

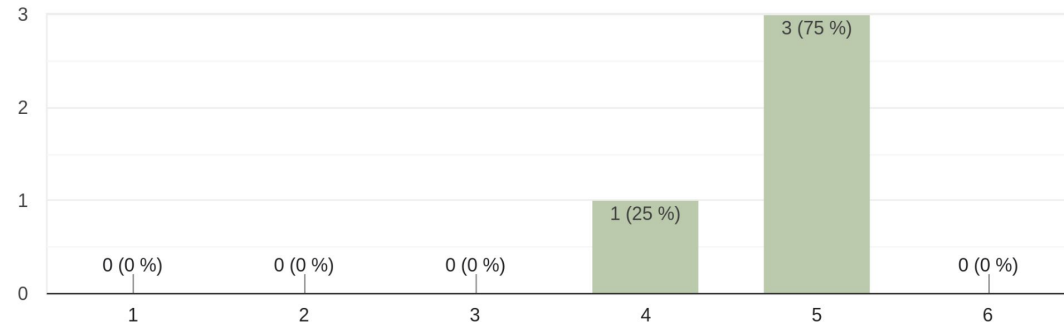
# Feedback of lecture 6

...in lecture 6 I felt run over. I did not get the golden thread. At first we talked about poisson and exp distributions and suddenly we had a look on several ML methods until we end up with the correlations. I know those topics are related to each other since we are dealing with statistical methods. Maybe one single "representative" ML method would have been enough (at least for me) to understand the statistical nature of data driven models. I think we had QSAR, Random Forest and Neuronal Networks.

Thanks so much for your effort and I hope this feedback helps!

How was your overall impression of the sixth lecture?

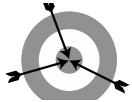
4 Antworten



# AMIDD 2024 Lecture 7: causal inference



Strength



Consistency



Specificity



Temporality



Dose response



Plausibility



Coherence



Experiment



Analogy



Reversibility

Adapted from “The Environment and Disease: Association or Causation?” (1965) by A. B. Hill.

*Dr. Jitao David Zhang, Computational Biologist*

*<sup>1</sup> Pharmaceutical Sciences, Pharma Research and Early Development, Roche Innovation Center Basel, F. Hoffmann-La Roche*

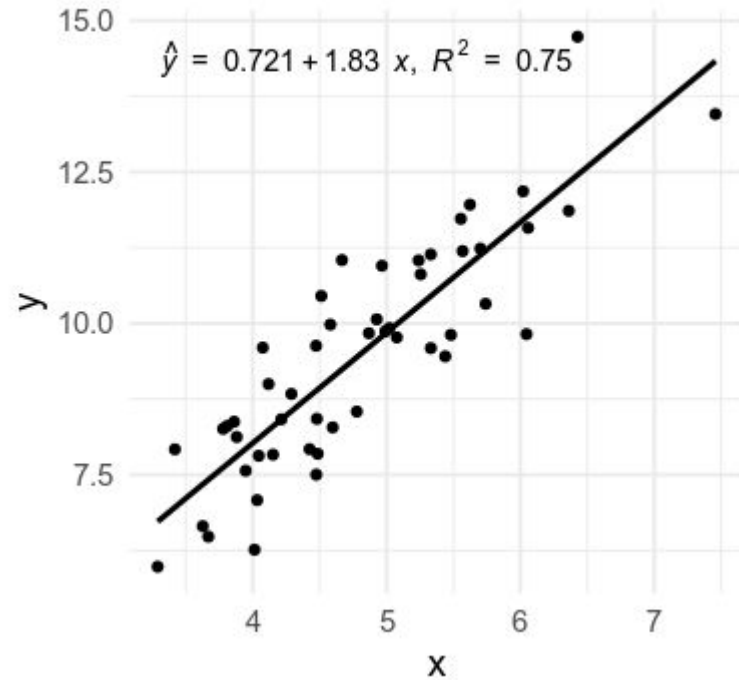
*<sup>2</sup> Department of Mathematics and Informatics, University of Basel*

# The simplest linear model has three components: the intercept, the slope, and a measure of fit

In this example, the coefficient of determination ( $R^2$ ) is used as the measure.

$R^2$  measures the relative fit of the linear model with regard to a baseline model, where the mean value of  $y$  is used as a fit.

	x	y
1	4.926791	10.067779
2	4.479734	8.424283
3	4.289686	8.835629
4	4.474023	9.630499
5	4.214551	8.416680
6	6.057431	11.578080
7	4.597903	8.283025
8	5.021571	9.922731
9	3.627323	6.651222
10	5.622794	11.959972
11	5.555025	11.727815
12	4.966007	10.951562
13	5.076791	9.768299



# Generative models shed light on correlation and causality



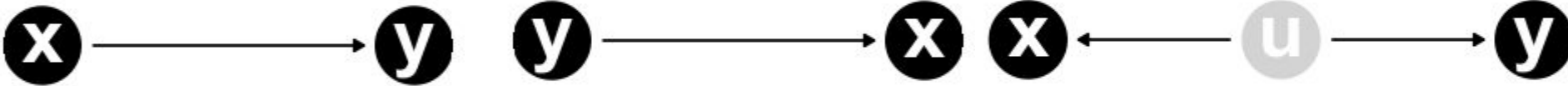
	x	y
1	4.926791	10.067779
2	4.479734	8.424283
3	4.289686	8.835629
4	4.474023	9.630499
5	4.214551	8.416680
6	6.057431	11.578080
7	4.597903	8.283025
8	5.021571	9.922731
9	3.627323	6.651222
10	5.622794	11.959972
11	5.555025	11.727815
12	4.966007	10.951562
13	5.076791	9.768299



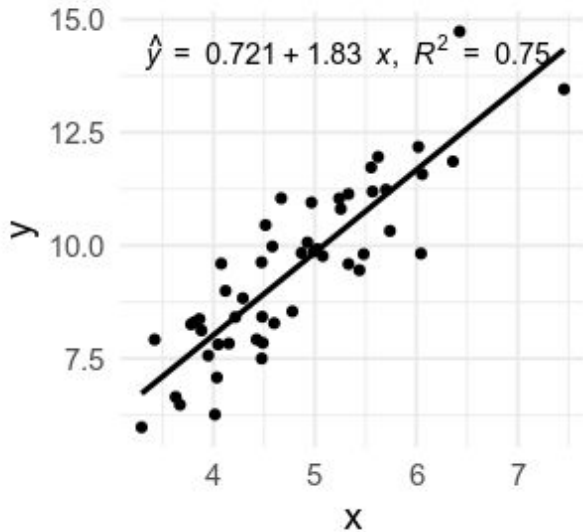
Assumptions of the **generative model**:

1. X is a random variable;
2. Every unit change of X induces a change of 2 units in Y.

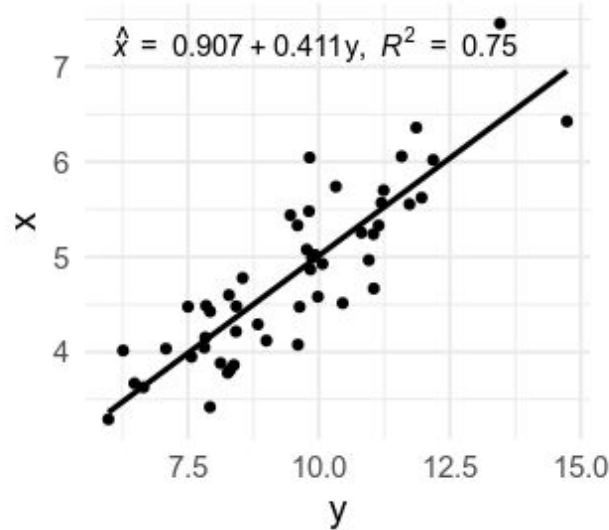
# Correlation may be coincidence, or causation, or confounding (common cause)



