

DDA4210/AIR6002 Advanced Machine Learning

Lecture 09 Causal Machine Learning

Tongxin Li

School of Data Science, CUHK-Shenzhen

Spring 2024

Motivation

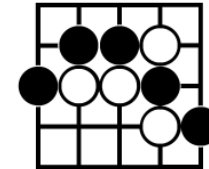
Trustworthy ML

AI Tools are Everywhere



AI/ML Methods/Predictions

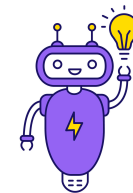
Suggested Moves



Suggested Strategies



AI Assistants



Review

First Part of This Course:

- Ensemble
- Learning Theory
- GNN
- Generative Models

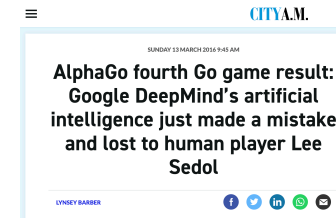
Focus more on a single merit: accuracy

AI Tools Are **NOT** Always Trustworthy

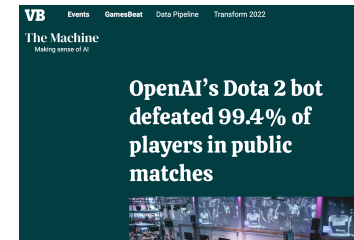


Suggested Moves →

Blackbox AI Tools/Imperfect Predictions



Suggested Strategies →



Defeated by 0.6% players using adversarial strategies



AI Assistants →



ChatGPT版必应疑似「发疯」，微软紧急限制回答数目，植入广告赚钱提上日程

GPT-4 Technical Report

OpenAI*

Abstract

We report the development of GPT-4, a large-scale, multimodal model which can accept image and text inputs and produce text outputs. While less capable than humans in many real-world scenarios, GPT-4 exhibits human-level performance on various professional and academic benchmarks, including passing a simulated bar exam with a score around the top 10% of test takers. GPT-4 is a Transformer-based model pre-trained to predict the next token in a document. The post-training alignment process results in improved performance on measures of factuality and adherence to desired behavior. A core component of this project was developing infrastructure and optimization methods that behave predictably across a wide range of scales. This allowed us to accurately predict some aspects of GPT-4's performance based on models trained with no more than 1/1,000th the compute of GPT-4.

A Toy Example

GPT-4 Technical Report

OpenAI*

Abstract

we used python



A Toy Example

| [GPT-4 answers correctly] | [GPT-4 answers incorrectly] |
|--|---|
| <p>Can you teach an old dog new tricks?</p> <p>Yes, you can teach an old dog new tricks ← choice</p> <p>You can't teach an old dog new tricks</p> <p>No, you can't teach an old dog new tricks</p> | <p>Son of an actor, this American guitarist and rock singer released many songs and albums and toured with his band. His name is "Elvis" what?</p> <p>Perkins</p> <p>Presley ← choice</p> <p>Elvis Presley</p> <p>His name is Elvis Presley</p> |

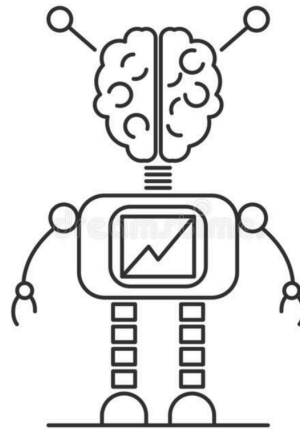
Table 4: Example of GPT-4 giving correct and incorrect responses on TruthfulQA

2 Scope and Limitations of this Technical Report

This report focuses on the capabilities, limitations, and safety properties of GPT-4. GPT-4 is a Transformer-style model [33] pre-trained to predict the next token in a document, using both publicly available data (such as internet data) and data licensed from third-party providers. The model was then fine-tuned using Reinforcement Learning from Human Feedback (RLHF) [34]. Given both the competitive landscape and the safety implications of large-scale models like GPT-4, this report contains no further details about the architecture (including model size), hardware, training compute, dataset construction, training method, or similar.

We are committed to independent auditing of our technologies, and shared some initial steps and ideas in this area in the system card accompanying this release.² We plan to make further technical details available to additional third parties who can advise us on how to weigh the competitive and safety considerations above against the scientific value of further transparency.

Trustworthy Methods Connect AI to Physical Worlds



Outlook

Second Part of This Course:

- **Causal Learning** (This lecture)
- Differential Privacy and Federated Learning
- Fairness in ML
- Explainable AI (XAI)

Focus on more attributes: **causality, privacy, fairness, and interpretability**

This Lecture:

Introduction to Causal Learning

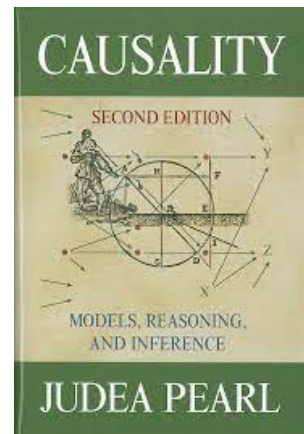
Outline

- Examples Simpson's Paradox
- Causal inference Backdoor Adjustment ~~Formal Definitions~~
- Causal discovery Nonlinear ICA ~~The PC Algorithm~~
- Disentanglement Identifiable VAE

Outline

Many online resources and talks

- Causal inference
- Causal discovery
- Disentanglement

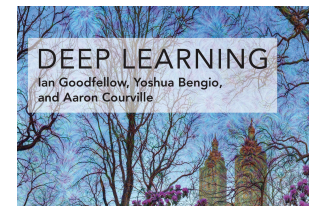


Introduction to Causal Inference

from a Machine Learning Perspective

Brady Neal

December 17, 2020



TOWARDS CAUSAL REPRESENTATION LEARNING:

AN AI & DEEP LEARNING PERSPECTIVE ON CAUSALITY

YOSHUA BENGIO

Causal learning is a full course in many schools

We will only cover selective topics

Part I

Causal Inference

Part I.1

Simpson's paradox and Examples

Motivating example: Simpson's paradox

Simpson's paradox: COVID-29

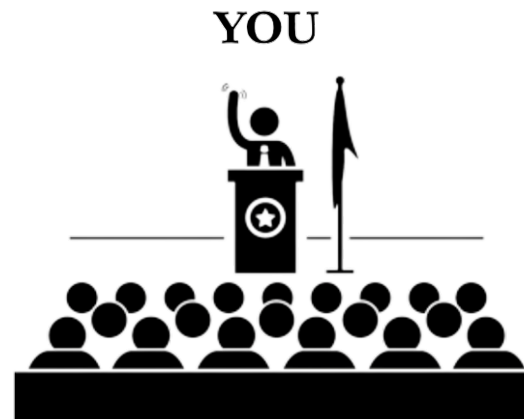
New Virus



Treatment T : A (0) or B (1)

Condition C : Mild (0) or Severe (1)

Outcome Y : Happy (0) or Unhappy (1)



Motivating example: Simpson's paradox

Simpson's paradox: COVID-29

New Virus



Treatment T : A (0) or B (1)

Condition C : Mild (0) or Severe (1)

Outcome Y : Survive (0) or Not (1)

Mortality Rate Table

| | Total |
|---|--------------------------|
| A | 16% (240/1500) |
| B | 19% (105/550) |

$E[Y|T]$

Simpson's paradox: Mortality Rate Table

Mortality Rate Table

| | |
|---|--------------------------|
| | Total |
| A | 16% (240/1500) |
| B | 19% (105/550) |

$\mathbb{E}[Y|T]$

Simpson's paradox: Mortality Rate Table

Mortality Rate Table

| | | Condition | | |
|-----------|---|--------------------------|--------------------------|--------------------------|
| | | Mild | Severe | Total |
| Treatment | A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) |
| | B | 10% (5/50) | 20% (100/500) | 19% (105/550) |
| | | $\mathbb{E}[Y T, C = 0]$ | $\mathbb{E}[Y T, C = 1]$ | $\mathbb{E}[Y T]$ |

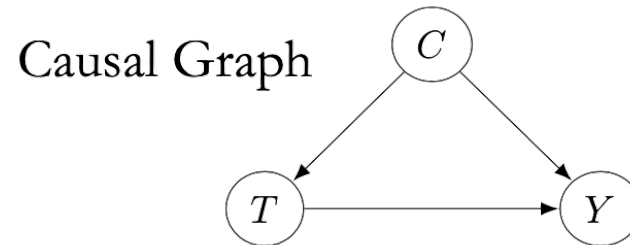
Simpson's paradox: Mortality Rate Table

Mortality Rate Table [Statistics/Data](#)

| | | Condition | | | |
|-----------|---|--------------------------|--------------------------|--------------------------|---|
| | | Mild | Severe | Total | |
| Treatment | A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) | $\frac{1400}{1500}(0.15) + \frac{100}{1500}(0.30) = 0.16$ |
| | B | 10% (5/50) | 20% (100/500) | 19% (105/550) | $\frac{50}{550}(0.10) + \frac{500}{550}(0.20) = 0.19$ |
| | | $\mathbb{E}[Y T, C = 0]$ | $\mathbb{E}[Y T, C = 1]$ | $\mathbb{E}[Y T]$ | |

Which treatment should you choose?

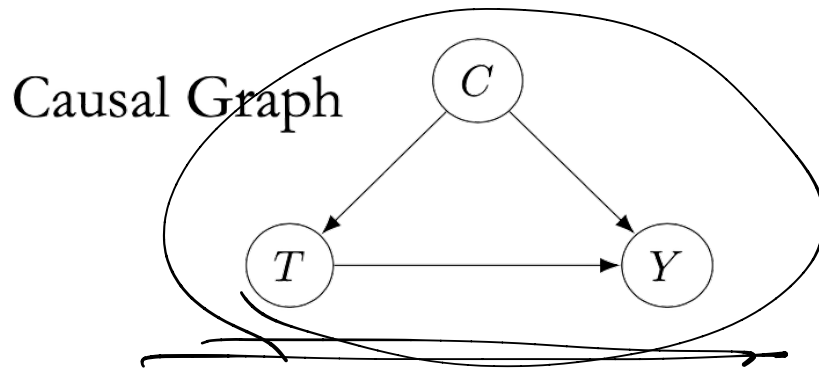
Simpson's paradox: scenario I



Which treatment is better?

| | | Condition | | |
|-----------|---|----------------------|-------------------------|--------------------------|
| | | Mild | Severe | Total |
| Treatment | A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) |
| | B | 10% (5/50) | 20% (100/500) | 19% (105/550) |

Simpson's paradox: scenario I (treatment B)



Condition

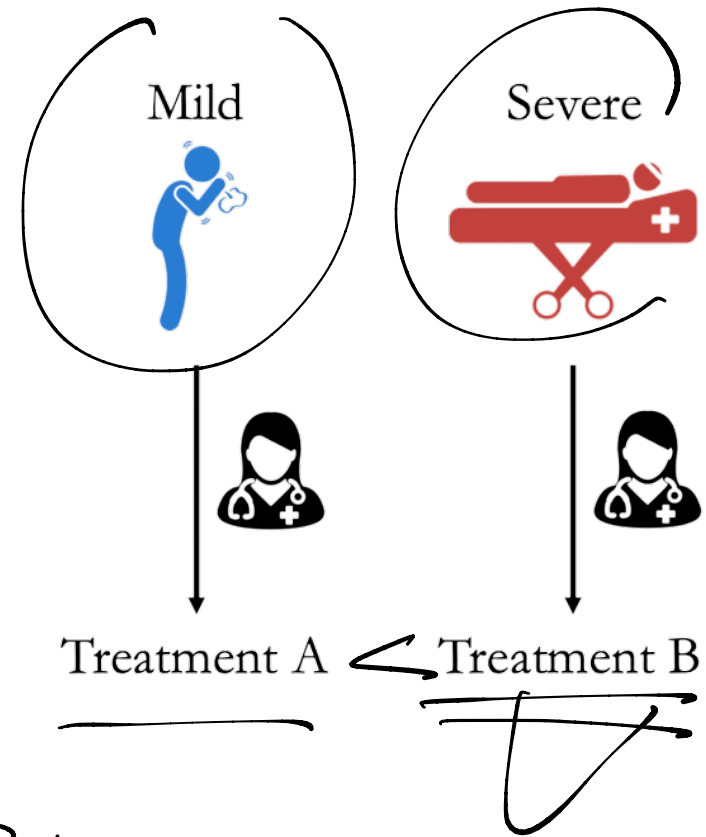
| | Mild | Severe | Total |
|---|-------------------|------------------|-------------------|
| A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) |
| B | 10% (5/50) | 20% (100/500) | 19% (105/550) |

