# Machine learning for causality?

@kordinglab

# The bitter Lesson (Sutton, 2019)

- Clever human solutions are eventually beat by general purpose machine learning
  - Chess
  - Go
  - Speech recognition
  - Image recognition
  - Natural Language Processing

## The Bitter Lesson

### Rich Sutton

**March 13, 2019**

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of computation. Most AI research has been conducted as if the computation available to the agent were constant (in which case leveraging human knowledge would be one of the only ways to improve performance) but, over a slightly longer time than a typical research project, massively more computation inevitably

# How good we are not after ML but after causality

- Medicine is almost always after causality
- And ML merely models correlations

# 3 Schools

- ## RCT or bust
    - Randomize, avoid threads
    - Very strong strategy for large audience, high value
    - Weak for underserved populations
    - Gold standard (also C. elegans)

- ## Observational
    - Fit complex models, Correct for various variables
    - Lack of causal validity

- ## Quasiexperimental
    - Find a place where the world contains randomization
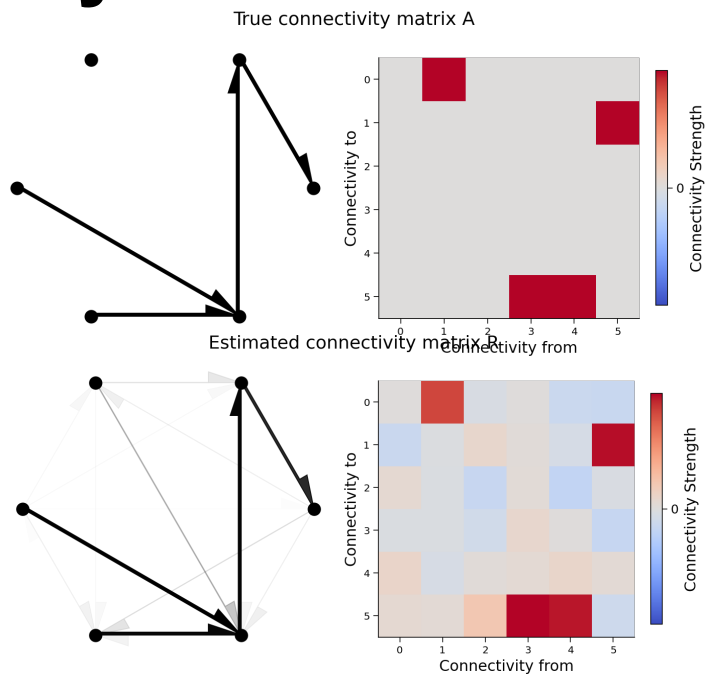    - Usually not possible

# Intuition: Simulate a trivial causal system

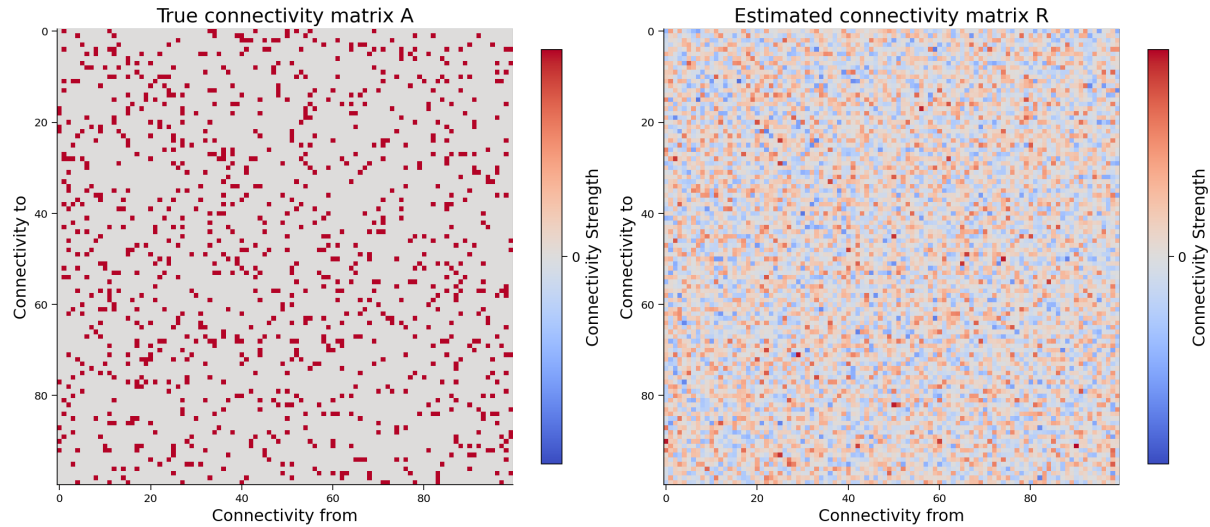$$\vec{x}_{t+1} = \sigma(A\vec{x}_t + \epsilon_t)$$