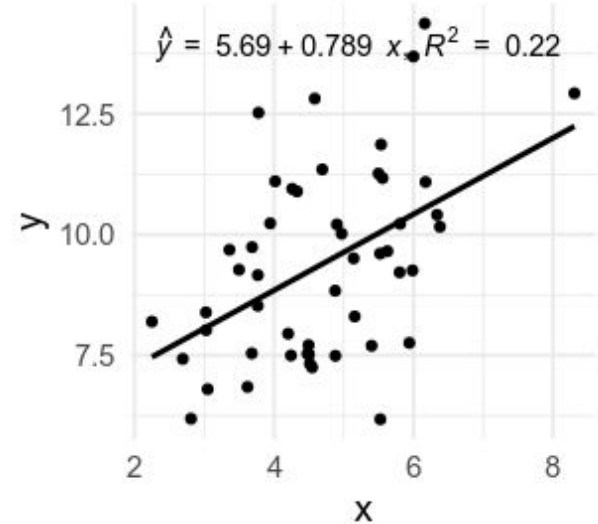
True effect: 2.0



The reverse fit

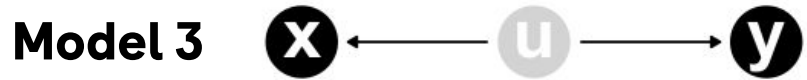


True effect: 0.0



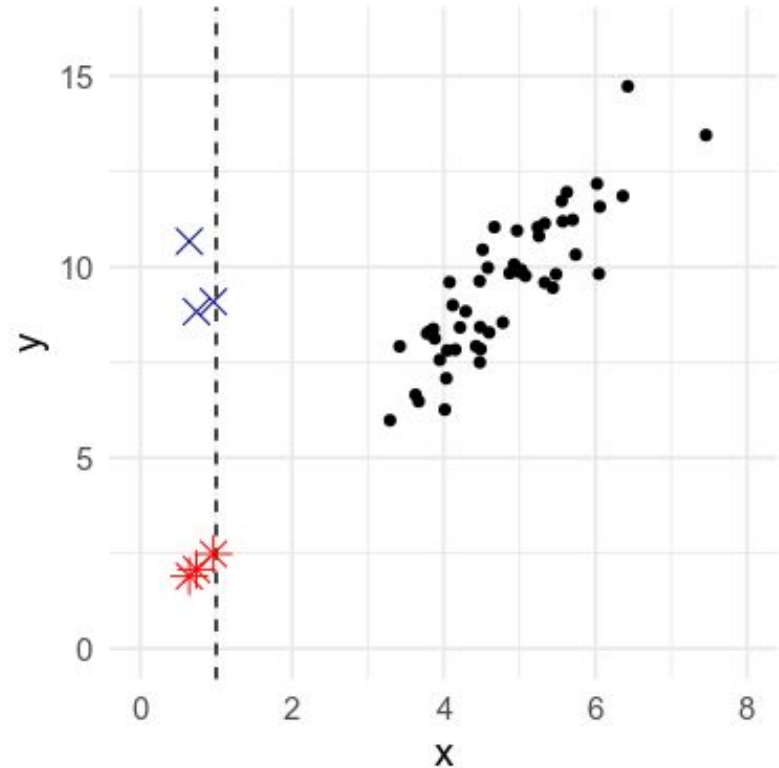
**Statistical models alone cannot derive causality from correlation**

# We learn causality by (1) listing models explicitly and (2) manipulating a variable and observe the outcomes



Assume that the data is generated by either Model 1, or Model 2, or Model 3. And assume that we can manipulate the value of X by setting it to 1.0 (the dash line).

**Question: which outcomes (red stars or blue crosses) would support which models? Why?**



## A very good question

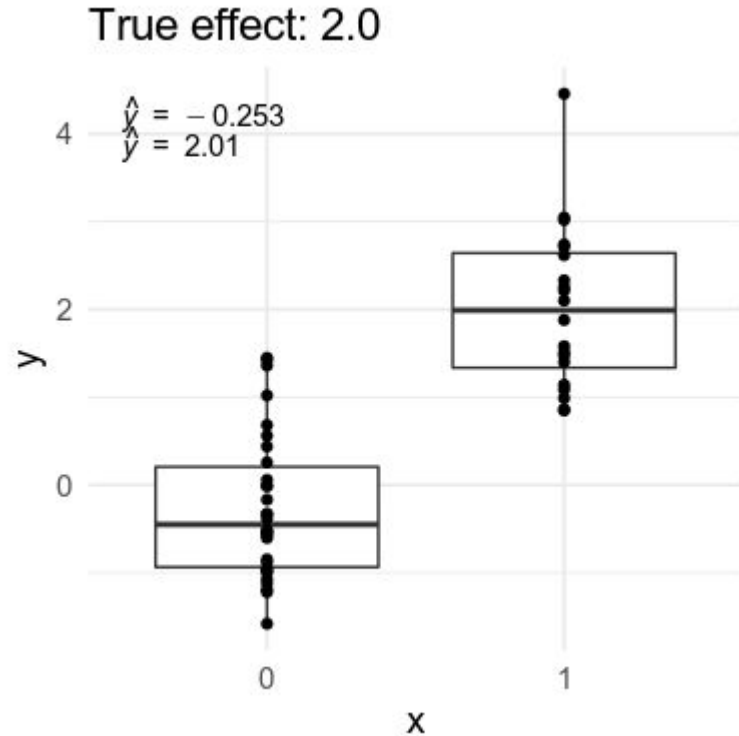
It is still unclear to me why sometimes the blue crosses are predicted so high when we set  $X=1$ . If there is the following relationship " $X \leftarrow U \rightarrow Y$ " wouldn't  $X$  and  $Y$  still be correlated? Let's use the ice cream/temperature/wildfire example where the relationship is "ice cream  $\leftarrow$  temperature  $\rightarrow$  wildfire". If  $X$  is lower,  $Y$  is also lower, even though it's not causal. If we have one wildfire, we can't conclude anything about the temperature at that time. But if we have an overall low number of wildfires ( $Y$ ) in a given time period, couldn't we assume the temperature during this given time period was low? And through correlation assume that the ice cream sales during this given time period was also low? Wouldn't this look like the red stars?

# Variables in models can be either continuous or discrete



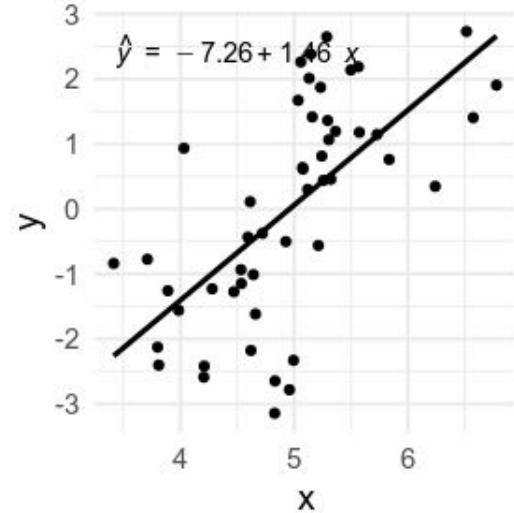
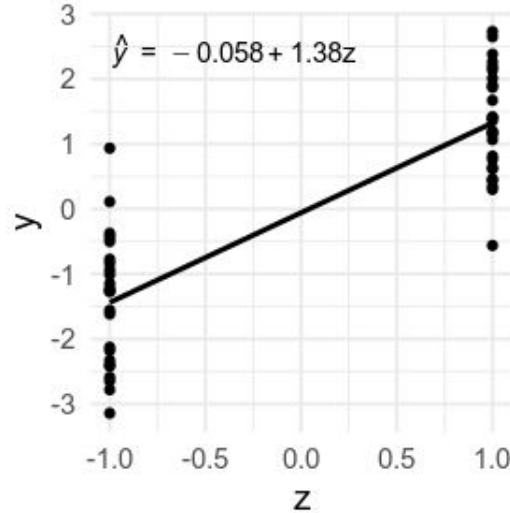
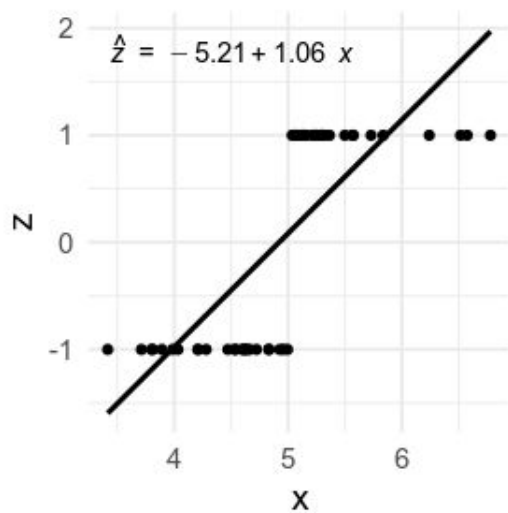
Assumptions of the **generative model**:

1. X is a random variable taking the value of either 0 or 1 with equal probability  $p=0.5$ .
2. Y is a random variable following Gaussian distributions:
  - a.  $\text{Mean}(Y|X=0)=0$
  - b.  $\text{Mean}(Y|X=1)=2.0$





# Common Directed Acyclic Graph (DAG) structures (1): the Pipe



Assumptions of the **generative model**:

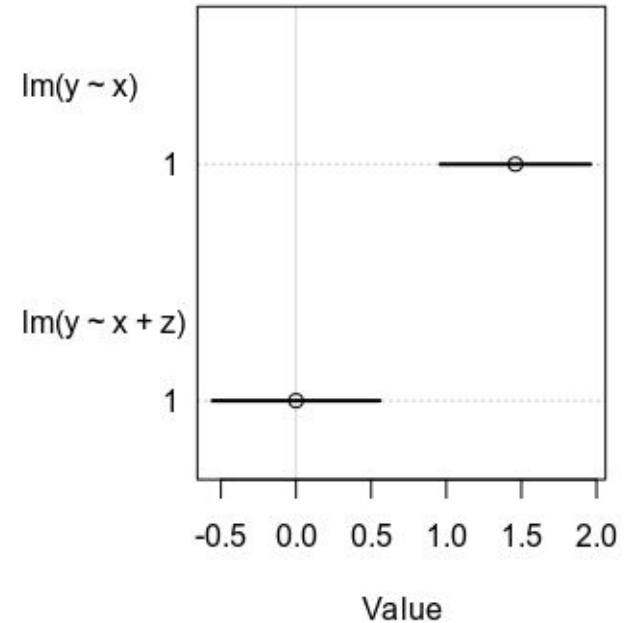
1.  $X$  is a random variable following Gaussian distribution  $N(5,1)$
2.  $Z$  takes the value of -1 if  $X$  is smaller than 5, and 1 if  $X$  is equal to or larger than 5.
3.  $Y$  is a random variable with mean defined by  $Z \cdot 1.5$ .

# Conditional on the mediator in a pipe, the effect of the cause is blocked

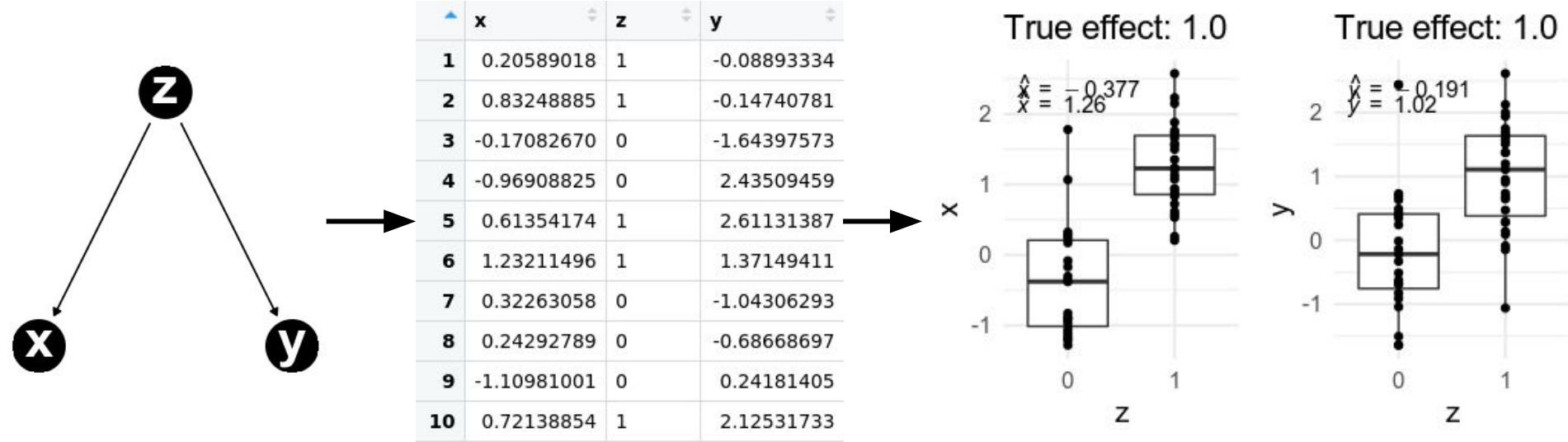


Assumptions of the **generative model**:

1. X is a random variable following Gaussian distribution  $N(5,1)$
2. Z takes the value of -1 if X is smaller than 5, and X is equal to or larger than 5.
3. Y is a random variable with mean defined by  $Z*1$



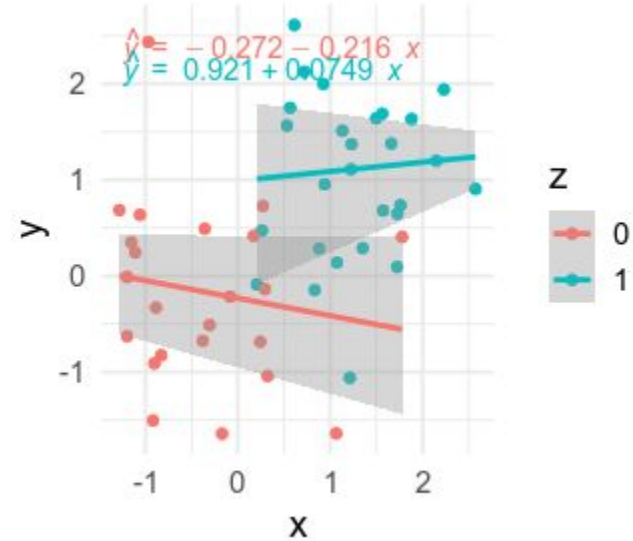
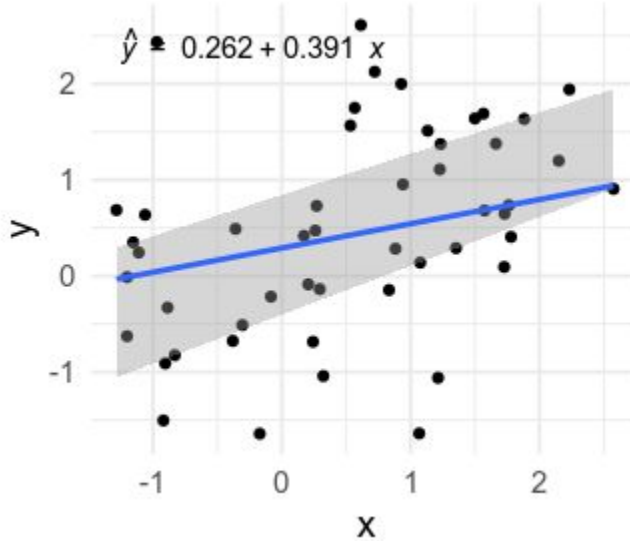
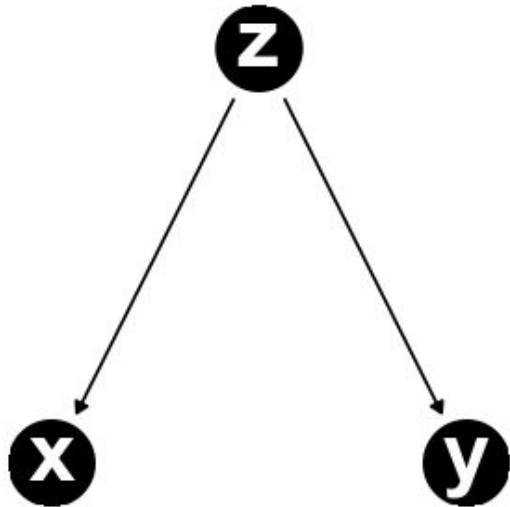
## Common DAG structures (2): The Fork



