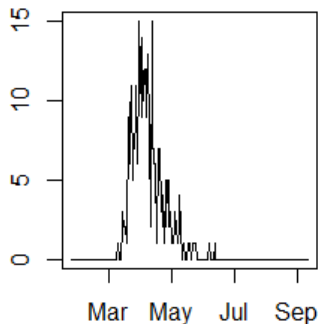
3. Did mortality rate for Covid19 improve?

# Shadows = Data



shadows cast on wall

# Let's look at Swiss Covid-19 data (2020)



Daily cases in TI

Daily deaths in TI

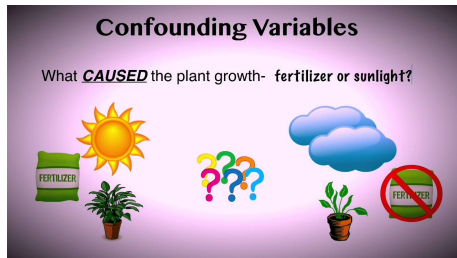# Fire = Noise, sampling bias, confounding, ...



The fire

# Sources of Uncertainty: "Noise"

- "Chance"
  - Measurement uncertainty;
  - Intrinsic system noise;
- "Sampling"
  - Experimental design
- "Confounding"

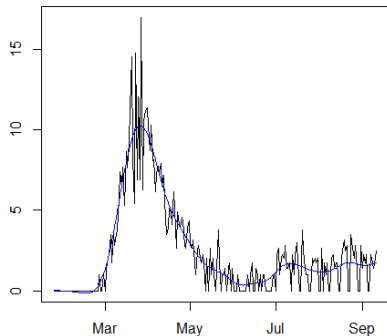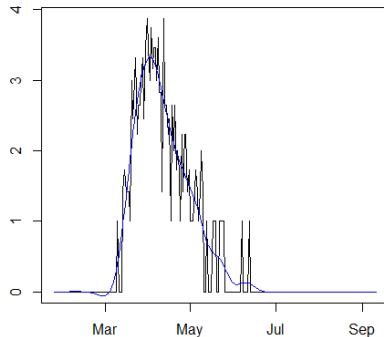# Sources of Uncertainty: "Noise"

- **"Chance"**
  - Measurement uncertainty;
  - Intrinsic system noise;
- "Sampling"
  - Experimental design
- "Confounding"

- "Chance"
  - Measurement uncertainty;
  - Intrinsic system noise;
- **"Sampling"**
  - Experimental design
- "Confounding"

- "Chance"
  - Measurement uncertainty;
  - Intrinsic system noise;
- "Sampling"
  - Experimental design
- **"Confounding"**



**Confounding Variables**

What **_CAUSED_** the plant growth- **fertilizer or sunlight?**

# How to deal with noise?



**Sqrt Cases TI**

**Sqrt Deaths TI**

We do 2 things:

- transform data with $\sqrt{\phantom{x}}$ (takes away large extremes)
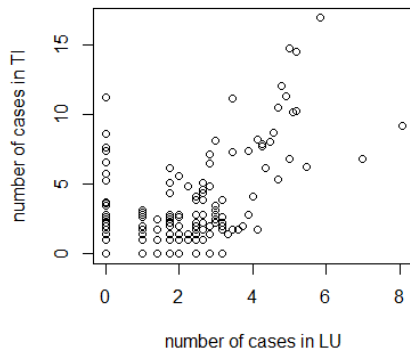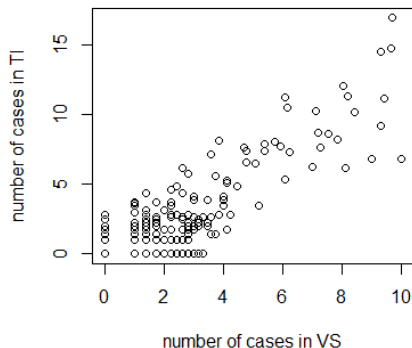- smooth data (spreads deaths/cases over empty weekends)