- $\vec{x}_t$ is an $n$-dimensional vector representing our $n$-neuron system at timestep $t$
- $\sigma$ is a sigmoid nonlinearity
- $A$ is our $n \times n$ *causal ground truth connectivity matrix* (more on this later)
- $\epsilon_t$ is random noise: $\epsilon_t \sim N(\vec{0}, I_n)$
- $\vec{x}_0$ is initialized to $\vec{0}$
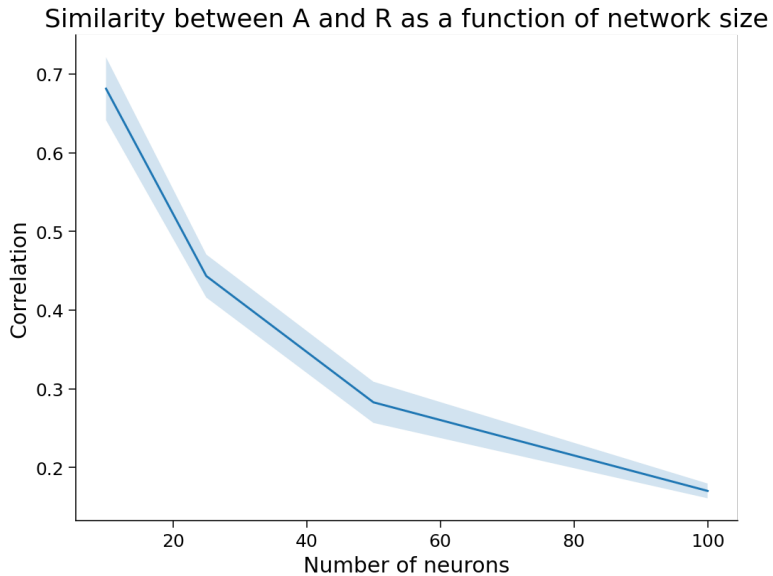
Is correlation ~ causation?

# Great in small system

# Not so great in big system



True connectivity matrix A

Estimated connectivity matrix R

# Delayed Correlation vs Causation



Similarity between A and R as a function of network size

# Omitted Variable Bias

$$y_i = x_i \beta + z_i \delta + u_i$$

$$\widehat{\beta} = (X'X)^{-1} X'(X\beta + Z\delta + U)$$

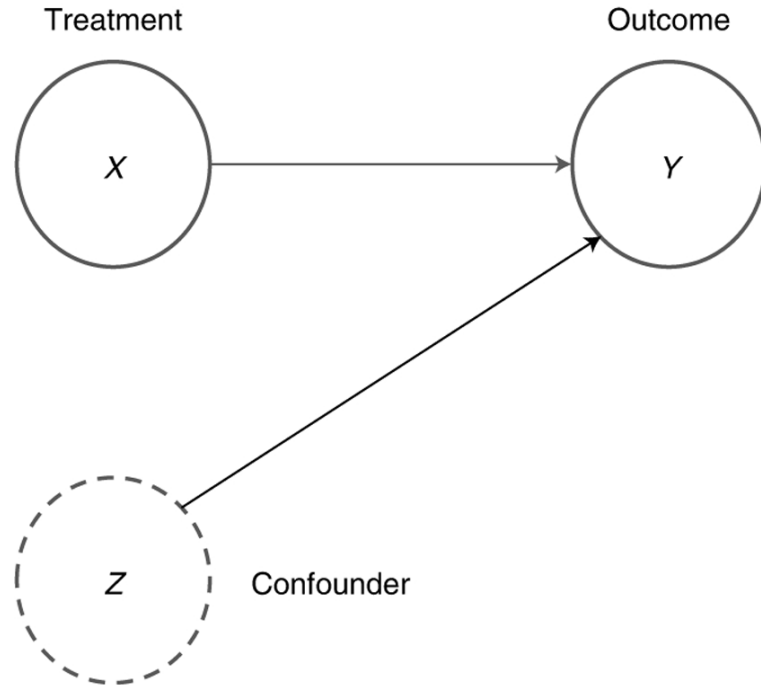$$E[\widehat{\beta} \mid X] = \beta + (X'X)^{-1} E[X'Z \mid X]\delta$$

The bias should be arbitrarily big relative to the signal
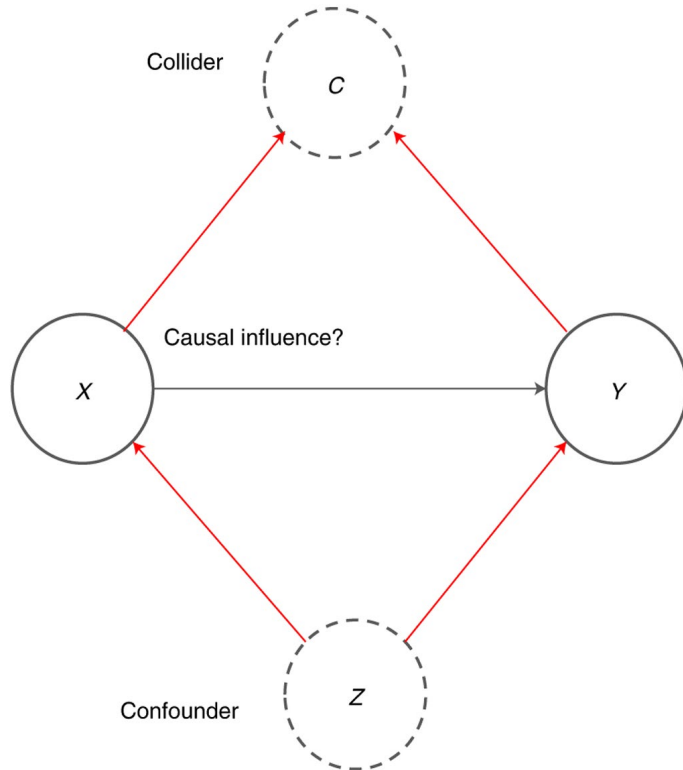This problem does not go away with more data