Assumptions of the **generative model**:

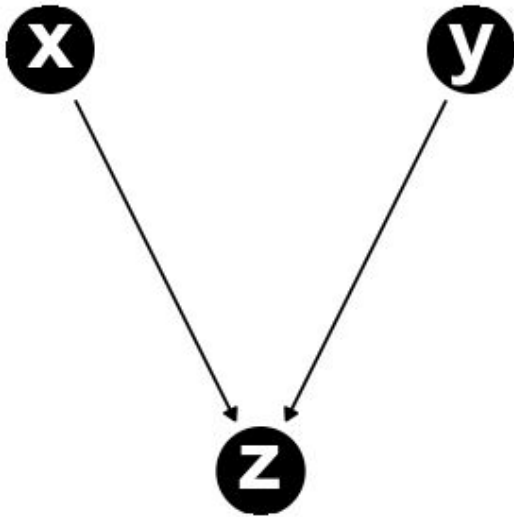
1. Z is a random variable taking the value of either 0 or 1.
2. Both X and Y are random variables following Gaussian distribution with mean equal to Z.

# Conditioning on the fork *breaks* the correlation

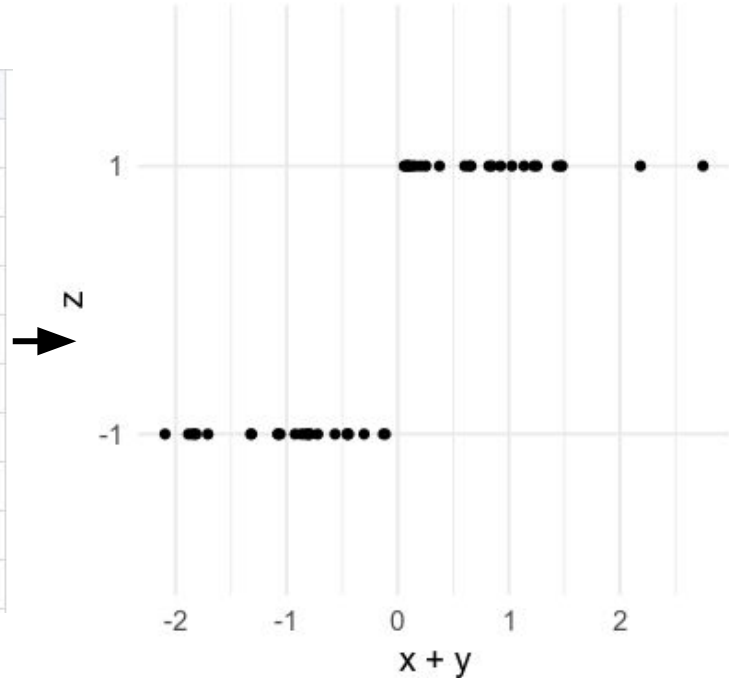


Given a fork structure, both children of the common cause are correlated. The correlation disappears when we condition on the common cause (i.e. stratification by the common cause in the case of discrete variables, or including the variable in the regression in the case of continuous variables).

# Common DAG structures (3): The Collider



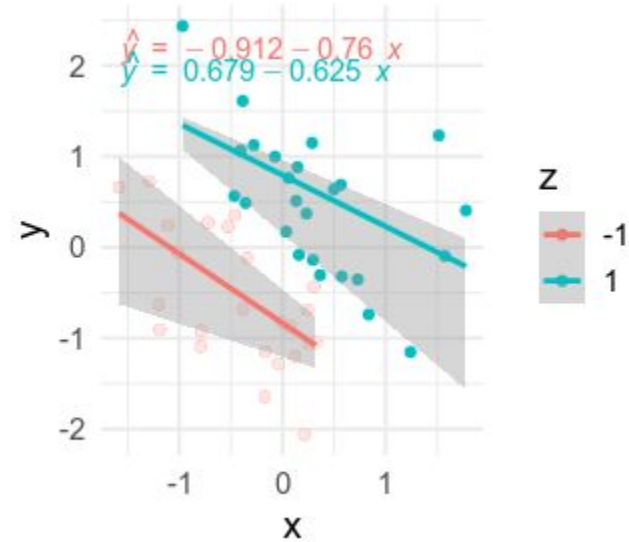
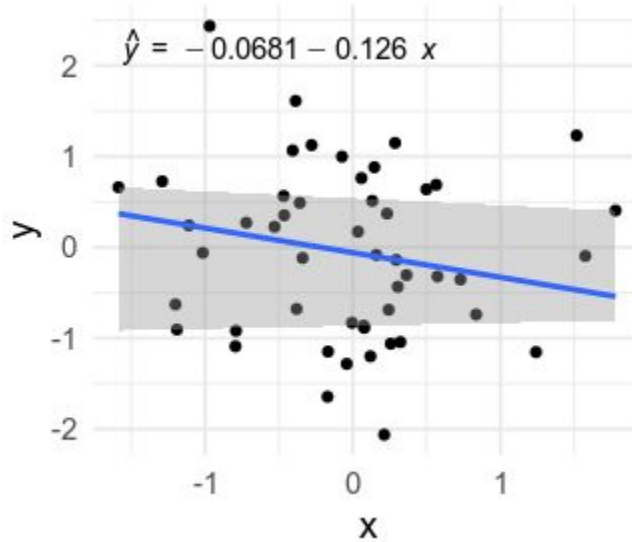
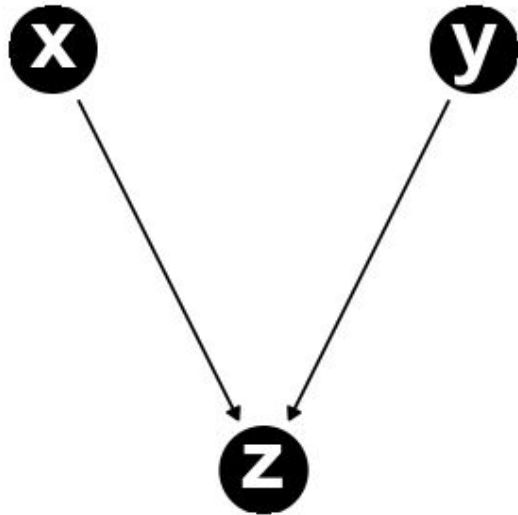
	x	z	y
1	0.835386320	1	-0.73897252
2	-0.005354014	-1	-0.82972315
3	0.058788286	1	0.76213369
4	-1.015602246	-1	-0.05951719
5	-0.339569780	-1	-0.11745910
6	-0.041077979	-1	-1.28243716
7	0.363740407	1	-0.30570762
8	0.119496314	-1	-1.19932461
9	0.257108454	-1	-1.06044066
10	0.304537158	-1	-0.43396492



Assumptions of the **generative model**:

1.  $X$  and  $Y$  are random variables following Gaussian distribution  $N(0,1)$
2. The value of  $Z$  is 1 if  $X+Y>0$ , and -1 if  $X-Y\leq 0$ .