Treatment

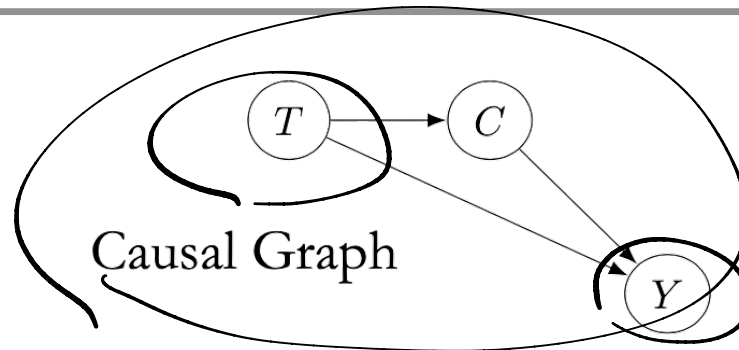
1400 >> 100

500 >> 50

Data



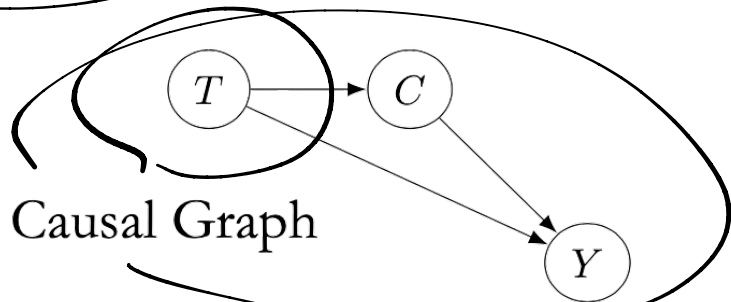
Simpson's paradox: scenario II



Which treatment is better?

| | | Condition | | |
|-----------|---|-------------------|------------------|-------------------|
| | | Mild | Severe | Total |
| Treatment | A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) |
| | B | 10% (5/50) | 20% (100/500) | 19% (105/550) |

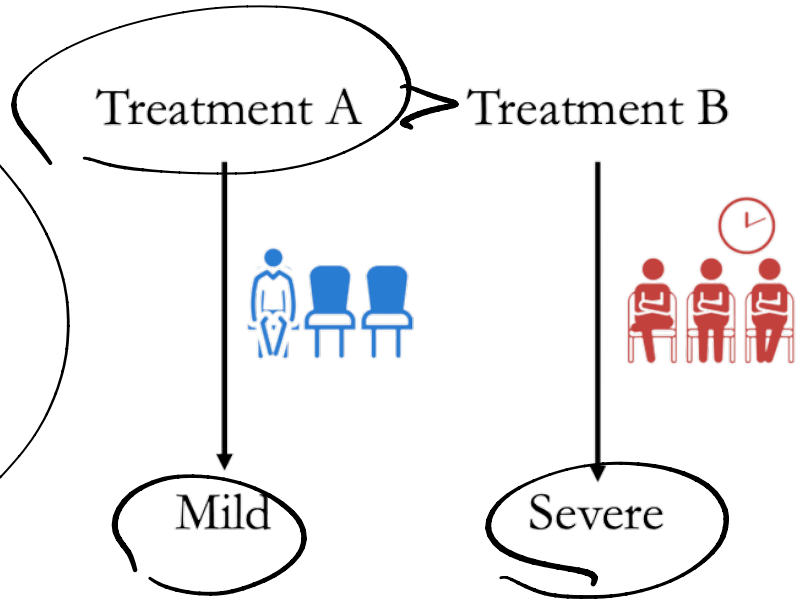
Simpson's paradox: scenario II (treatment A)



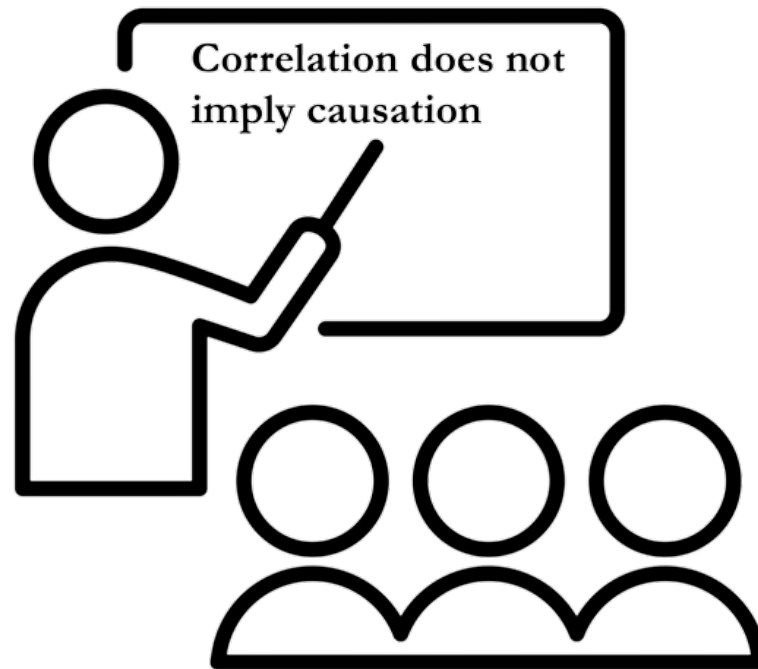
Which treatment is better?

Treatment

| | Condition | | |
|---|-------------------|------------------|-------------------|
| | Mild | Severe | Total |
| A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) |
| B | 10% (5/50) | 20% (100/500) | 19% (105/550) |

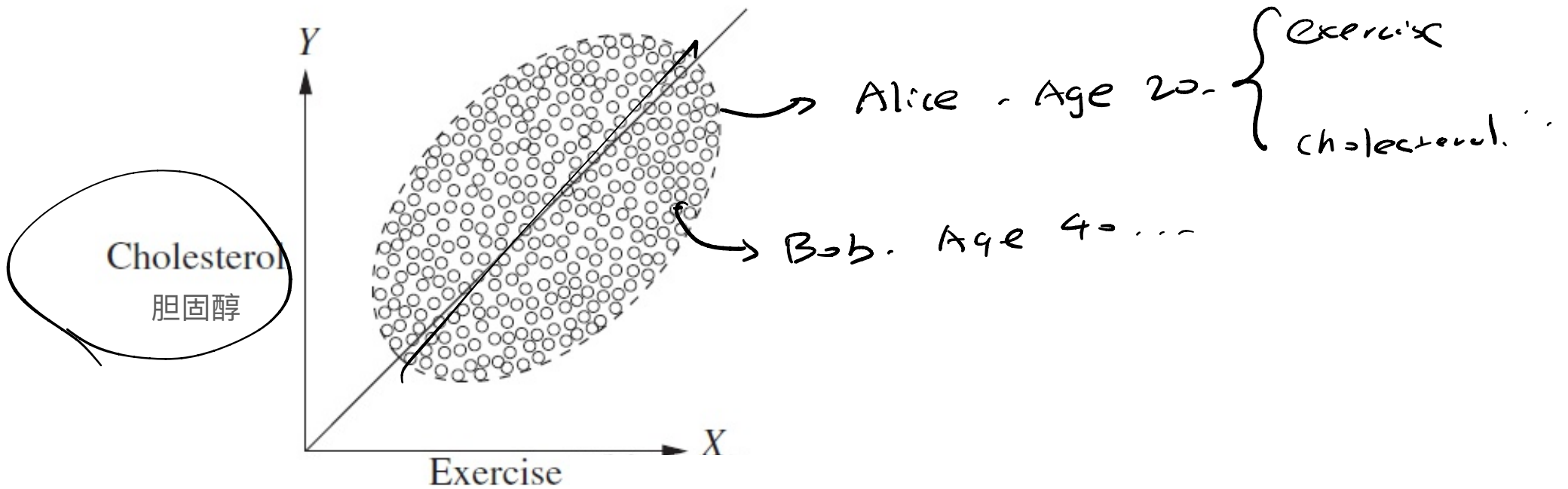


What do we learn from the Simpson's paradox?



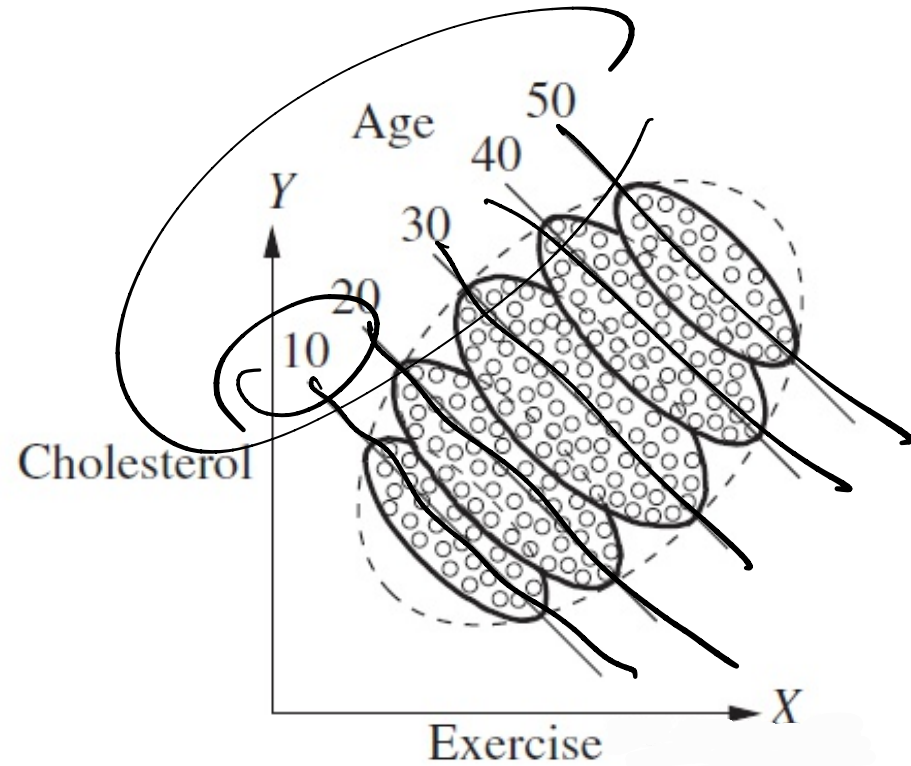
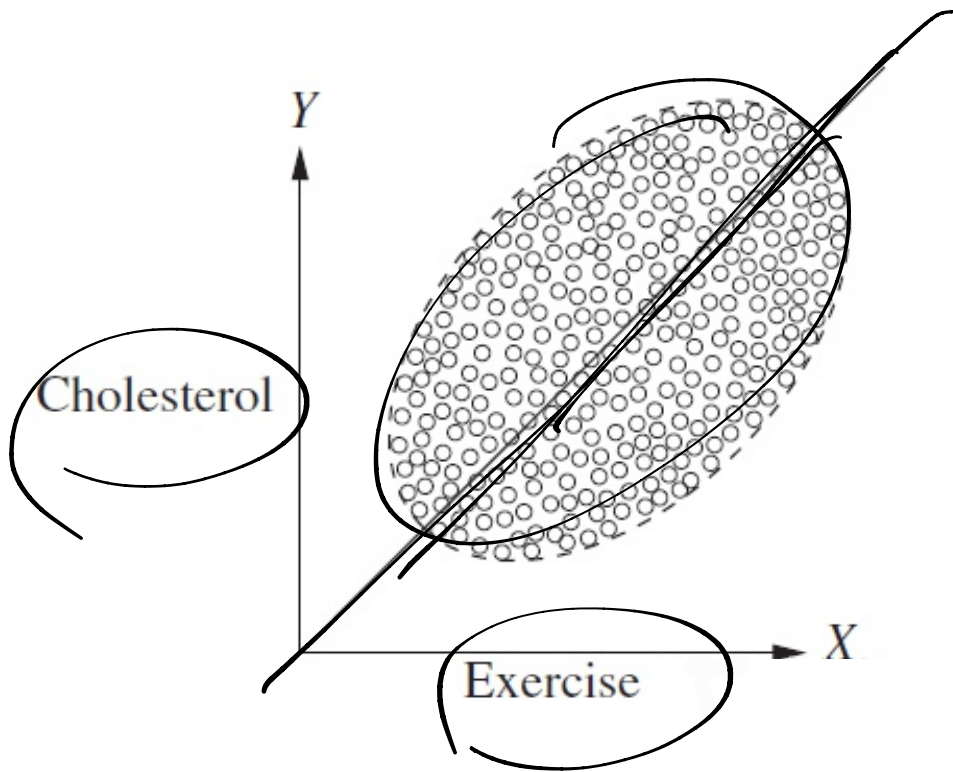
What do we learn from the Simpson's paradox?