Measuring and interpreting neuronal correlations

Marlene R. Cohen[1] and Adam Kohn[2]

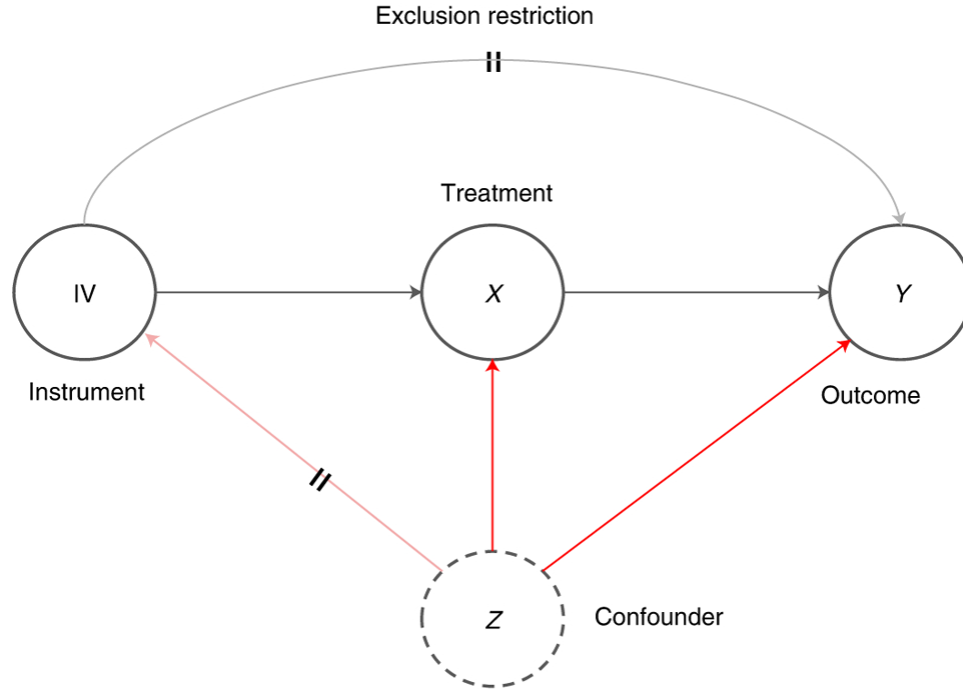# An ideal RCT allows unbiased measurements of causal effects

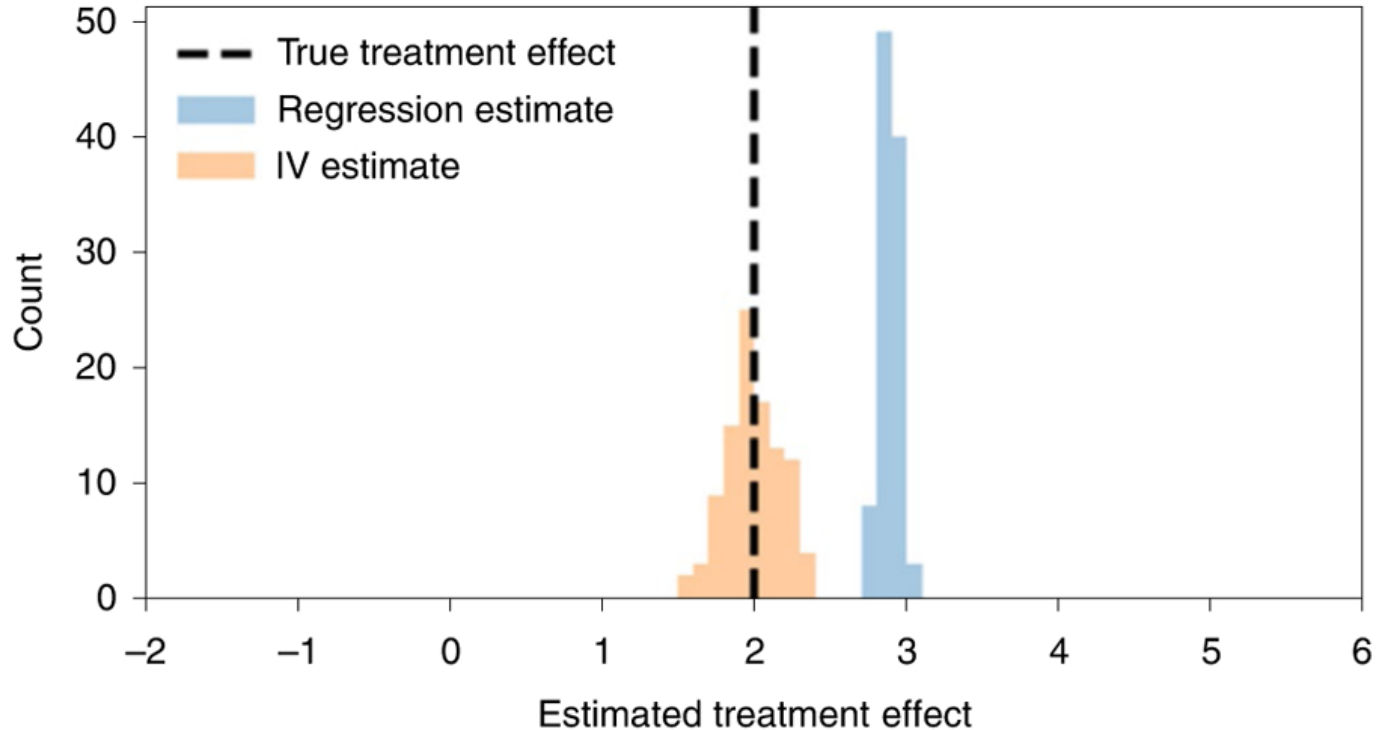# But alas, the world is not perfect: observational studies