# Conditioning on the collider introduces *spurious correlations*

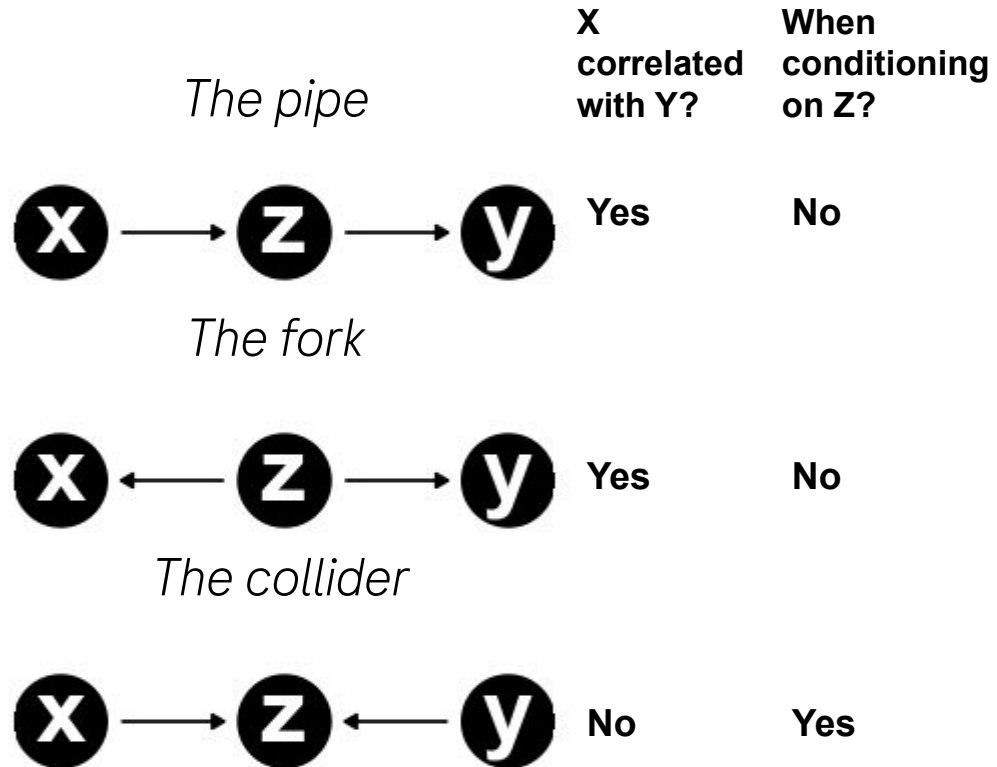


In a collider structure, the parents of the collider can be independent from each other. However, they become correlated when we condition on the collider.

**Collider is everywhere!**

## A summary so far

- Data alone cannot tell causality, though in most cases we are interested in causal questions.
- Correlation between two variables can be caused by coincidence, causality, or common cause.
- Most common structures in a graph causal model are pipes, forks, and colliders. Stratifying by or regressing out variables may **remove** or **create** correlation.

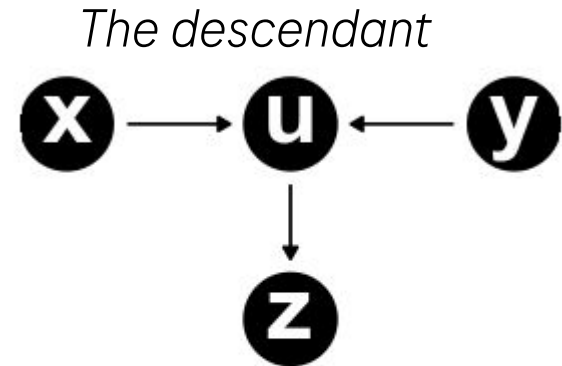


# Stop exploitative data analysis, build generative models



Biomarker, tox study, pathology, omics data, real-world data, EHR, ...

	x	z	y
1	0.835386320	1	-0.73897252
2	-0.005354014	-1	-0.82972315
3	0.058788286	1	0.76213369
4	-1.015602246	-1	-0.05951719
5	-0.339569780	-1	-0.11745910
6	-0.041077979	-1	-1.28243716
7	0.363740407	1	-0.30570762
8	0.119496314	-1	-1.19932461
9	0.257108454	-1	-1.06044066
10	0.304537158	-1	-0.43396492

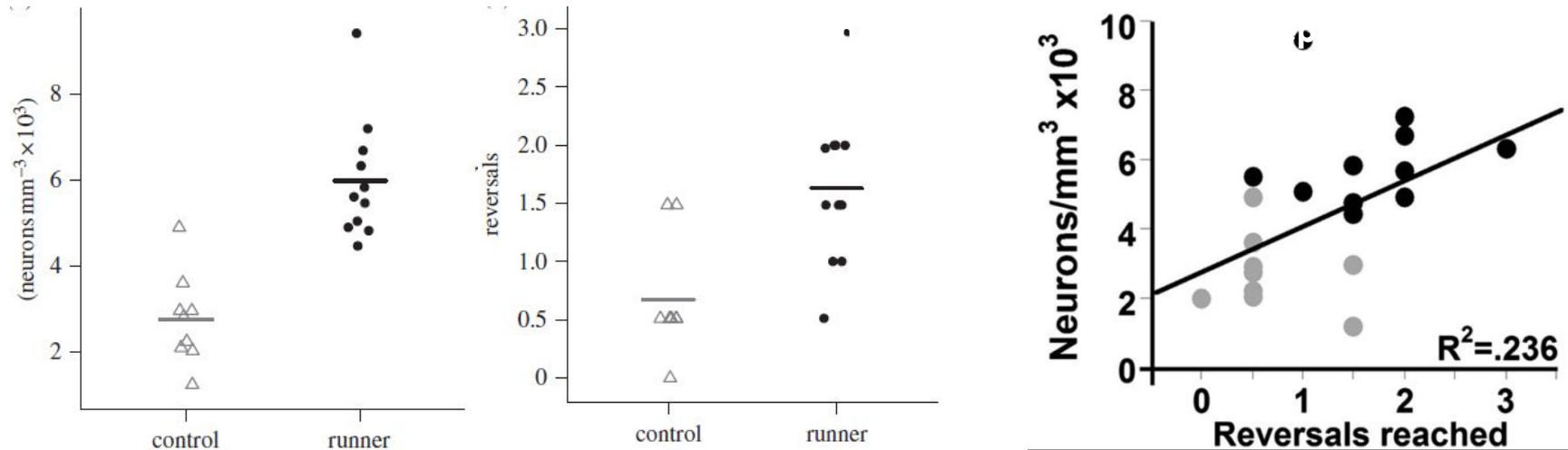


**We need to build models (knowledge + assumptions) to infer causality**



# Claim: running enhances spatial pattern separation in mice

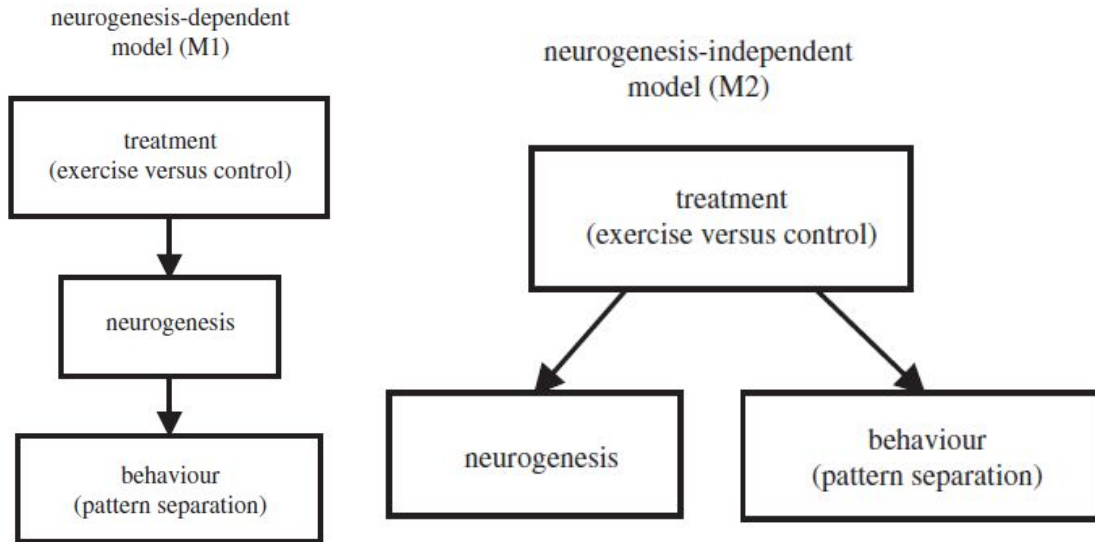
Creer et al., PNAS 2010



Creer, David J., Carola Romberg, Lisa M. Saksida, Henriette van Praag, and Timothy J. Bussey. "Running Enhances Spatial Pattern Separation in Mice." *Proceedings of the National Academy of Sciences* 107, no. 5 (February 2, 2010): 2367–72. <https://doi.org/10.1073/pnas.0911725107>.