# TRUTH

# How did infections spread through Switzerland?

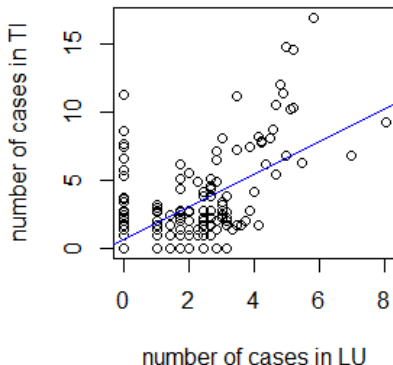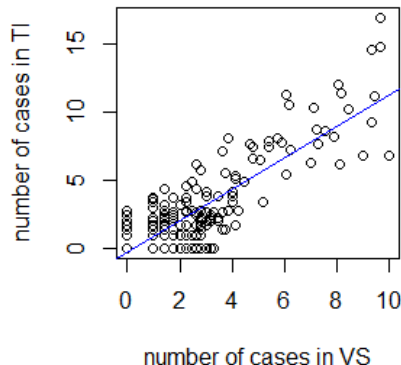Consider relationship between TI and other cantons:



**Question:** Who affected who?
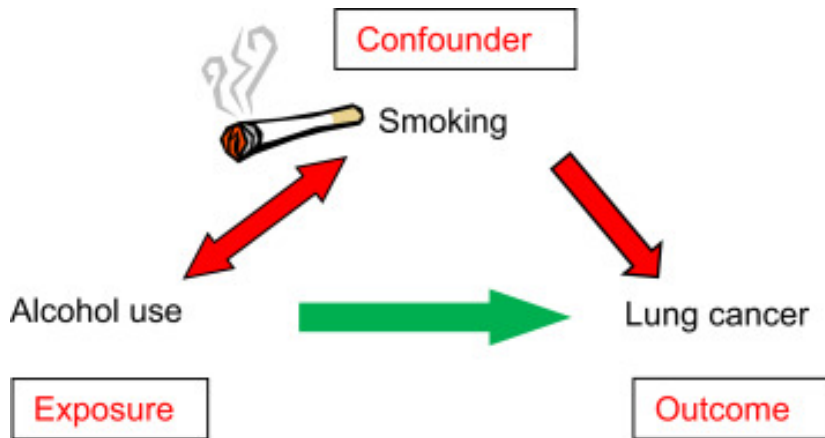
# Linear Regression: modelling relationships

We can model relationships between TI and other cantons:

$$\sqrt{\text{cases in TI}} = \alpha_1 + \beta_1 \sqrt{\text{cases in VS}} + \text{noise}$$
$$\sqrt{\text{cases in TI}} = \alpha_2 + \beta_2 \sqrt{\text{cases in LU}} + \text{noise}$$