Collider

C

Causal influence?

X → Y

Confounder

Z

**Why is it so hard to obtain causality?**

# Instrumental variables



Quantifying causality in data science with quasi-experiments

Tony Liu,[1] Lyle Ungar,[1] and Konrad Kording[2,3,✉]

# If exclusion restriction correct

# If it is violated

# There are lots of quasiexperiments

- Check out Tony's tutorials

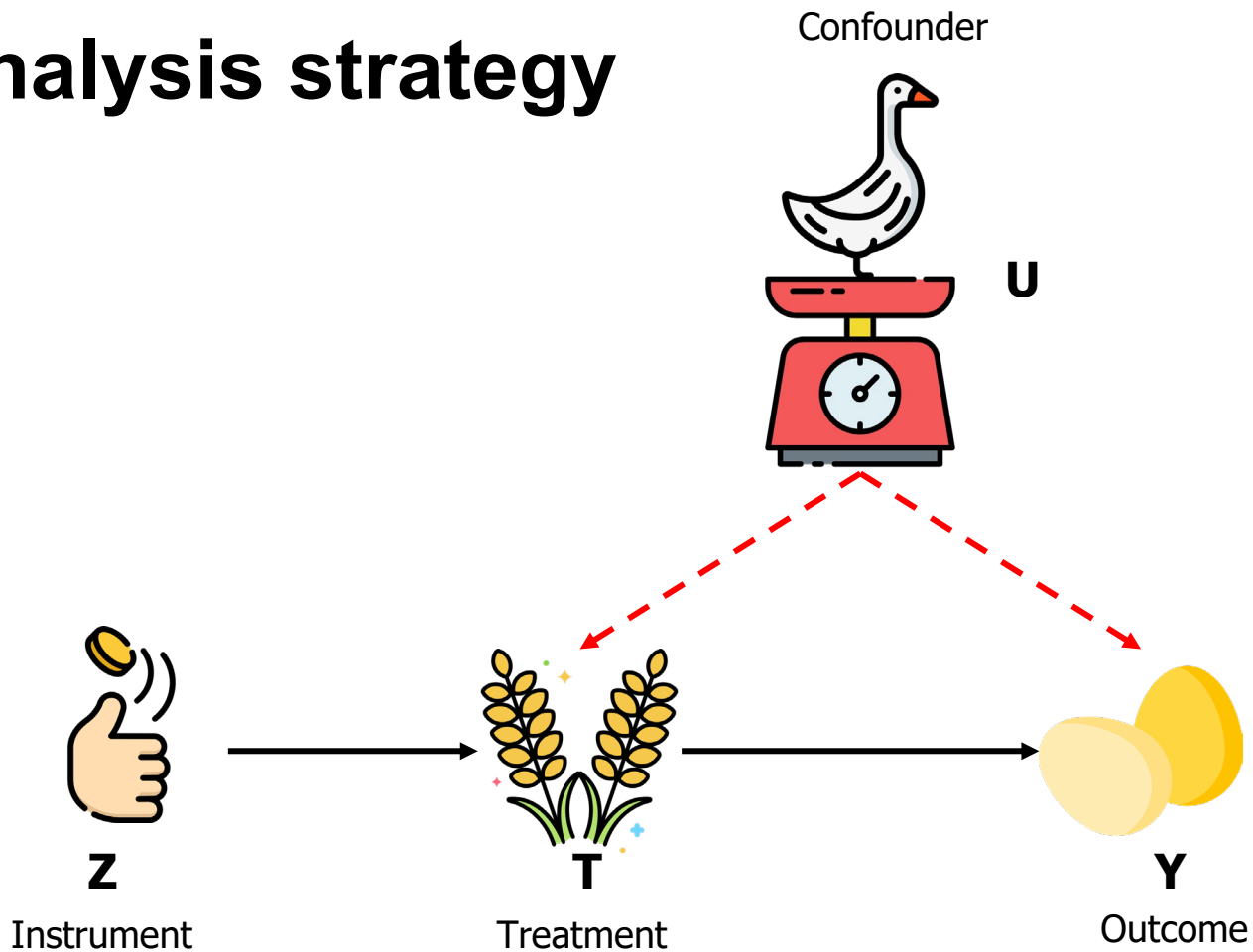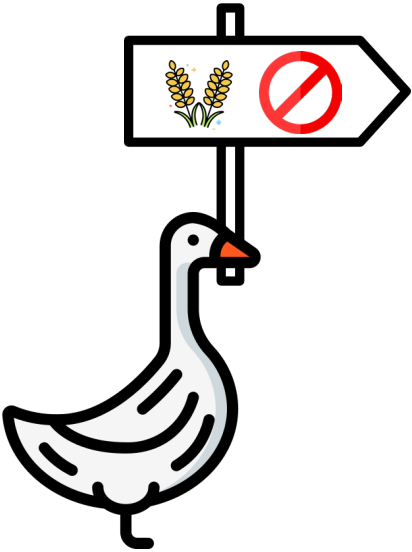https://github.com/tliu526/causal-data-science-perspective