Lazic Stanley E. "Using Causal Models to Distinguish between Neurogenesis-Dependent and -Independent Effects on Behaviour." *Journal of The Royal Society Interface* 9, no. 70 (May 7, 2012): 907–17. <https://doi.org/10.1098/rsif.2011.0510>.

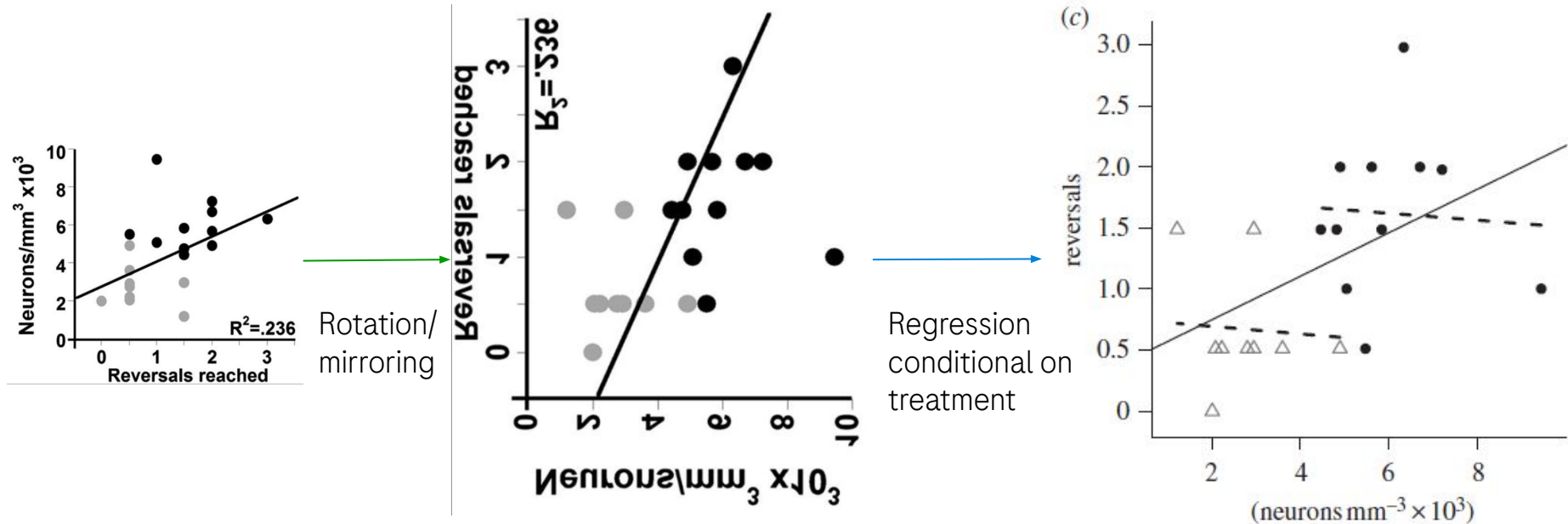
# Question: does pharmaceutical modulation of neurogenesis benefit pattern separation?



M1 (the **pipe** model) suggests that conditioned on neurogenesis, exercise and behaviour are independent (not correlated).

M2 (the **fork** model) suggests that conditioned on exercise, neurogenesis and behaviour are independent.

# Behaviour and neurogenesis even shows *negative correlation conditional on exercise*- an example of Simpson's Paradox

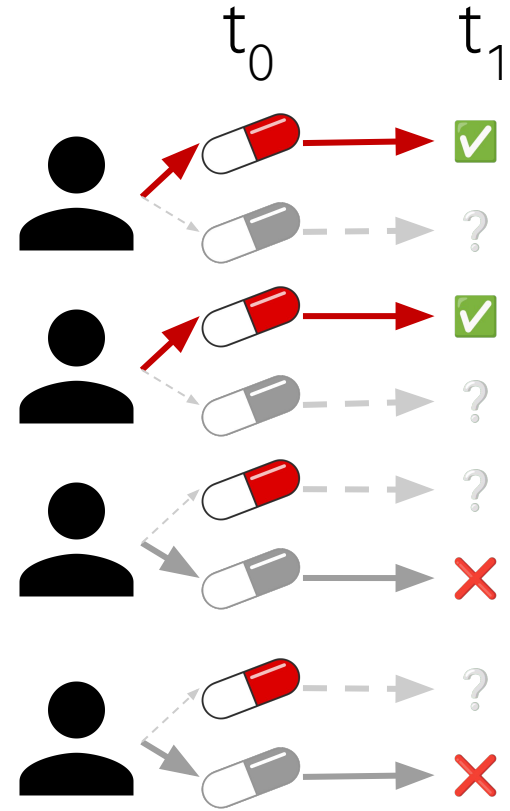


Based on the analysis, I believe model M2 is more likely to be true than M1.

**Causal inference reduces bias in analysis by listing models explicitly**

# Causal inference is important for both randomized experiments and observational studies

- In drug discovery and development, we often care about **potential outcomes** or **counterfactuals**: what had if the patient received the alternative treatment, keeping everything else constant?
- **Randomized experiments** and **controlled trials** are gold-standard methods to address causal questions. Non-compliance and intermittent events call for causal analysis of the data even in randomized trials.
- Given causal models, it is sometimes possible to learn causal relationships from observational data as well.



# Causal inference is a missing data problem

Individual	Treatment	Value (AU)
1	Control	75
2	Control	73
3	Control	74
4	Treatment	55
5	Treatment	45
6	Treatment	60

A classical textbook

Individual	Value (AU) with Control	Value (AU) with Treatment
1	75	?
2	73	?
3	74	?
4	?	55
5	?	45
6	?	60

A classical textbook in 50 years

# Assignment mechanism determines which data are missing, and determine the statistical technique to be used

Individual	Value (AU) with Control	Value (AU) with Treatment
1	75	?
2	73	?
3	74	?
4	?	55
5	?	45
6	?	60

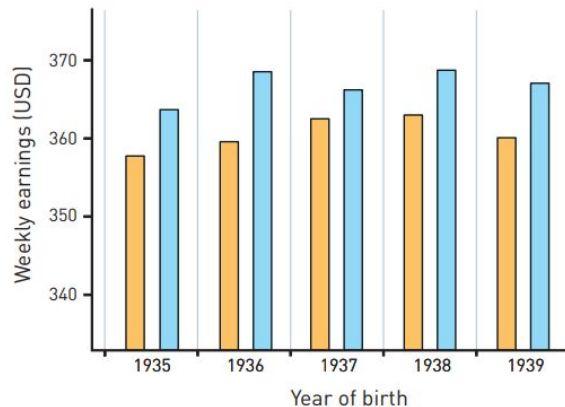
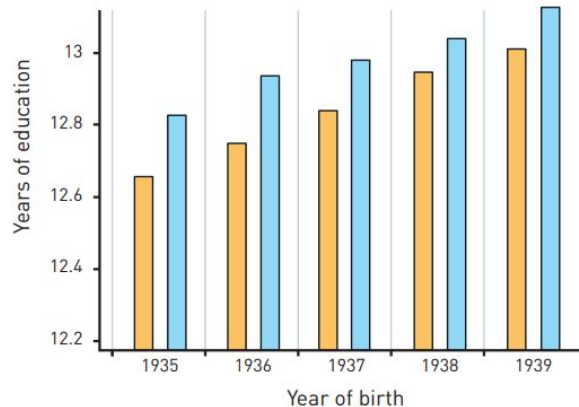
- **Classical Randomized Experiments:** we control the assignment mechanism;
- **Regular assignment (*observational studies*):** we know part but not all of the assignment mechanism;
- **Regular assignment with non-compliance:** we need an *instrumental variable*.