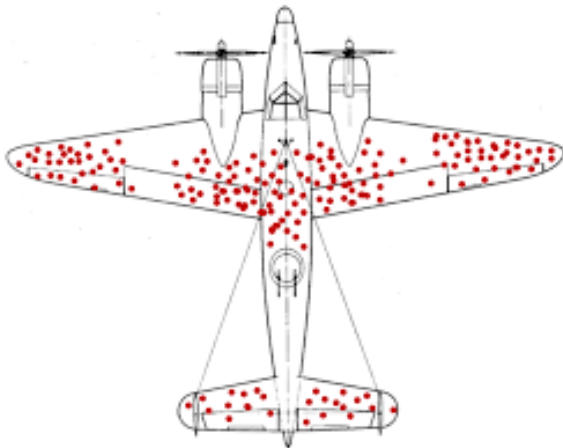# Sources of Uncertainty: Confounding

# Confounding: what data do you get to see?

Different types

- Non-response bias
- Healthy user bias
- Berkson's fallacy
- Overmatching
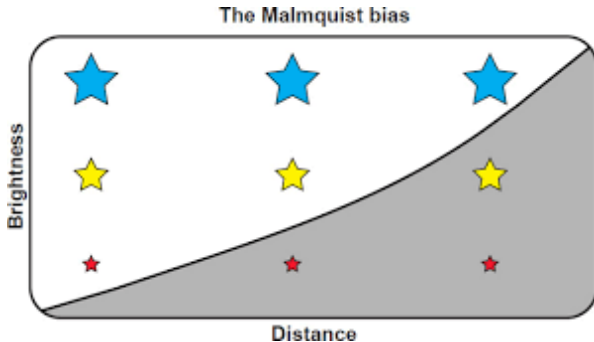- **Survivorship bias**
- Malmquist bias

Different types

- Non-response bias
- Healthy user bias
- Berkson's fallacy
- Overmatching
- Survivorship bias
- **Malmquist bias**



The Malmquist bias

(Brightness vs Distance)

# Linear Regression: modelling complex relationships

We can model relationships between TI and other cantons TOGETHER:

$$\sqrt{\text{cases in TI}} \;=\; \beta_0 + \beta_1\sqrt{\text{cases in VS}} + \beta_2\sqrt{\text{cases in LU}} + \text{noise}$$
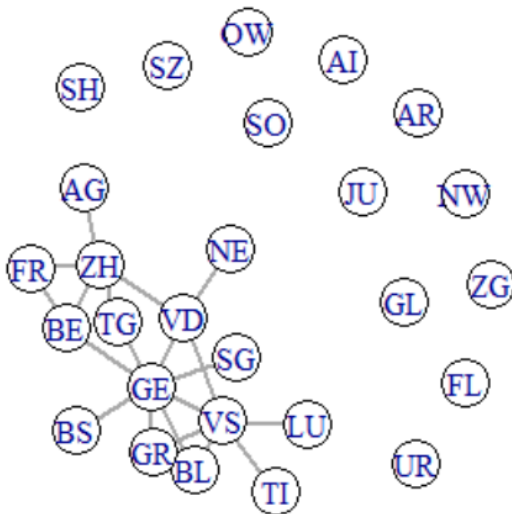
We can model relationships between TI and other cantons TOGETHER:

$$\sqrt{\text{cases in TI}} = \beta_0 + \beta_1\sqrt{\text{cases in VS}} + \beta_2\sqrt{\text{cases in LU}} + \text{noise}$$
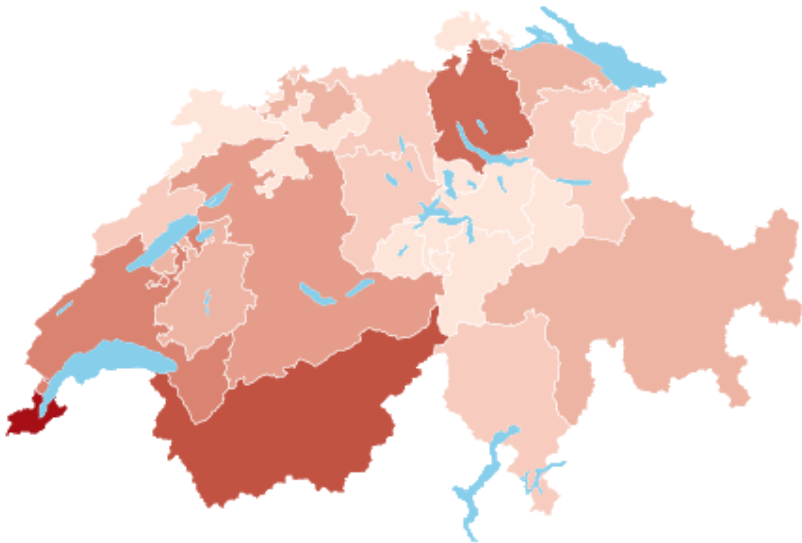
|  | Estimate | Std. Error | t value | Pr($>$|t|) |
|---|---|---|---|---|
| $\beta_0$ | -0.26 | 0.15 | -1.66 | 0.09 |
| $\beta_1$ (VS) | 1.11 | 0.07 | 16.3 | $< 0.01$ |
| $\beta_2$ (LU) | 0.07 | 0.09 | 0.74 | 0.48 |

So, VS affected TI, but LU didn't!