# Does eating grains make ducks lay more eggs?



**T**
Treatment

**Y**
Outcome

# IV analysis strategy
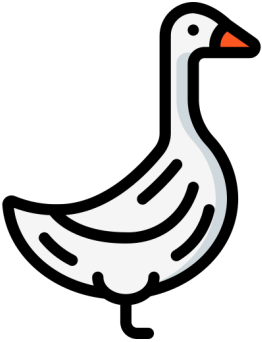


Confounder

U

Z
Instrument

T
Treatment

Y
Outcome

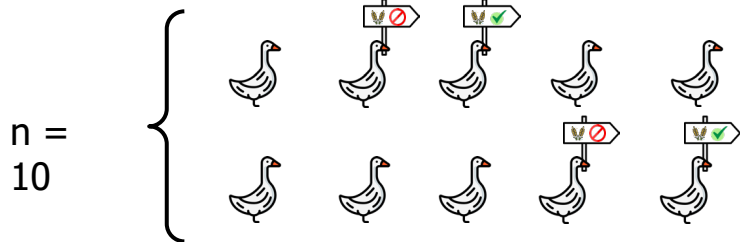# Problem: non-compliance



Never-takers
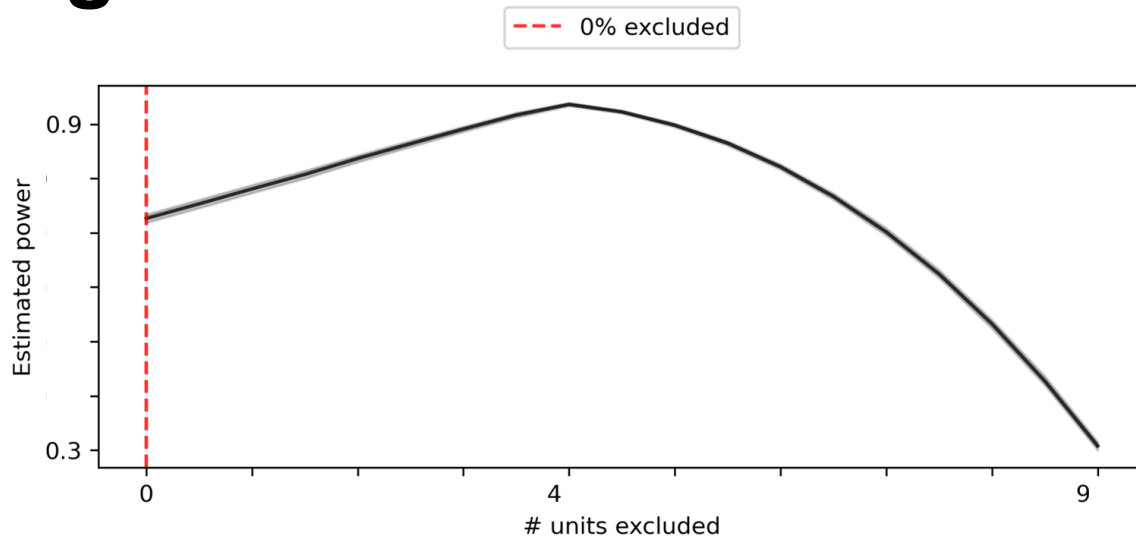
Always-takers

Compliers

# Improving exclusion criteria for IV studies

# Improving exclusion criteria for IV studies



$$N_{\text{eff}} = N \times p^2_{\text{comply}}$$

$n_{\text{eff}} = 3.6$

# Improving exclusion criteria for IV studies



$n_{eff} = 6$

$$N_{\mathrm{eff}} = N \times p^2_{\mathrm{comply}}$$

# Data-driven exclusion criteria

Use observed features and machine learning to predict out-of-sample compliance

# Data-driven exclusion improves power

# Applying this to diabetis

- Optum data
- Above A1C=6.5 you should be diabetic
- But it is fuzzy (hence IV)
- Follow-up A1C should be better