# Instrumental variable helps to dissect causal from confounding effects

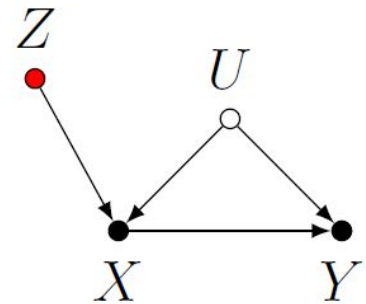
People born late in the year have more years of education and higher incomes

Additional years of education have a positive effect on income. The figure uses data from Angrist and Krueger (1991).

■ Born in first quarter   
 ■ Born in fourth quarter

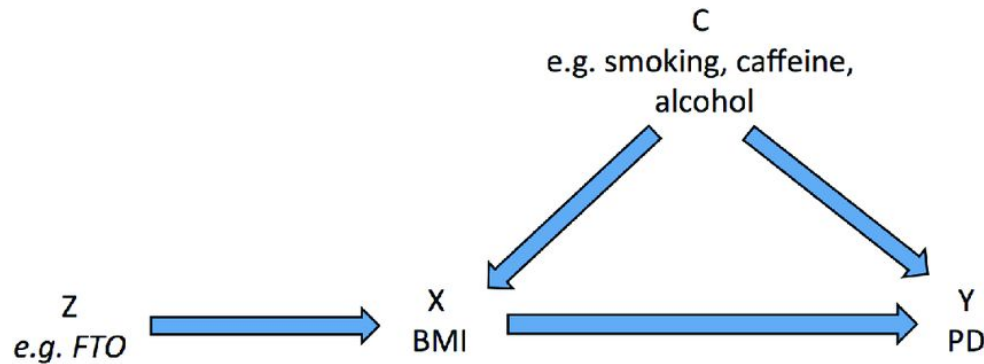


<https://www.nobelprize.org/prizes/economic-sciences/2021/popular-information/>

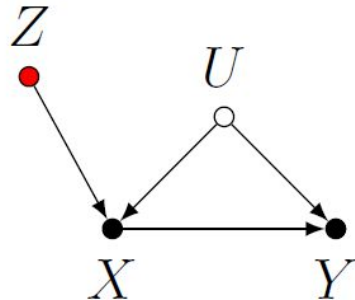


Z: Birthday  
 X: Education  
 Y: Income  
 U: Socioeconomic and individual factors

# Instrumental variable is critical for both Mendelian Randomization (MR) and handling with non-compliance



Noyce, Alastair J., Demis A. Kia, Gibran Hemani, Aude Nicolas, T. Ryan Price, Eduardo De Pablo-Fernandez, Philip C. Haycock, et al. "Estimating the Causal Influence of Body Mass Index on Risk of Parkinson Disease: A Mendelian Randomisation Study." *PLoS Medicine* 14, no. 6 (June 2017): e1002314. <https://doi.org/10.1371/journal.pmed.1002314>.



Z: Assignment (Placebo/Drug)  
 X: Treatment  
 Y: Value of interest  
 U: unobserved factors

Cinelli, Carlos, Andrew Forney, and Judea Pearl. "A Crash Course in Good and Bad Controls." SSRN Scholarly Paper. Rochester, NY: Social Science Research Network, September 9, 2020. <https://doi.org/10.2139/ssrn.3689437>.



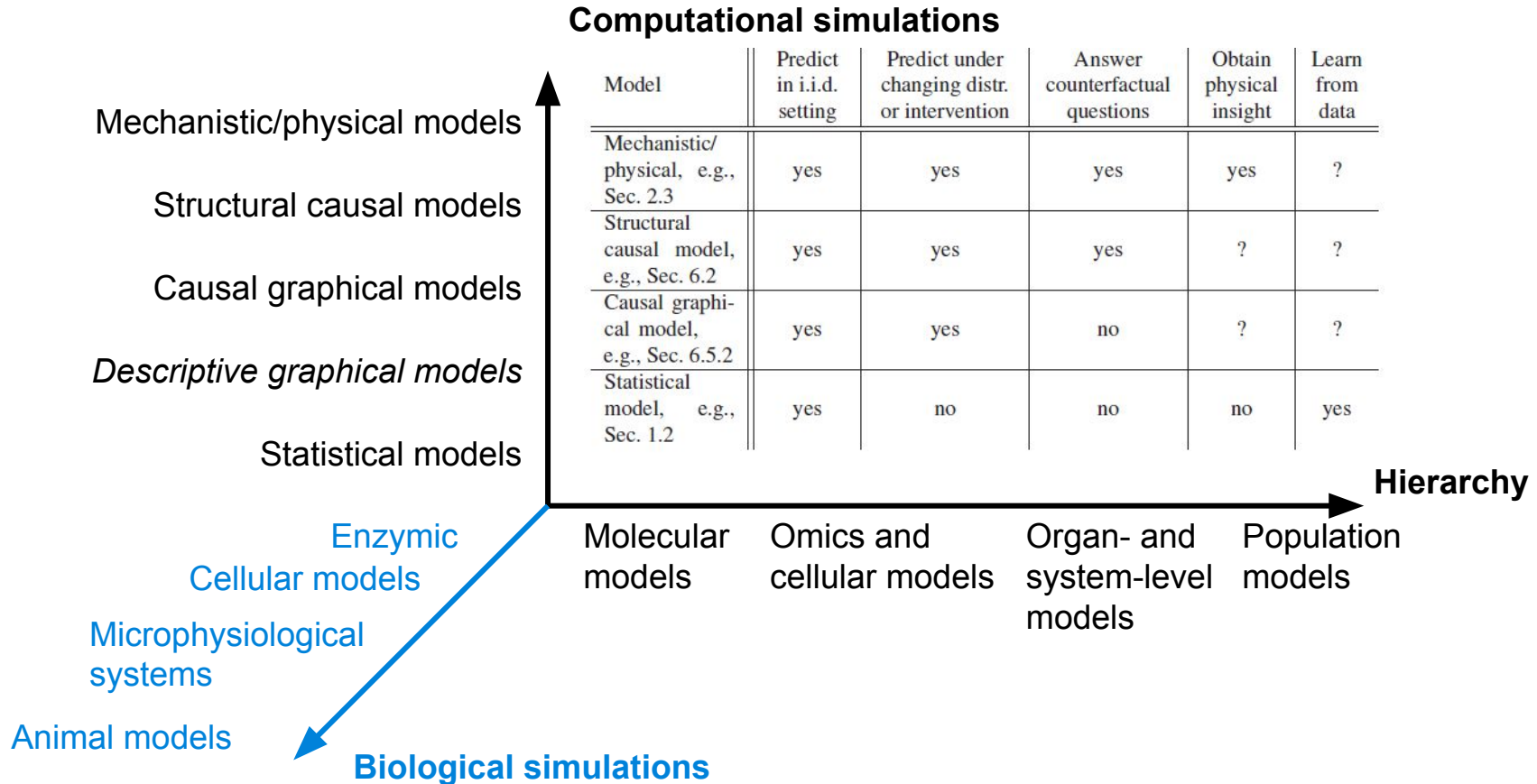
# Consequences

1. **Data alone does not answer causal questions:** whenever we are interested in interventions (modulating a target, changing the structure of a molecule, *etc.*), predictive tools such as linear regression, machine learning, and artificial intelligence models must be embedded in the causal framework.
2. **Addressing causal questions when necessary:**
  - a. Derive causal models using science, making assumptions transparent
  - b. Program the model as a generative simulation
  - c. Design research and validate statistical analysis using (b)
  - d. Confront the model with data, share both wins and losses transparently with others
  - e. Revise and repeat
3. **Model first, data second:** From **DA** (**D**ata and **A**nalytics) to **MADAM** (**M**odel construction, **A**nalysis of the model, **D**ata collection, **A**nalysis of the data with the model, and **M**odel refinement)

## Ten Simple Rules of Causal Inference

1. Clarify whether correlation or causation is of interest.
2. Draw models of knowledge, assumptions, and data-generation processes as graphs.
3. Formulate the causal question by identifying the target estimand.
4. Collect, check the quality of, and filter data.
5. Estimate the causal effect with software.
6. Challenge and refute the causal model.
7. Compare results with estimates from alternative methods.
8. Share model, data, and analysis.
9. Design, perform, and analyze new experiments.
10. Apply learnings from causal inference in the real world.