# ... Now do this for all cantons simultaneously!

# Direction of arrow?

CORRELATION IS NOT CAUSATION!

ICE CREAM SALES
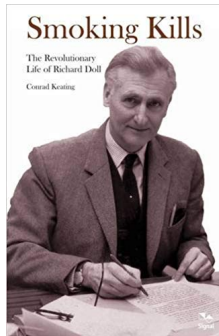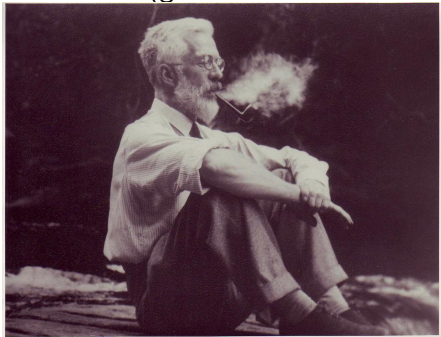SHARK ATTACKS

JAN    MAR    MAY    JUL    SEP    NOV

Both ice cream sales and shark attacks increase when the weather is hot
and sunny, but they are not caused by each other (they are caused by
good weather, with lots of people at the beach, both eating ice cream
and having a swim in the sea)

Events are connected by a **common cause**: confounding

# R.A. Fisher vs. Richard Doll

R.A. Fisher (geneticist and statistician) was a fervent smoker.
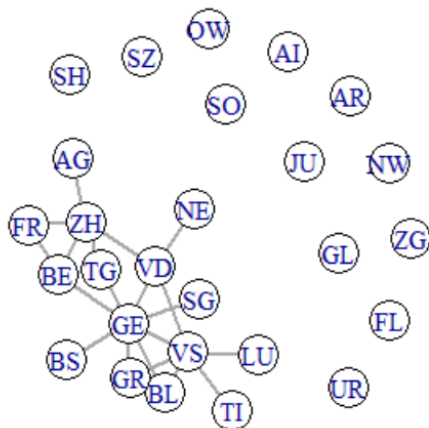


"Smoking and lung cancer are confounded"    "Control for *all* possible confounders"

Sir Richard Doll conducted:

- 1950. Lung cancer study in 20 London hospitals.
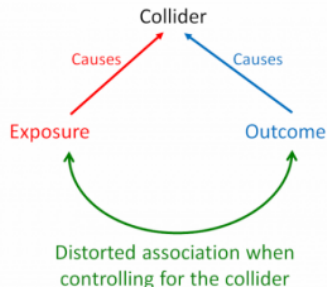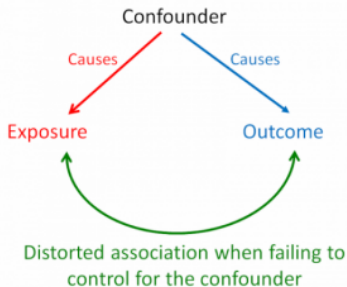- 1954–2001 British Doctor Study to eliminate confounders.

**COLLIDER RULE:** If

- **conditional** dependence between A-C and B-C and A-B
- no dependence between A-B

Then

- A and B cause C.