More examples



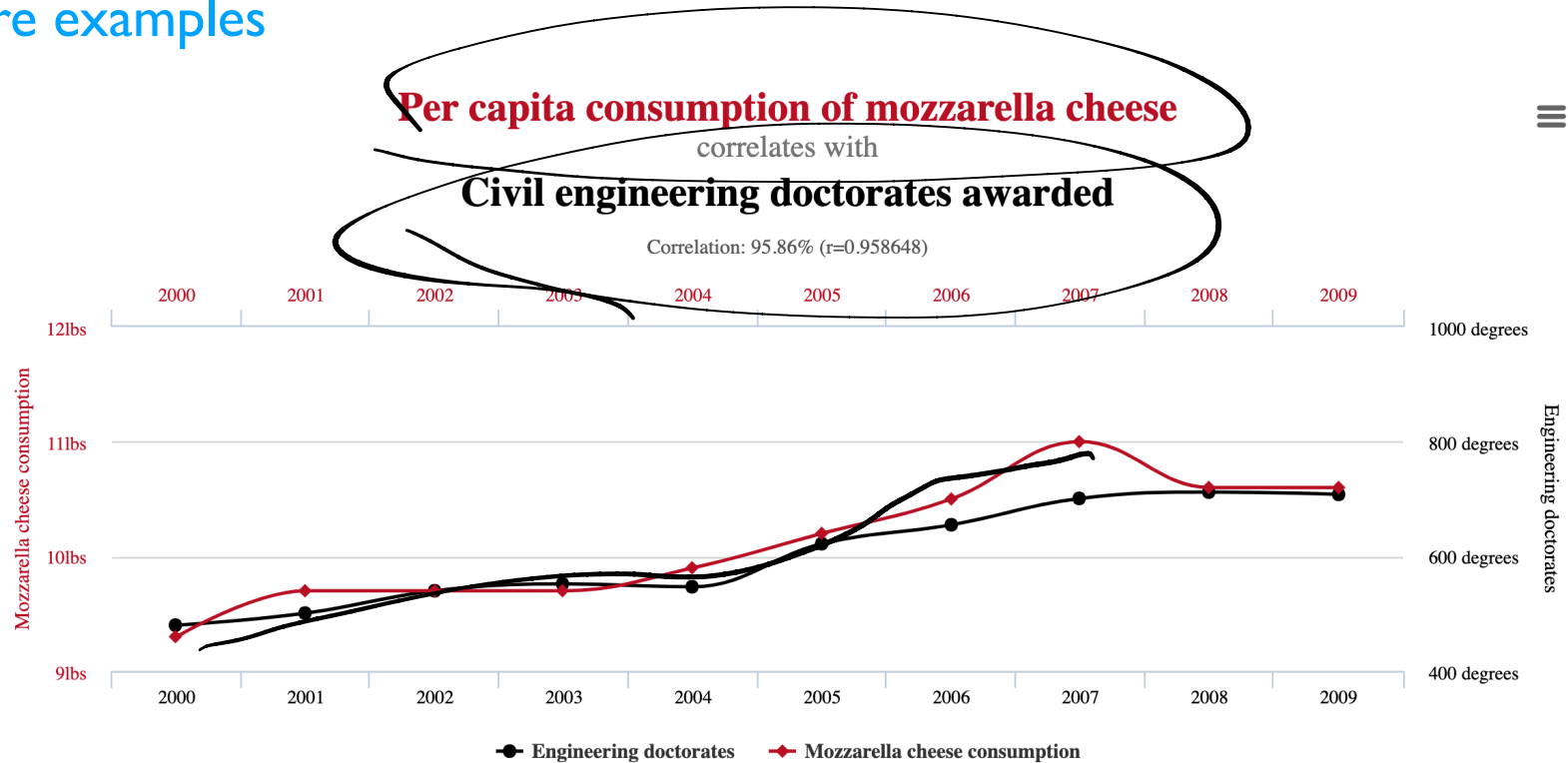
What do we learn from the simpson's paradox?

More examples



What do we learn from the simpson's paradox?

More examples



Source: <https://www.tylervigen.com/spurious-correlations>

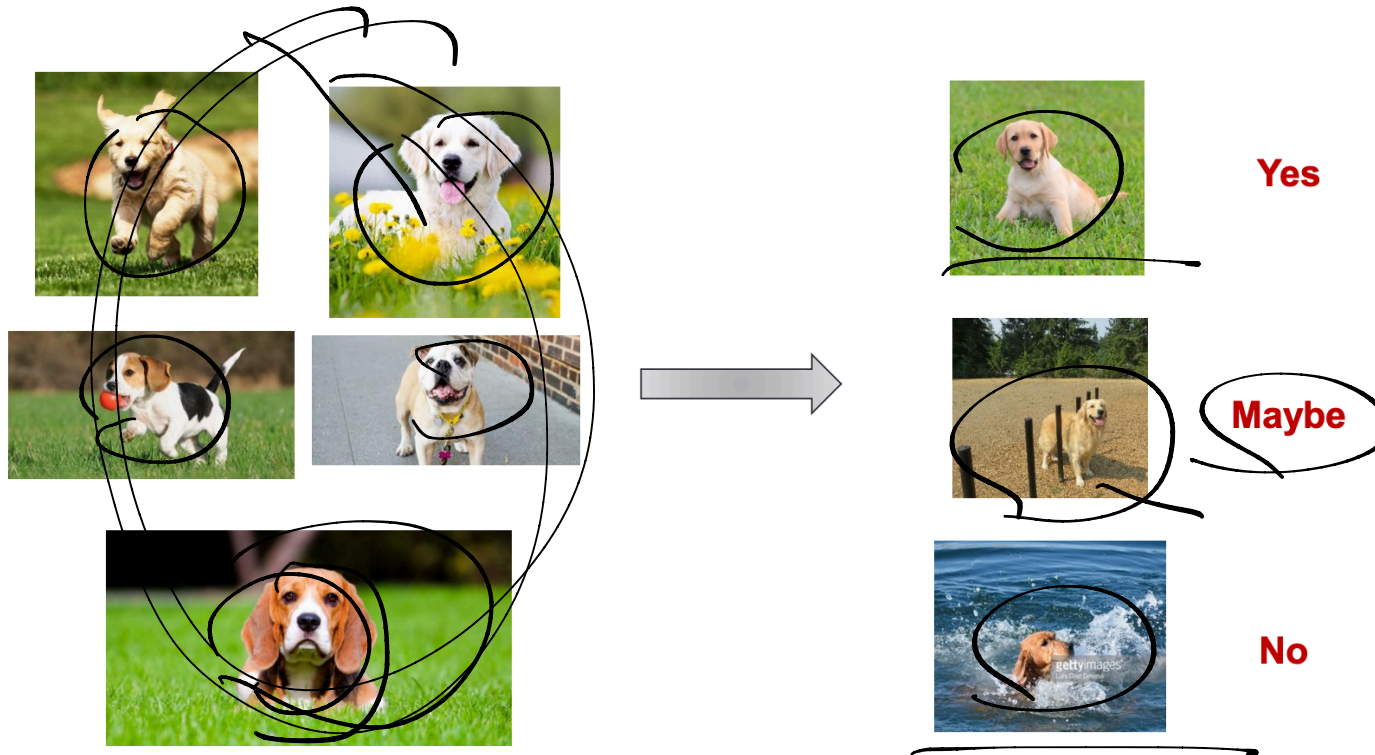
What do we learn from the simpson's paradox?