# Let us use regression discontinuity design

**Full sample**

**Data-driven exclusion sample**

30 day diabetes diagnosis rate

0.065 ± 0.006

0.095 ± 0.007

- - - diabetic criteria

A1C %

A1C %

|  | Full sample | Data-driven exclusion sample |

# A nonstandard logic

- ML tells us whom we are talking about
- Our predictions then only apply to these items
- We can then purposefully say nothing about the expected noncompliers
- (which seems unavoidably true)

# Find all the RDDs

- Currently working to
  - Apply this idea to every threshold
  - Continuously sweep to find all thresholds
  - Optum and other datasets

# A bitter lesson approach

## Learning domain-specific causal discovery from time series

**Xinyue Wang**                                          *wsinyue@seas.upenn.edu*
*Department of Bioengineering*
*University of Pennsylvania*

**Konrad Kording**                                          *koerding@gmail.com*
*Department of Bioengineering*
*University of Pennsylvania*

# Supervised learning

- Know inputs and true output
- Get good at mapping from input to output
- Causal inference is not like this!

# Where do CI/ CD algorithms come from?



Human Knowledge    Data

Algorithm    Inference    Causality

Just like in ML

# Learning causal discovery



Causal inference is supervised learning, after all.

# We need something to test this idea

- Something with known causality
- And nontrivial data

# A causal discovery problem with the microprocessor



A. Causal pairs

B. Non-causal pairs

- Read through two channels. Predict if they are connected

# Perturbations to define causality



A. Average treatment effect of lesioning transistor 1

B. Average treatment effect of lesioning transistor 990

C. Average treatment effect of lesioning transistor 3057

# Honest validation

- Use 1 half of the data to train a causality detector
- Use other half to test if it works

# Standard Transformer ML system

## A. Learning Schematic

P(Causal|X)

Supervised Learning

Global Pooling Head

Sequence Encoder

Window Embedding Layer

Pairwise Sequences

Adjacency Matrix

Causal Effect Calculation

Single Element Perturbation

# Works pretty great

# Also robust to noise

# Generalizes well across games: train on Donkey Kong, test others

# Focuses attention where it matters (transients)

# Simulated fMRI

Table 3: AUPRC comparison on different simulations of NetSim (mean ± std).