# Models in disease understanding and drug discovery

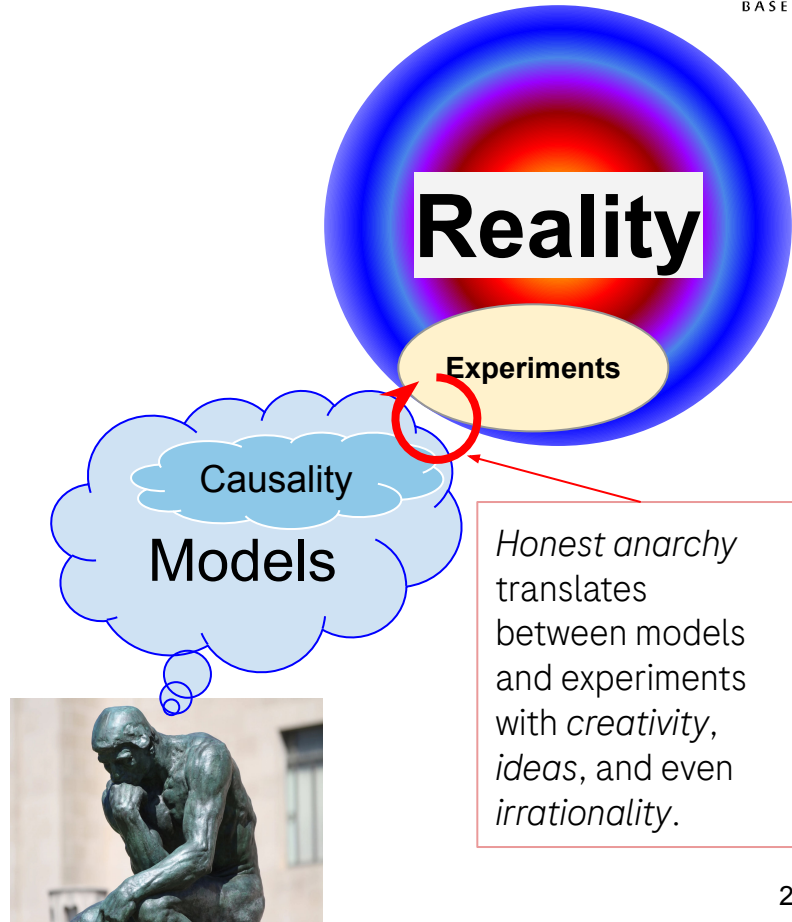


# Causal inference is not always easy



Credit: Ulrich Certa

*Science is not about reality.  
Science is always about models.*



*“There is no method for making  
causal models other than **science**.  
There is no method to science other  
than **honest anarchy**.”*

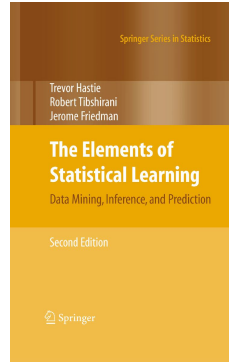
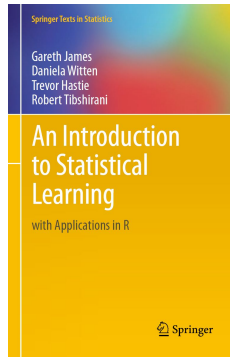
- Richard McElreath

# Conclusions

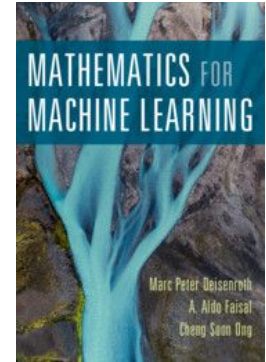
1. Statistical and machine learning models can model linear and nonlinear relationships between variables.
2. Applying statistical and ML models in drug discovery needs to consider the facts that we always work on something new, structure similarity does not warrant activity similarity, and correlation is not causation.
3. Causal inference combines prior knowledge and statistical/ML modelling to answer *what-if* questions.

# Resources for learning about machine learning

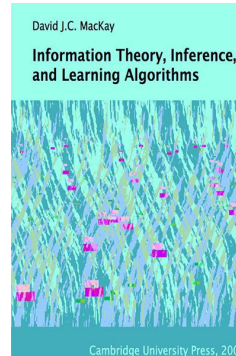
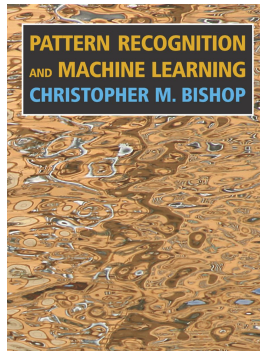
ESL and ISL: From a frequentist view (almost)



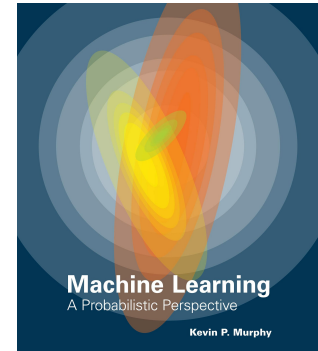
Mathematical foundations



PRML and ITILA: From a Bayesian view



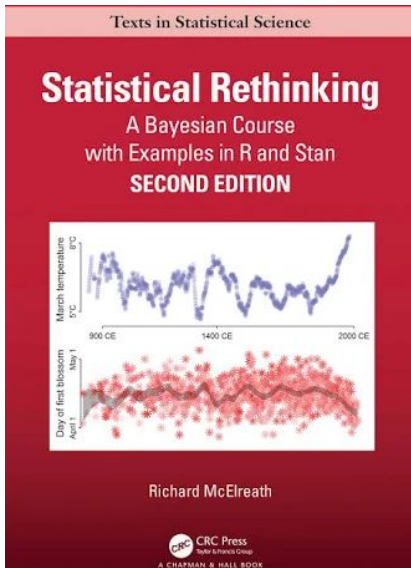
MLaPP: Application oriented, more accessible, and balanced views



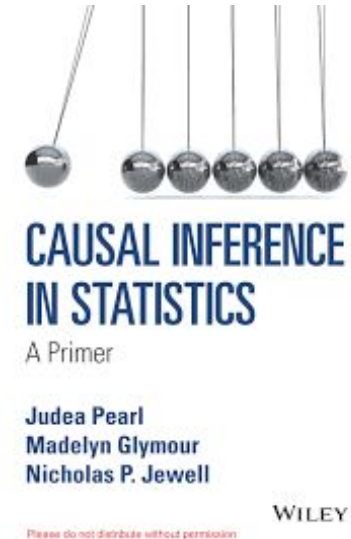
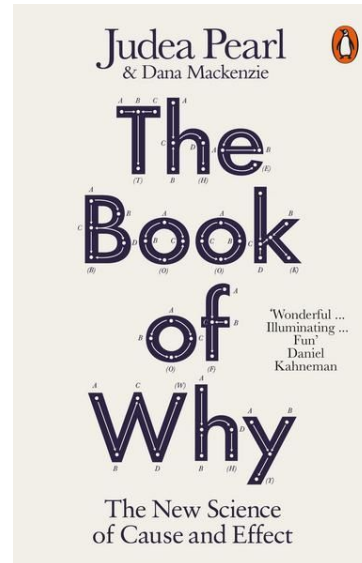


# Resources for learning about causal inference

[Causal inference in drug discovery and development](#), Michael and Zhang, 2022



[Lectures available on YouTube](#)



# Answers

Red stars are supported by Model 1.

Blue crosses are supported by both Model 2 and Model 3.

Reason: causality ( $C \rightarrow E$ , from cause to effect) is directional. Manipulating  $C$  has an effect on  $E$ , while manipulating  $E$  has no effect on  $C$ . Blue crosses are around mean values of  $Y$ . If  $Y$  causes  $X$ , manipulating  $X$  has no effect on  $Y$ . Then the most likely values of  $Y$  will be around the mean of existing samples.