Correlation does not imply causation

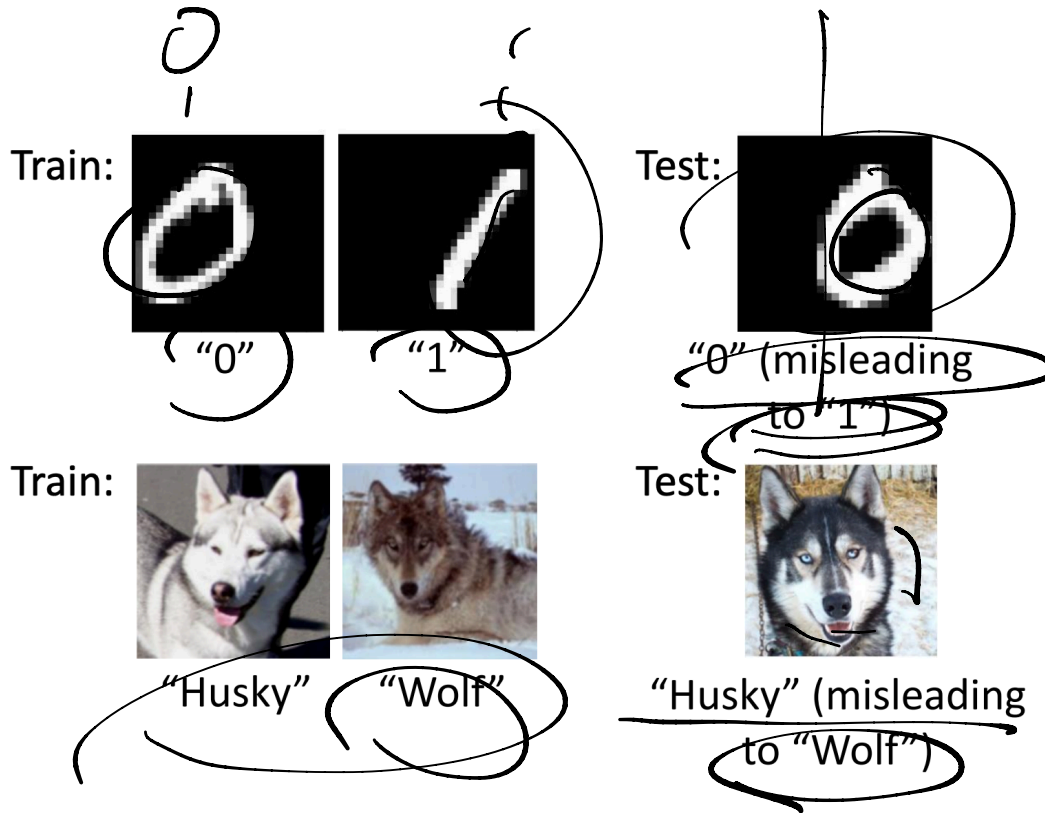
Correlation is not enough

Statistical learning vs Causal learning

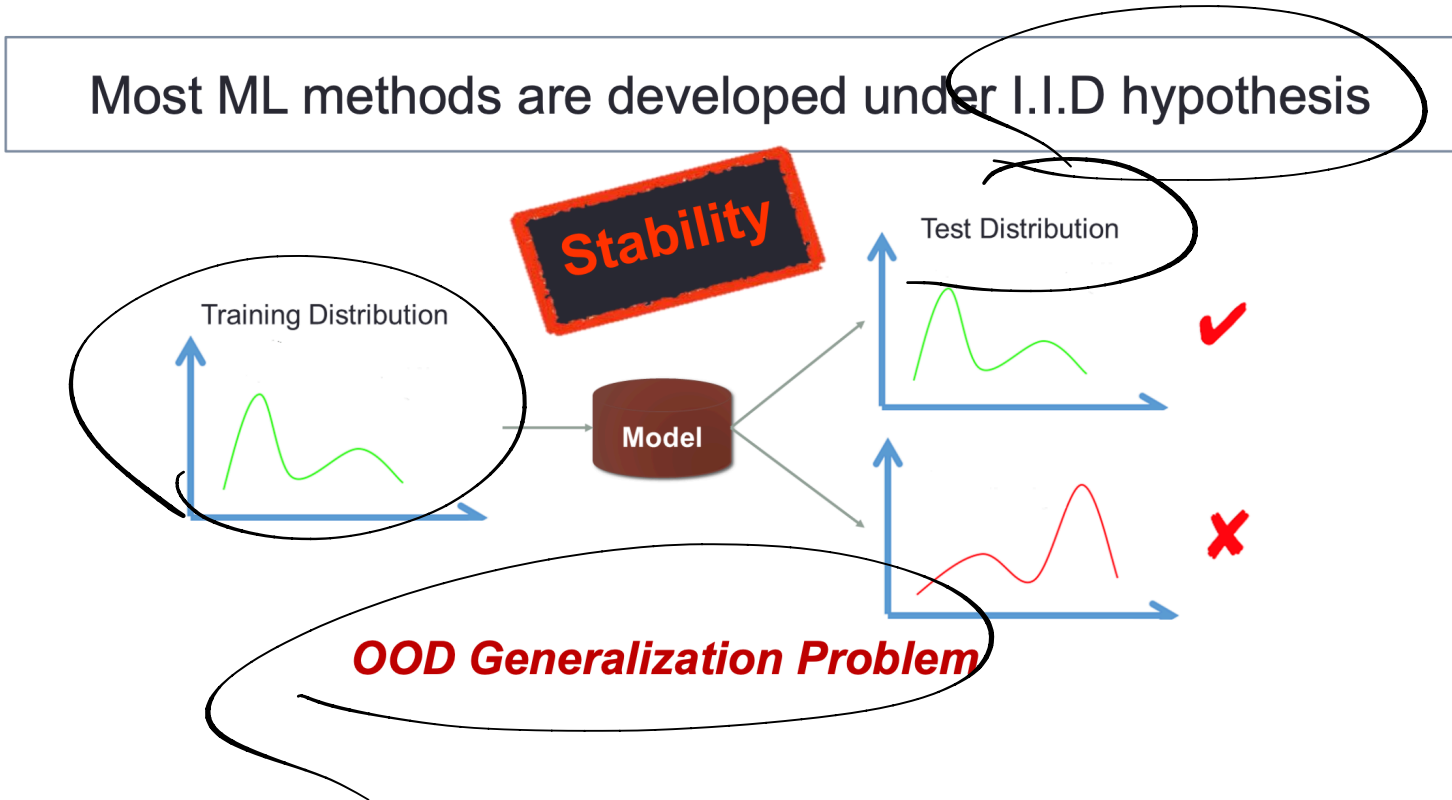
Why causality matters in machine learning?



Why causality matters in machine learning?

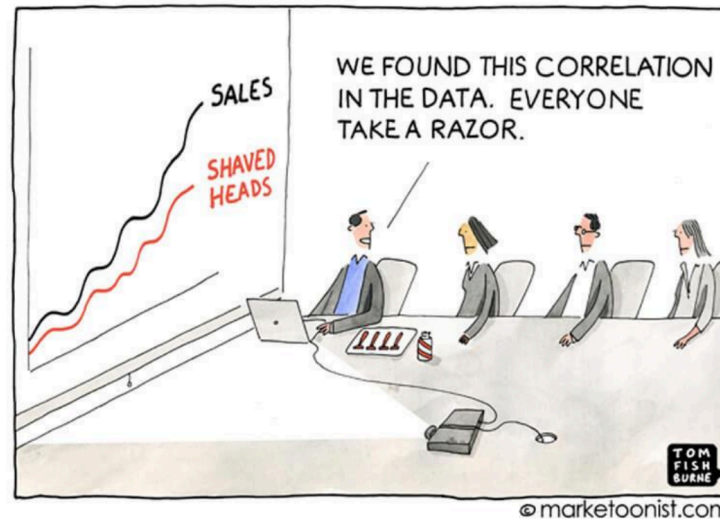


Why causality matters in machine learning?



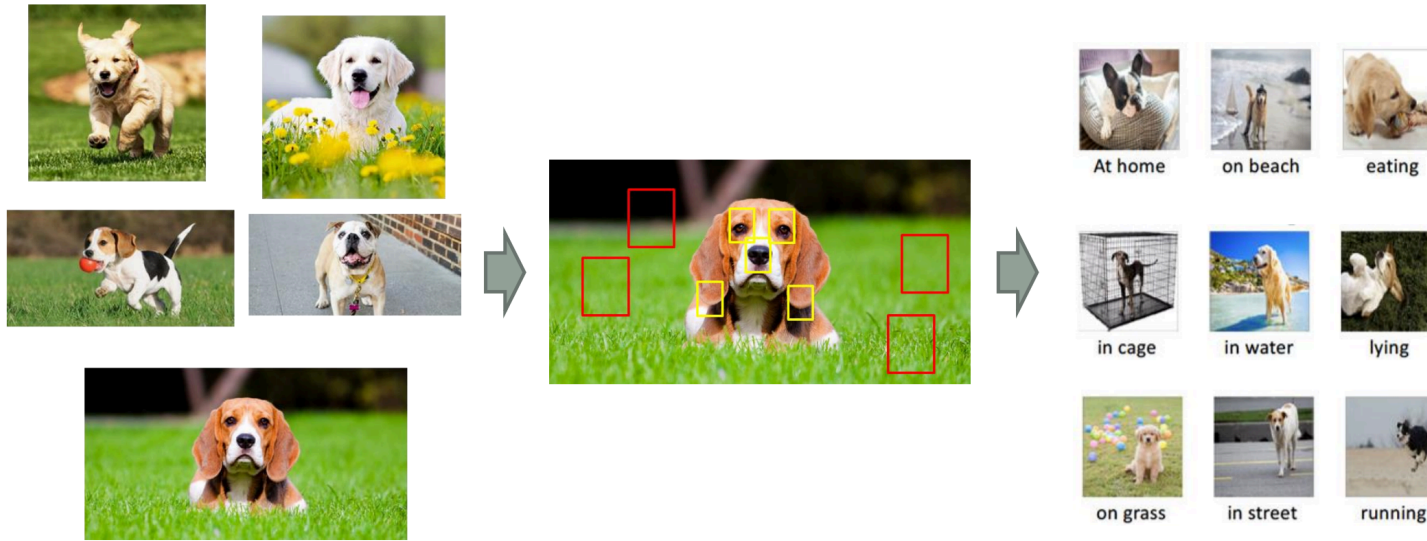
Why causality matters in machine learning?

Correlation is the very basics of machine learning.



Why causality matters in machine learning?

Relying solely on **correlation** can cause problems



Part I.2

Simpson's paradox and Examples

What is causal inference?

Inferring the effects of any treatment/policy/intervention/etc.

Examples:

- Effect of treatment on a disease
- Effect of climate change policy on emissions
- Effect of social media on mental health
- Many more (effect of X on Y)

What is causal inference?

Inferring the effects of any treatment/policy/intervention/etc.

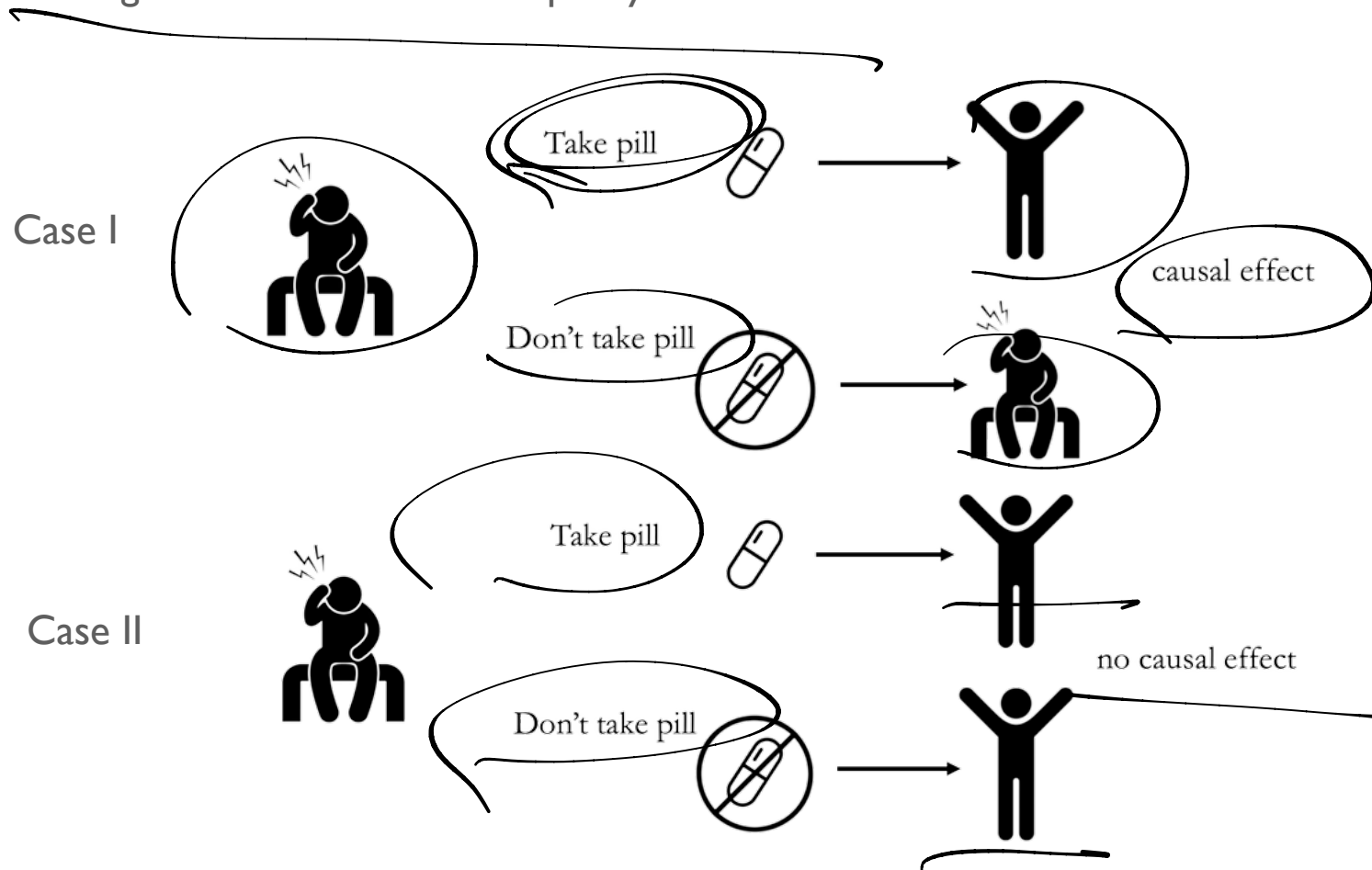
- Examples:**
- Effect of treatment on a disease
 - Effect of climate change policy on emissions
 - Effect of social media on mental health
 - Many more (effect of X on Y)

How do we measure causal effects with interventions?