Applying PC algorithm to Covid network in CH, we find:



directed causal graph

# BONUS

## Did mortality rates improve during Covid-19 pandemic?

Ratio of deaths over cases should tell us something about death rate...
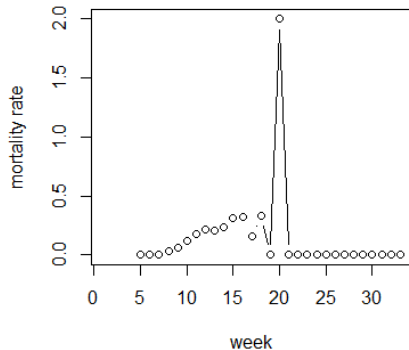
# Did mortality rate improve?

$$\text{mortality}_0(t) = \frac{\text{deaths}(t)}{\text{cases}(t)}$$

# Did mortality rate improve?

$$\text{mortality}_0(t) = \frac{\text{deaths}(t)}{\text{cases}(t)}$$



No delay

# Did mortality rate improve?

$$\text{mortality}_0(t) = \frac{\text{deaths}(t)}{\text{cases}(t)} \qquad \text{mortality}_2(t) = \frac{\text{deaths}(t)}{\text{cases}(t-2)}$$
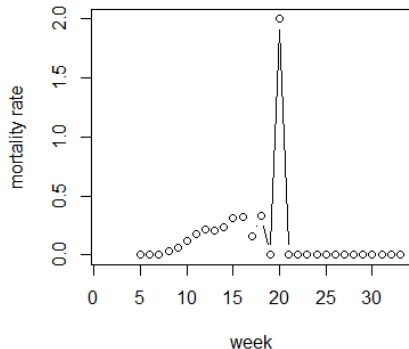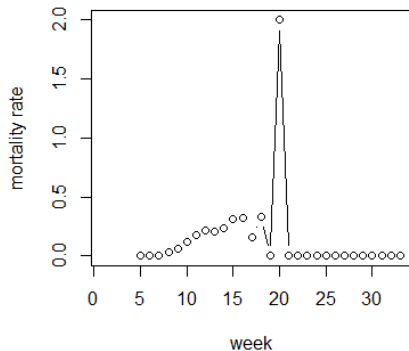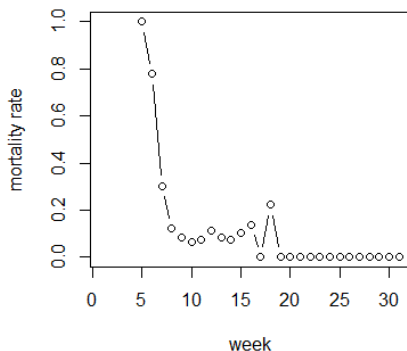


**No delay**

# Did mortality rate improve?

$$\text{mortality}_0(t) = \frac{\text{deaths}(t)}{\text{cases}(t)}$$

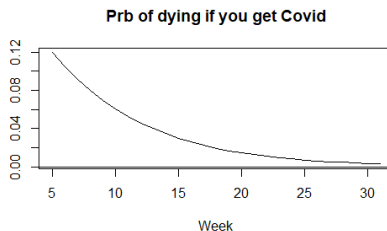$$\text{mortality}_2(t) = \frac{\text{deaths}(t)}{\text{cases}(t-2)}$$

# Logistic regression model

We fit mortality regression with time:

$$\text{odds of dying of covid-19 (in week } t) = e^{\alpha + \beta t}$$

|             | Estimate | Std. Error | z value | Pr(>|z|) |
|------------:|---------:|-----------:|--------:|---------:|
| (Intercept) | -1.262   | 0.073      | -17.35  | 0.0000   |
| week        | -0.147   | 0.007      | -21.51  | 0.0000   |



Prb of dying if you get Covid

Week

Probability of dying of Covid-19 reduced approx 14% each week:

$$e^{-0.147} = 0.86$$

# Conclusions

- **Never take data for granted!**
  (There may be all kind of errors!)
- **Value of data lies entirely in its collection!**
  (Randomized designs are more valuable than observational ones.)
- **Modelling needs to capture underlying process, but also design!**
  (The weaker the design, the stronger the modelling needs to be...)