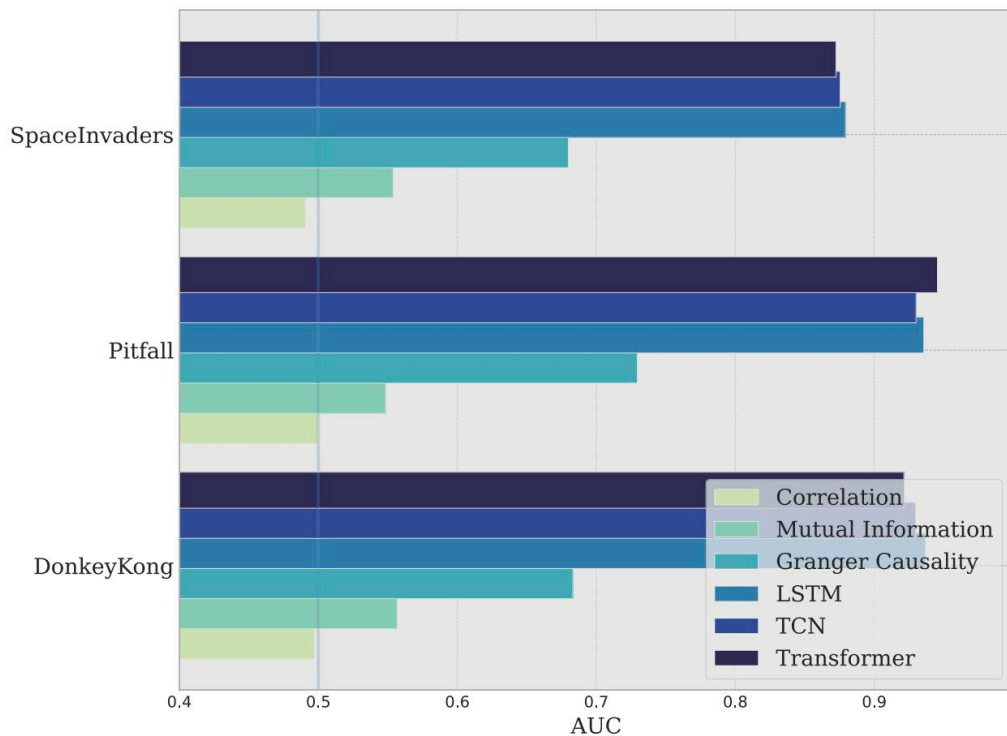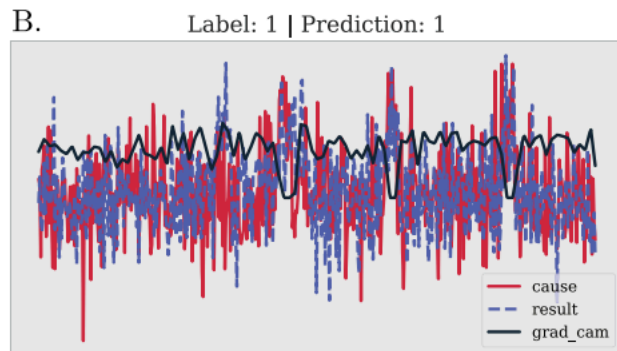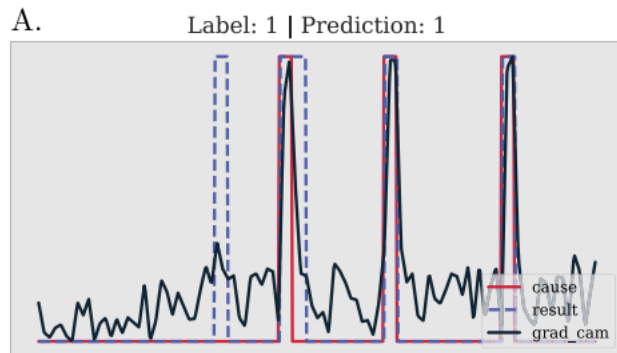| Dataset | Corr | DYNO | GC | MI | PCMCI+ | LiNGAM | cMLP | cLSTM | eSRU | SRU | Transformer |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Sim1 | $0.47 \pm 0.02$ | $0.41 \pm 0.08$ | $0.40 \pm 0.08$ | $0.39 \pm 0.08$ | $0.39 \pm 0.09$ | $0.43 \pm 0.15$ | $0.42 \pm 0.15$ | $0.41 \pm 0.14$ | $0.40 \pm 0.14$ | $0.39 \pm 0.14$ | $\mathbf{0.97 \pm 0.06}$ |
| Sim2 | $0.45 \pm 0.03$ | $0.33 \pm 0.12$ | $0.32 \pm 0.12$ | $0.31 \pm 0.11$ | $0.29 \pm 0.11$ | $0.30 \pm 0.11$ | $0.29 \pm 0.11$ | $0.28 \pm 0.11$ | $0.27 \pm 0.11$ | $0.26 \pm 0.10$ | $\mathbf{0.94 \pm 0.05}$ |
| Sim3 | $0.44 \pm 0.03$ | $0.32 \pm 0.13$ | $0.29 \pm 0.14$ | $0.27 \pm 0.12$ | $0.26 \pm 0.12$ | $0.27 \pm 0.13$ | $0.26 \pm 0.12$ | $0.24 \pm 0.12$ | $0.23 \pm 0.12$ | $0.22 \pm 0.12$ | $\mathbf{0.89 \pm 0.06}$ |
| Sim8 | $0.41 \pm 0.07$ | $0.36 \pm 0.08$ | $0.38 \pm 0.11$ | $0.37 \pm 0.10$ | $0.36 \pm 0.10$ | $0.40 \pm 0.14$ | $0.40 \pm 0.14$ | $0.39 \pm 0.14$ | $0.39 \pm 0.14$ | $0.38 \pm 0.13$ | $\mathbf{0.85 \pm 0.13}$ |
| Sim10 | $0.48 \pm 0.02$ | $0.38 \pm 0.10$ | $0.39 \pm 0.12$ | $0.40 \pm 0.11$ | $0.40 \pm 0.12$ | $0.42 \pm 0.16$ | $0.42 \pm 0.16$ | $0.42 \pm 0.15$ | $0.42 \pm 0.15$ | $0.42 \pm 0.15$ | $\mathbf{0.96 \pm 0.03}$ |
| Sim11 | $0.28 \pm 0.03$ | $0.26 \pm 0.04$ | $0.26 \pm 0.06$ | $0.26 \pm 0.06$ | $0.25 \pm 0.07$ | $0.25 \pm 0.08$ | $0.25 \pm 0.08$ | $0.24 \pm 0.08$ | $0.24 \pm 0.08$ | $0.23 \pm 0.08$ | $\mathbf{0.71 \pm 0.10}$ |
| Sim12 | $0.43 \pm 0.02$ | $0.36 \pm 0.08$ | $0.33 \pm 0.11$ | $0.31 \pm 0.10$ | $0.29 \pm 0.11$ | $0.30 \pm 0.11$ | $0.28 \pm 0.11$ | $0.27 \pm 0.11$ | $0.26 \pm 0.11$ | $0.26 \pm 0.11$ | $\mathbf{0.90 \pm 0.06}$ |
| Sim13 | $0.48 \pm 0.04$ | $0.47 \pm 0.05$ | $0.48 \pm 0.07$ | $0.49 \pm 0.10$ | $0.47 \pm 0.10$ | $0.48 \pm 0.10$ | $0.47 \pm 0.10$ | $0.47 \pm 0.10$ | $0.47 \pm 0.11$ | $0.46 \pm 0.11$ | $\mathbf{0.76 \pm 0.10}$ |
| Sim14 | $0.48 \pm 0.02$ | $0.41 \pm 0.08$ | $0.41 \pm 0.09$ | $0.40 \pm 0.08$ | $0.38 \pm 0.09$ | $0.42 \pm 0.13$ | $0.41 \pm 0.13$ | $0.40 \pm 0.13$ | $0.39 \pm 0.13$ | $0.39 \pm 0.13$ | $\mathbf{0.93 \pm 0.08}$ |
| Sim15 | $0.45 \pm 0.03$ | $0.38 \pm 0.07$ | $0.40 \pm 0.09$ | $0.41 \pm 0.08$ | $0.41 \pm 0.10$ | $0.48 \pm 0.21$ | $0.47 \pm 0.20$ | $0.45 \pm 0.20$ | $0.44 \pm 0.19$ | $0.43 \pm 0.19$ | $\mathbf{0.72 \pm 0.09}$ |
| Sim16 | $0.48 \pm 0.01$ | $0.44 \pm 0.05$ | $0.45 \pm 0.07$ | $0.44 \pm 0.06$ | $0.44 \pm 0.06$ | $0.46 \pm 0.10$ | $0.46 \pm 0.10$ | $0.45 \pm 0.10$ | $0.45 \pm 0.10$ | $0.45 \pm 0.10$ | $\mathbf{0.96 \pm 0.03}$ |
| Sim17 | $0.47 \pm 0.01$ | $0.39 \pm 0.09$ | $0.36 \pm 0.10$ | $0.36 \pm 0.09$ | $0.35 \pm 0.10$ | $0.42 \pm 0.19$ | $0.40 \pm 0.19$ | $0.37 \pm 0.19$ | $0.35 \pm 0.19$ | $0.34 \pm 0.19$ | $\mathbf{0.98 \pm 0.02}$ |
| Sim18 | $0.48 \pm 0.03$ | $0.42 \pm 0.07$ | $0.42 \pm 0.12$ | $0.41 \pm 0.10$ | $0.40 \pm 0.11$ | $0.43 \pm 0.16$ | $0.42 \pm 0.16$ | $0.41 \pm 0.16$ | $0.40 \pm 0.15$ | $0.39 \pm 0.15$ | $\mathbf{0.99 \pm 0.02}$ |
| Sim21 | $0.48 \pm 0.03$ | $0.42 \pm 0.08$ | $0.41 \pm 0.08$ | $0.40 \pm 0.08$ | $0.38 \pm 0.09$ | $0.42 \pm 0.15$ | $0.41 \pm 0.14$ | $0.40 \pm 0.14$ | $0.39 \pm 0.14$ | $0.38 \pm 0.13$ | $\mathbf{0.95 \pm 0.06}$ |
| Sim22 | $\mathbf{0.41 \pm 0.04}$ | $0.38 \pm 0.06$ | $0.38 \pm 0.08$ | $0.39 \pm 0.07$ | $0.37 \pm 0.08$ | $0.37 \pm 0.09$ | $0.35 \pm 0.09$ | $0.34 \pm 0.09$ | $0.34 \pm 0.09$ | $0.34 \pm 0.09$ | $0.31 \pm 0.11$ |
| Sim23 | $0.40 \pm 0.04$ | $0.35 \pm 0.06$ | $0.40 \pm 0.12$ | $0.39 \pm 0.10$ | $0.41 \pm 0.14$ | $\mathbf{0.47 \pm 0.21}$ | $0.45 \pm 0.20$ | $0.43 \pm 0.19$ | $0.42 \pm 0.19$ | $0.41 \pm 0.19$ | $0.41 \pm 0.05$ |
| Sim24 | $0.36 \pm 0.06$ | $0.31 \pm 0.07$ | $0.34 \pm 0.10$ | $0.35 \pm 0.10$ | $0.35 \pm 0.11$ | $0.35 \pm 0.11$ | $0.34 \pm 0.11$ | $0.34 \pm 0.11$ | $0.34 \pm 0.11$ | $0.33 \pm 0.11$ | $\mathbf{0.37 \pm 0.11}$ |