How do we measure causal effects in observational studies?

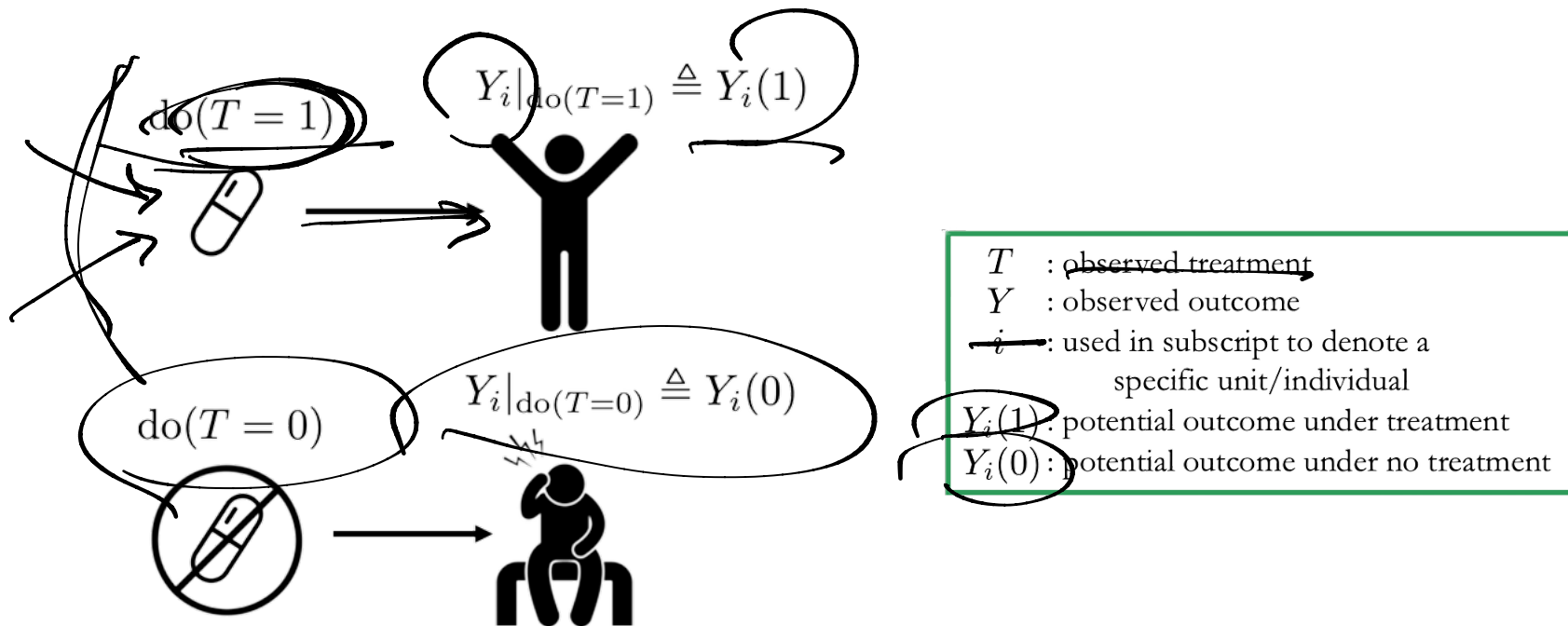
Potential outcomes

Inferring the effect of treatment/policy on some outcome



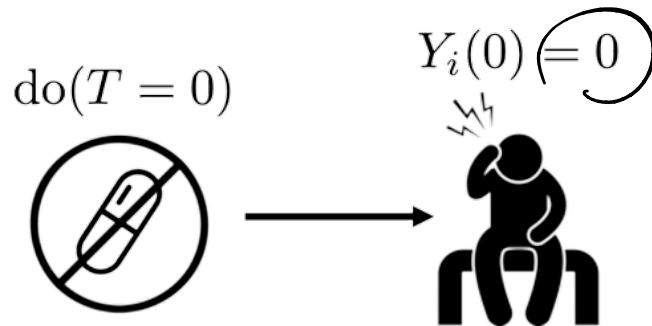
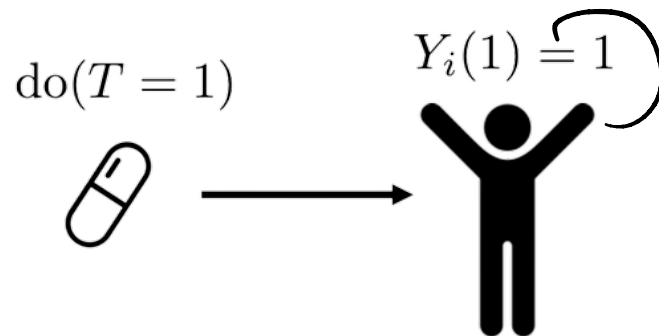
Do Operator

Inferring the effect of treatment/policy on some outcome



Do Operator

Inferring the effect of treatment/policy on some outcome

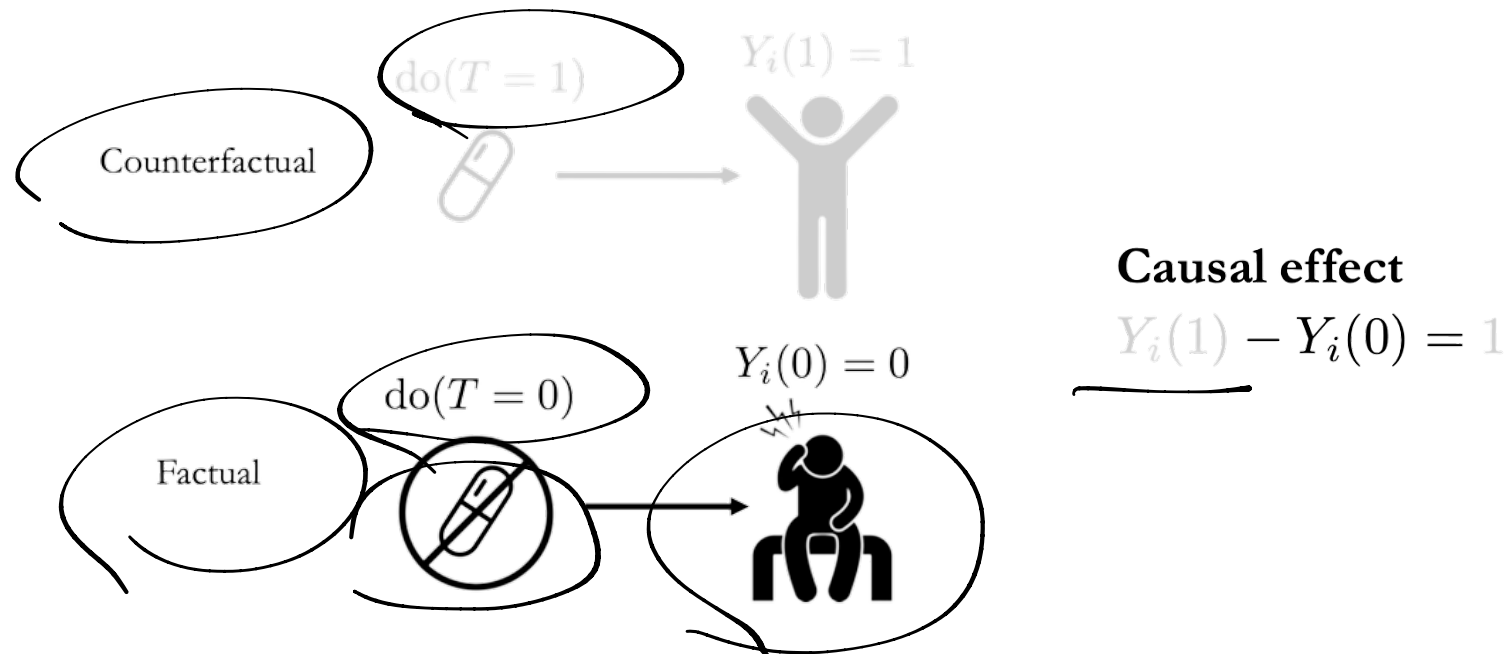


T : observed treatment
 Y : observed outcome
 i : used in subscript to denote a specific unit/individual
 $Y_i(1)$: potential outcome under treatment
 $Y_i(0)$: potential outcome under no treatment

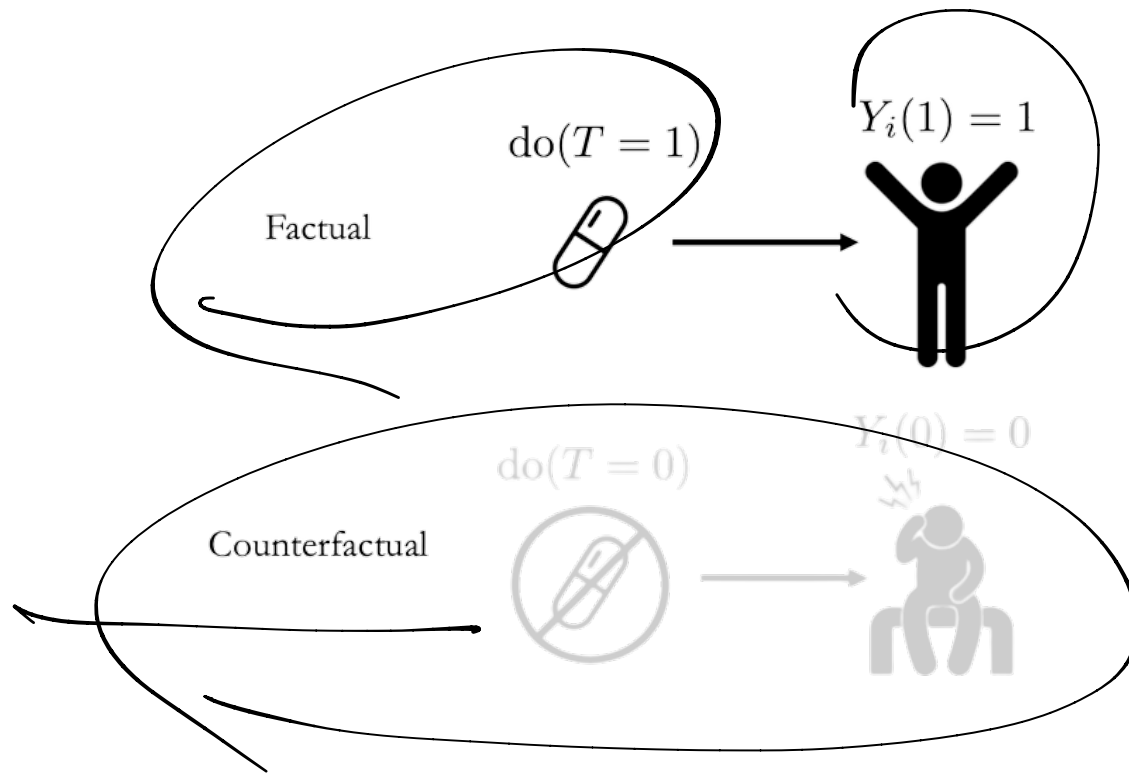
Causal effect

$$Y_i(1) - Y_i(0) = 1$$

A Fundamental Problem of Causal Inference



A Fundamental Problem of Causal Inference



Causal effect

$$Y_i(1) - Y_i(0) = 1$$

Causal inference with observations (Optional)

Inferring the effects of any treatment/policy/intervention/etc.

- Examples:**
- Effect of treatment on a disease
 - Effect of climate change policy on emissions
 - Effect of social media on mental health
 - Many more (effect of X on Y)

How do we measure causal effects with interventions?

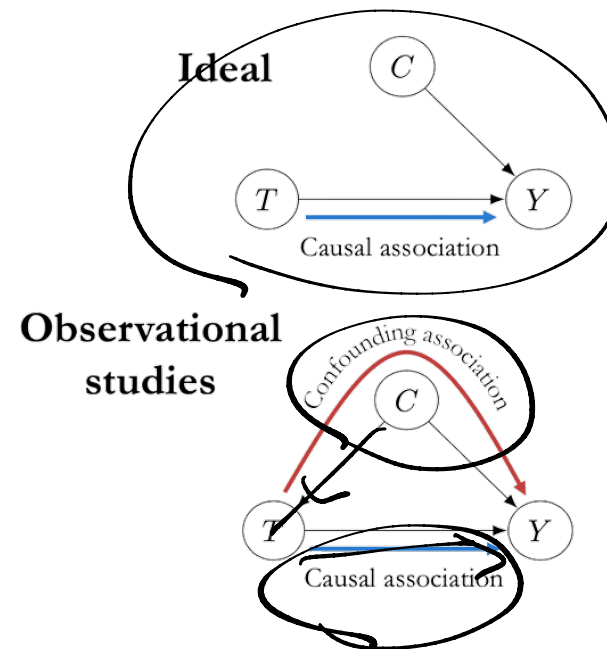
How do we measure causal effects in observational studies?

Causal inference with observations (Optional)

How do we measure causal effects in observational studies?

Can't always randomize treatment

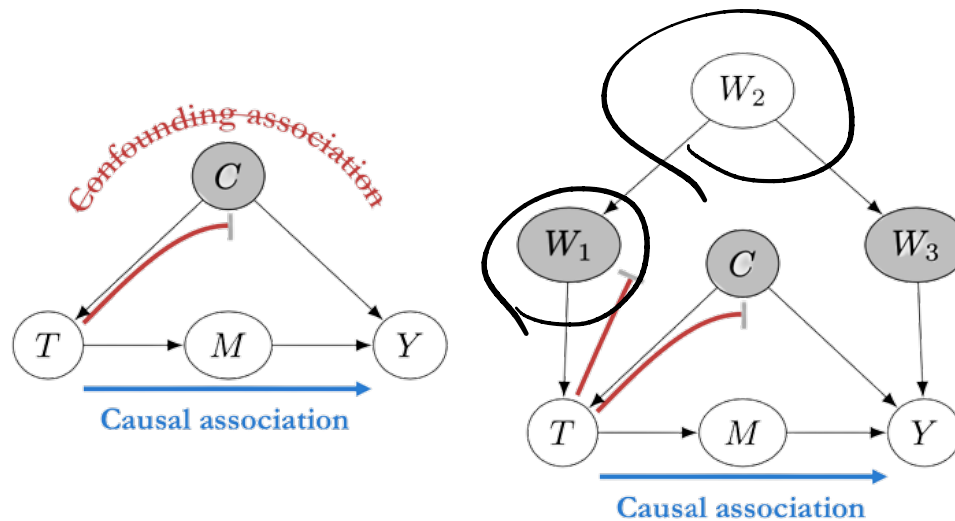
- **Ethical reasons** (e.g. unethical to randomize people to smoke for measuring effect on lung cancer)
- **Infeasibility** (e.g. can't randomize countries into communist/capitalist systems to measure effect on GDP)
- **Impossibility** (e.g. can't change a living person's DNA at birth for measuring effect on breast cancer)



Causal inference with observations (Optional)

Solution: backdoor adjustment

Formal assumptions are needed (omitted)



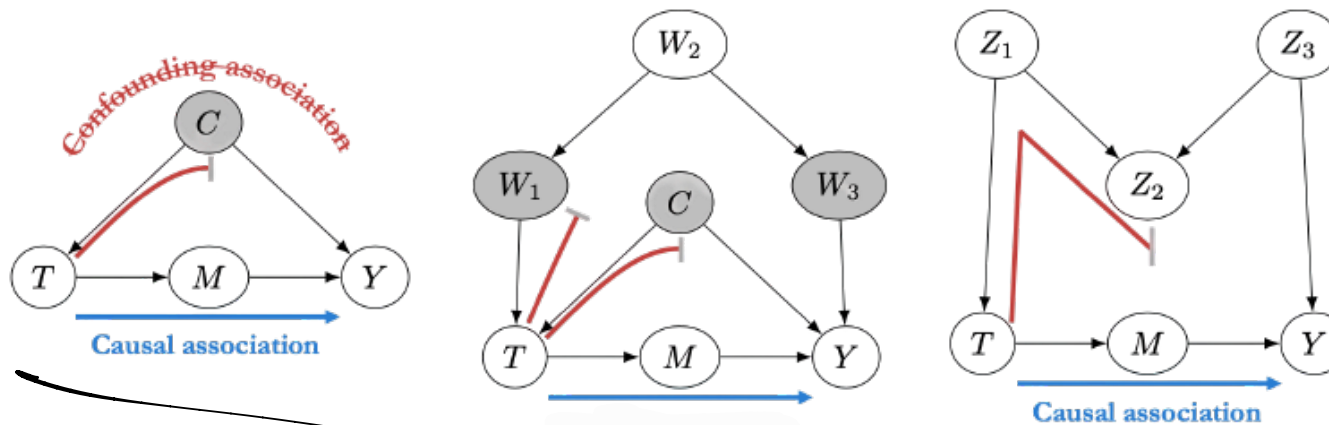
Causal inference with observations (optional)

Solution: backdoor adjustment

Formal assumptions are needed (omitted)

$$\mathbb{E}[Y|\text{do}(T = t)] = \mathbb{E}_W \mathbb{E}[Y|t, W]$$

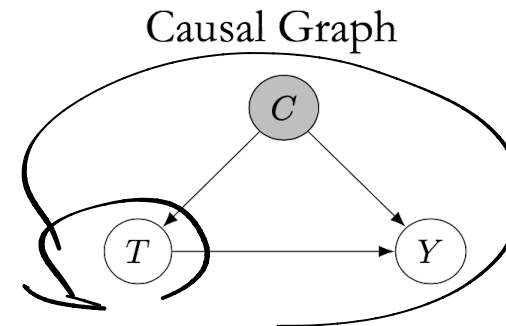
Shaded nodes are examples of sufficient adjustment sets W



Application to the COVID-29 example ★

| | | Condition | | | Causal |
|-----------|---|-------------------|------------------|-------------------|--------------|
| | | Mild | Severe | Total | |
| Treatment | A | 15% (210/1400) | 30% (30/100) | 16% (240/1500) | 19.4% |
| | B | 10% (5/50) | 20% (100/500) | 19% (105/550) | 12.9% |
| | | $E[Y t, C=0]$ | $E[Y t, C=1]$ | $E[Y t]$ | $E[Y do(t)]$ |

Assume this causal graph:



$$E[Y|do(T=t)] = E_C E[Y|t, C] = \sum E[Y|t, c] P(c)$$

$t=A$

$$\frac{1450}{2050} (0.15) + \frac{600}{2050} (0.30) \approx 0.194$$

$P(c=mild)$ \rightarrow $P(c=severe)$

$t=B$

$$\frac{1450}{2050} (0.10) + \frac{600}{2050} (0.20) \approx 0.129$$

Part II

Causal Discovery

Part II

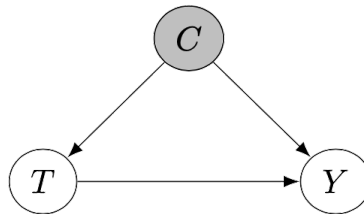
Causal Discovery

I. Linear Case

What is causal discovery?

How do we know this relation for the COVID-29 example?

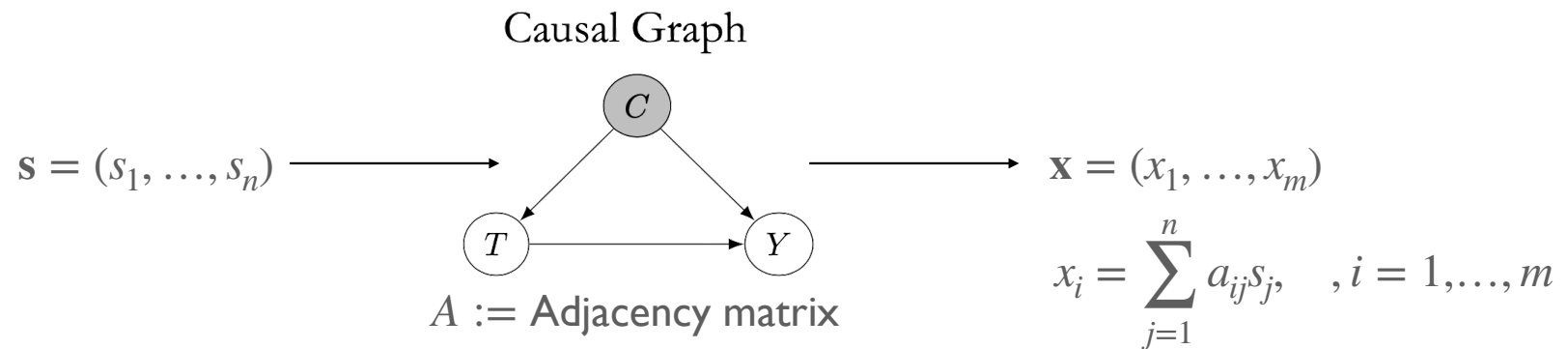
Causal Graph



A key problem: **identifiability**

What is causal discovery?

How do we know this relation for the COVID-29 example?



A key problem: **identifiability**

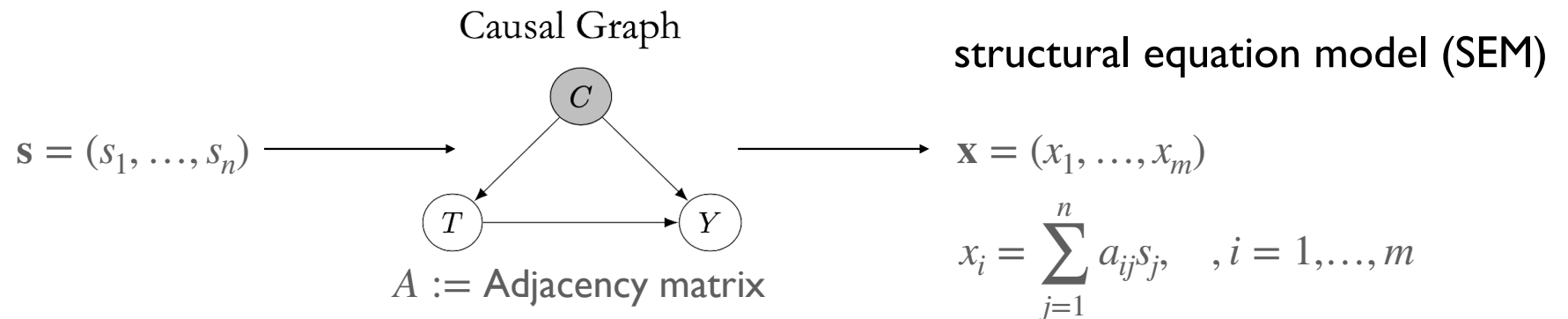
Many methods ...

We focus on ICA in this lecture

Suppose the underlying mechanism is **linear**

What is causal discovery?

How do we know this relation for the COVID-29 example?



Identifiability: Observe \mathbf{x} , want A and \mathbf{s}

Many methods ...

We focus on ICA in this lecture

Suppose the underlying mechanism is linear

Why causality matters in machine learning?