# Gene networks

Table 5: AUPRC comparison on the Dream3

| Dataset | Corr | DYNO | GC | MI | PCMCI+ | cMLP | cLSTM | eSRU | SRU | Transformer |
|---------|------|------|------|------|--------|------|-------|------|------|-------------|
| Ecoli2 | 0.01 | 0.01 | 0.01 | 0.02 | 0.01 | 0.01 | 0.02 | 0.01 | 0.01 | **0.04** |
| Yeast2 | 0.01 | 0.01 | 0.02 | 0.01 | 0.01 | 0.04 | 0.04 | 0.04 | 0.04 | **0.06** |
| Yeast3 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.06 | **0.07** | 0.06 | 0.06 | 0.06 |

# What does this mean

- ## Cons
  - ○ Requires a lot of ground truth data
  - ○ Impossible for humans to understand
  - ○ Biased

- ## Pros
  - ○ If we have enough ground truth it will beat all humans
  - ○ Meaningful proofs (inherited from empirical risk minimization)
  - ○ A new approach to get into messy observational spaces

# Take home message

- Causality and ML really have overlapping problems
- So much to learn
- Check out Clear conference

- Exclusion can be a meaningful data problem
- Maybe we should look for more ways to learn how to do CD/ CI, and prepare ourselves for a future where the bitter lesson is key to medicine