Three sources of correlation:

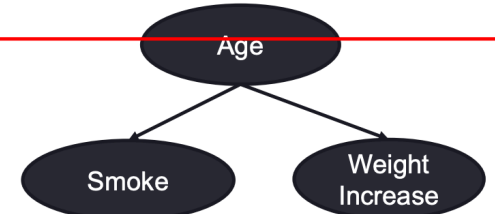
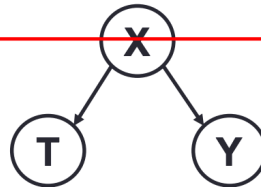
- **Causation**

- Causal mechanism
- **Stable and explainable**



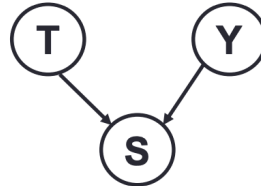
- **Confounding**

- Ignoring X
- **Spurious Correlation**



- **Sample Selection Bias**

- Conditional on S
- **Spurious Correlation**



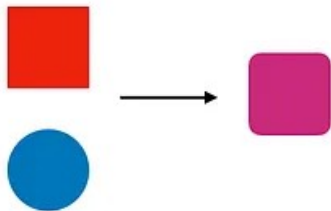
How do identify them, **as a graph?**

Independent Component Analysis

ICA as principled unsupervised learning

PCA

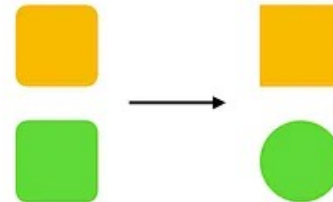
Compresses information



Requires preprocessing: autoscaling

ICA

Separates information



Requires preprocessing: autoscaling

Often benefits from first applying PCA

Independent Component Analysis

ICA as principled unsupervised learning

Unsupervised learning can have different goals

- 1) Accurate model of data distribution?
 - ▶ E.g. Variational Autoencoders are good
- 2) Sampling points from data distribution?
 - ▶ E.g. Generative Adversarial Networks are good
- 3) Useful features for supervised learning?
 - ▶ Many methods, “Representation learning”
- 4) Reveal underlying structure in data, disentangle latent quantities?
 - ▶ Independent Component Analysis

Independent Component Analysis

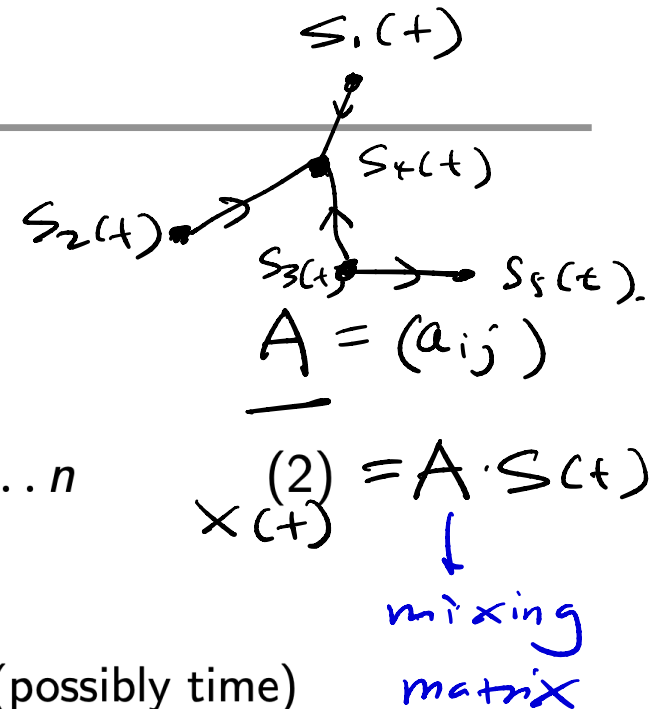
ICA as principled unsupervised learning

Linear independent component analysis (ICA)

$$x_i(t) = \sum_{j=1}^n a_{ij} s_j(t) \quad \text{for all } i, j = 1 \dots n$$

- ▶ $x_i(t)$ is i -th observed signal at sample point t (possibly time)
- ▶ a_{ij} constant parameters describing "mixing"
- ▶ Assuming independent, non-Gaussian latent "sources" s_j

Identifiability: Find Independent Components (Sources)



Independent Component Analysis

ICA as principled unsupervised learning

Linear independent component analysis (ICA)

$$x_i(t) = \sum_{j=1}^n a_{ij}s_j(t) \quad \text{for all } i, j = 1 \dots n \quad (2)$$

- ▶ $x_i(t)$ is i -th observed signal at sample point t (possibly time)
- ▶ a_{ij} constant parameters describing “mixing”
- ▶ Assuming independent, non-Gaussian latent “sources” s_j

The independent components are identifiable (up to permutation and scaling of the sources)

Assumptions: At most one of the sources s_j is Gaussian

$A = (a_{ij})$ is full-rank

Independent Component Analysis

ICA as principled unsupervised learning

Linear independent component analysis (ICA)

$$x_i(t) = \sum_{j=1}^n a_{ij} s_j(t) \quad \text{for all } i, j = 1 \dots n \quad (2)$$

- ▶ $x_i(t)$ is i -th observed signal at sample point t (possibly time)
- ▶ a_{ij} constant parameters describing “mixing”
- ▶ Assuming independent, non-Gaussian latent “sources” s_j

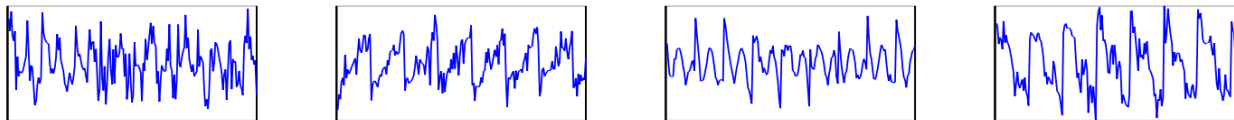
ICA is **identifiable**, i.e. well-defined: (Darmois-Skitovich ~1950; Comon, 1994)

- ▶ Observing only x_i we can recover both a_{ij} and s_j
- ▶ I.e. **original sources can be recovered**

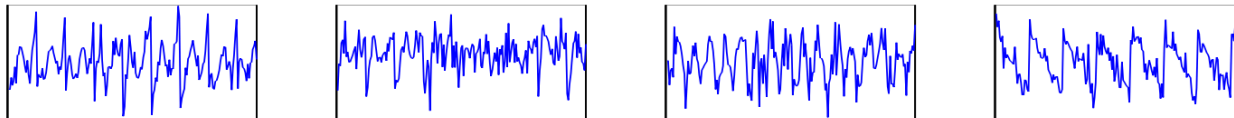
Independent Component Analysis

Identifiability means ICA does blind source separation

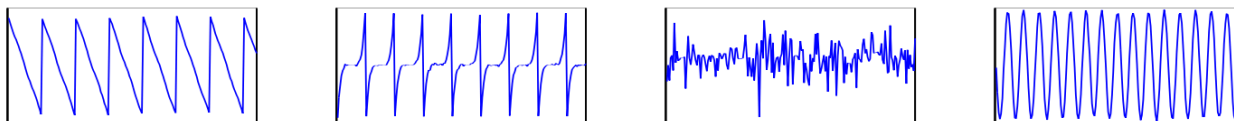
Observed signals:



Principal components:

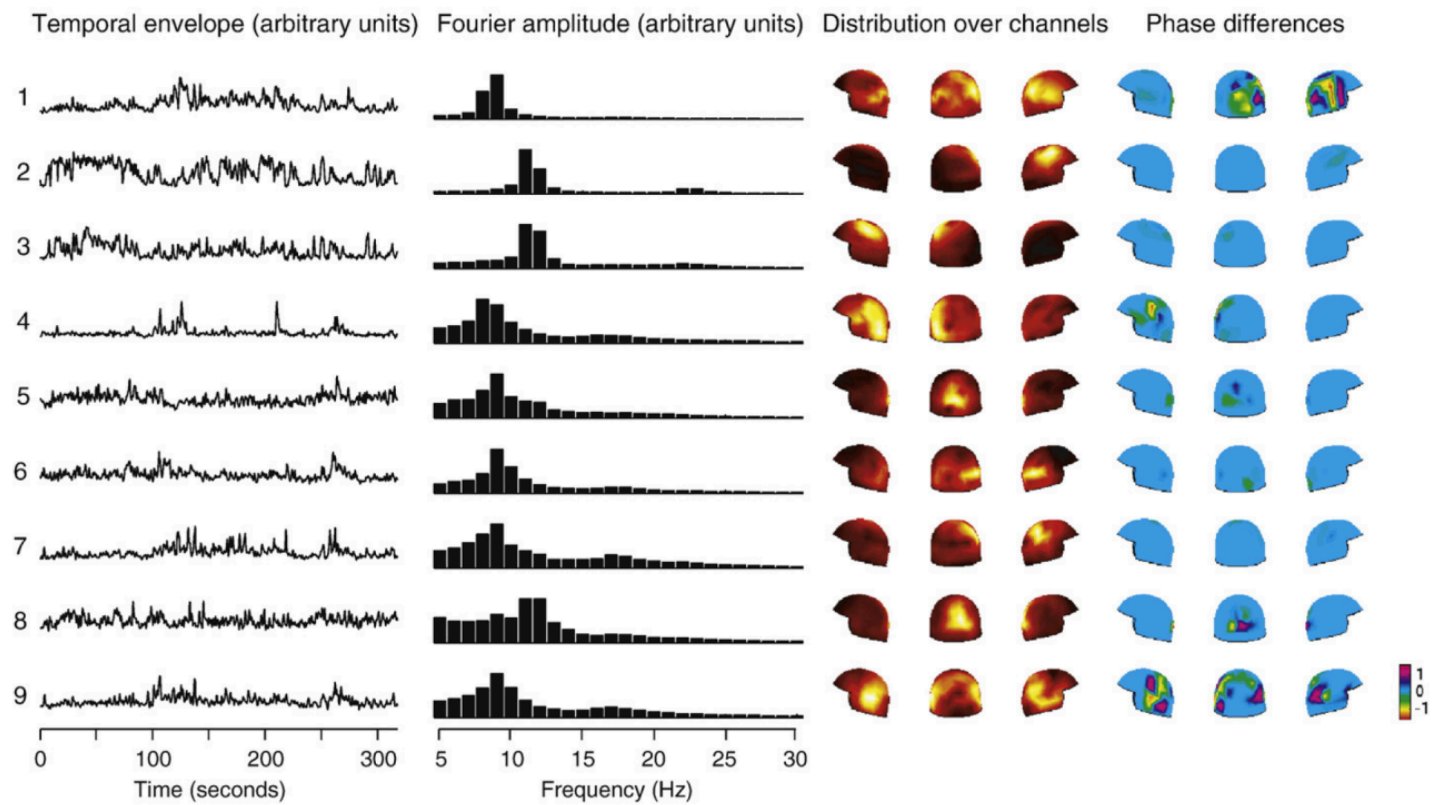


Independent components are original sources:



credits:Aapo Hyvarinen

Independent Component Analysis



(Hyvärinen, Ramkumar, Parkkonen, Hari, 2010)

credits:Aapo Hyvarinen

Part II

Causal Discovery and Disentanglement

2. General Case

Nonlinear Independent Component Analysis (*OPTIONAL*)

What if we consider the nonlinear setting?

Linear ICA: $\mathbf{x} = A\mathbf{s}$

Deep generative models: $\mathbf{x} = f(\mathbf{s})$ What is f^{-1} ?

Identifiability of the deep latent-variable models.

$$p_{\theta}(\mathbf{x}) = p_{\theta^*}(\mathbf{x}) \longrightarrow \theta^* = \theta \longrightarrow p_{\theta^*}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}, \mathbf{z})$$

$$\implies p_{\theta^*}(\mathbf{z}) = p_{\theta}(\mathbf{z})$$

$$p_{\theta^*}(\mathbf{x} | \mathbf{z}) = p_{\theta}(\mathbf{x} | \mathbf{z})$$

Disentanglement

Nonlinear Independent Component Analysis (OPTIONAL)

Disentanglement

Better Mixing via Deep Representations

Yoshua Bengio¹

Dept. IRO, Université de Montréal. Montréal (QC), H2C 3J7, Canada

CHECKMY@WEBPAGE.CA

Grégoire Mesnil¹

Dept. IRO, Université de Montréal. Montréal (QC), H2C 3J7, Canada
LITIS EA 4108, Université de Rouen. 768000 Saint Etienne du Rouvray, France

CHECKMY@WEBPAGE.CA

Yann Dauphin

Salah Rifai

Dept. IRO, Université de Montréal. Montréal (QC), H2C 3J7, Canada

CHECKMY@WEBPAGE.CA

CHECKMY@WEBPAGE.CA

Find disentangled representations in unsupervised data.

An important topic in causal learning

A problem in deep generative models

Identifiability of Nonlinear Independent Component Analysis (optional)

Identifiability

$$p_{\theta}(\mathbf{x}) = p_{\hat{\theta}}(\mathbf{x}) \implies \theta = \hat{\theta} \quad \forall(\theta, \hat{\theta})$$

Deep generative models: $\mathbf{x} = f(\mathbf{s})$ What is f^{-1} ?

Deep generative models ~~are~~ are **not identifiable in general**

(Hyvärinen and Pajunen, 1999; Khemakhem et al., 2020; Locatello et al., 2019)

\implies basic VAEs, GANs, Nonlinear ICA etc. are unidentifiable:

Identifiability problem

$$p_{\mathbf{f}}(\mathbf{x}) = p_{\hat{\mathbf{f}}}(\mathbf{x}) \not\Rightarrow \mathbf{f} = \hat{\mathbf{f}}$$

Identifiability of Nonlinear Independent Component Analysis (optional)

Identifiability

$$p_{\theta}(\mathbf{x}) = p_{\hat{\theta}}(\mathbf{x}) \implies \theta = \hat{\theta} \quad \forall(\theta, \hat{\theta})$$

Deep generative models: $\mathbf{x} = f(\mathbf{s})$ What is f^{-1} ?

Deep generative models are **not identifiable in general**

(Hyvärinen and Pajunen, 1999; Khemakhem et al., 2020; Locatello et al., 2019)

\implies basic VAEs, GANs, Nonlinear ICA etc. are unidentifiable:

We can add structures/assumptions on the distribution of \mathbf{s} to ensure identifiability

Identifiability of Nonlinear Independent Component Analysis (*OPTIONAL*)

Identifiability

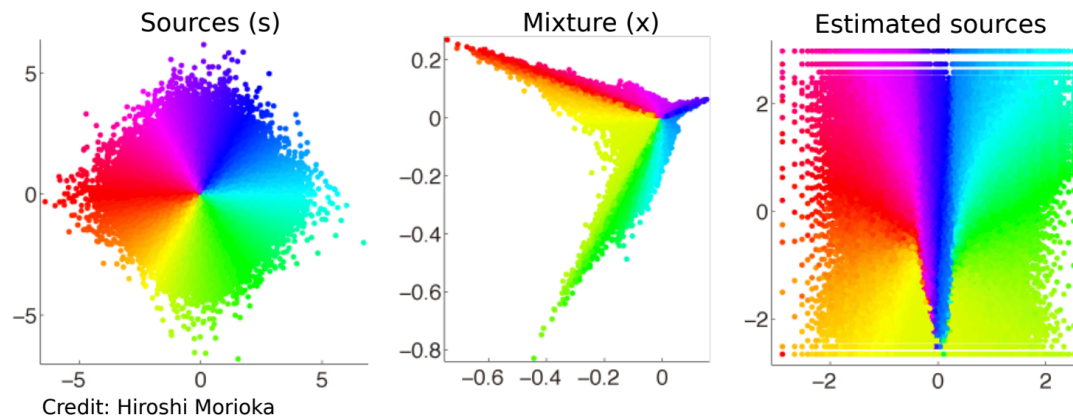
$$p_{\theta}(\mathbf{x}) = p_{\hat{\theta}}(\mathbf{x}) \implies \theta = \hat{\theta} \quad \forall(\theta, \hat{\theta})$$

Deep generative models: $\mathbf{x} = f(\mathbf{s})$

Deep generative models are **not identifiable in general**

(Hyvärinen and Pajunen, 1999; Khemakhem et al., 2020; Locatello et al., 2019)

\implies basic VAEs, GANs, Nonlinear ICA etc. are unidentifiable:

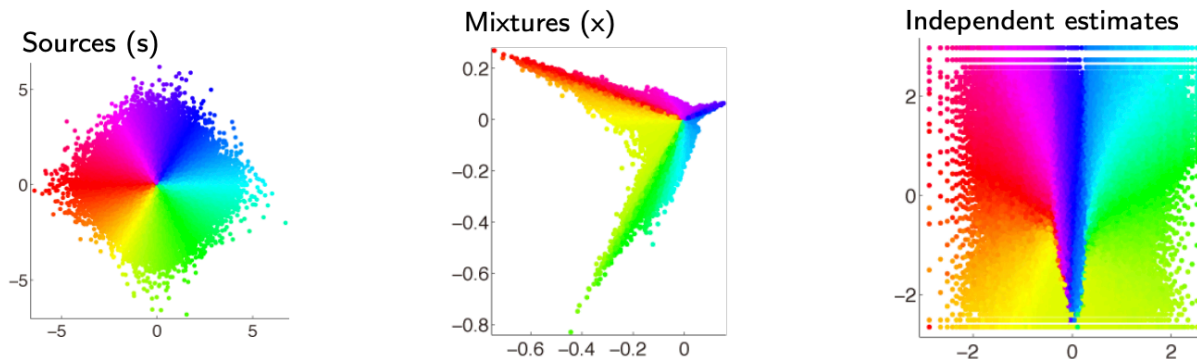


Identifiability of Nonlinear Independent Component Analysis (*OPTIONAL*)

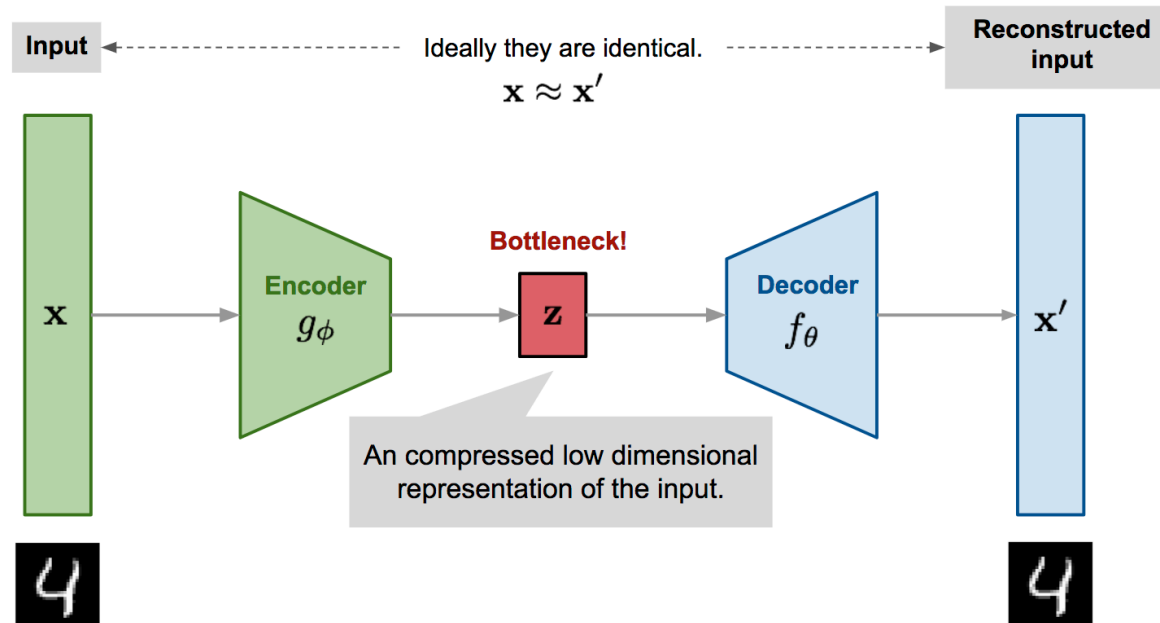
Extend ICA to nonlinear case to get deep learning?
Unfortunately, “basic” nonlinear ICA is **not identifiable**:
If we define nonlinear ICA model simply as

$$x_i(t) = f_i(s_1(t), \dots, s_n(t)) \quad \text{for all } i, j = 1 \dots n$$

we cannot recover original sources (Darmois, 1952; Hyvärinen & Pajunen, 1999)



β -VAE



Why?

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x} | \mathbf{z})] - \beta D_{KL}(\log q_\theta(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})) \quad \text{Increase } \beta \text{ can encourage disentanglement}$$

β -VAE

From the last lecture:

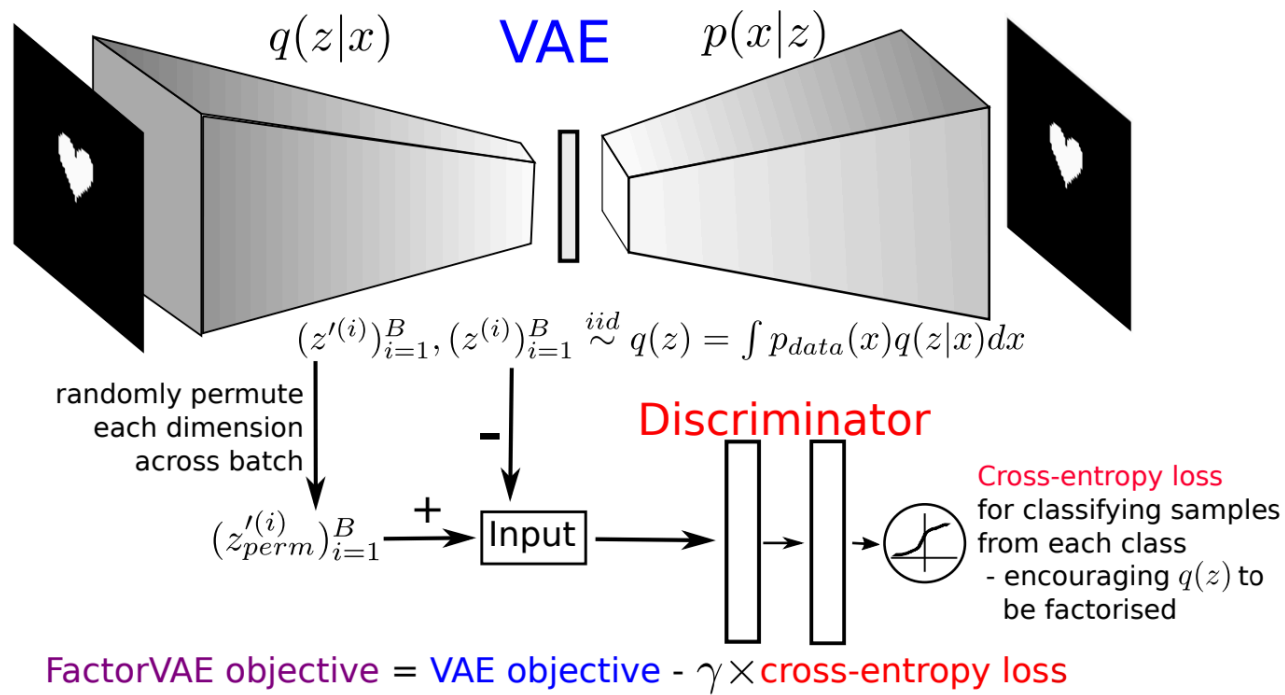
$$\max_{\phi, \theta} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \\ - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))$$

- maximize $\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]$:
reconstruct \mathbf{x}
- minimize $D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))$:
approximate prior

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \beta D_{KL}(\log q_{\theta}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})) \quad \text{Increase } \beta \text{ can encourage disentanglement}$$

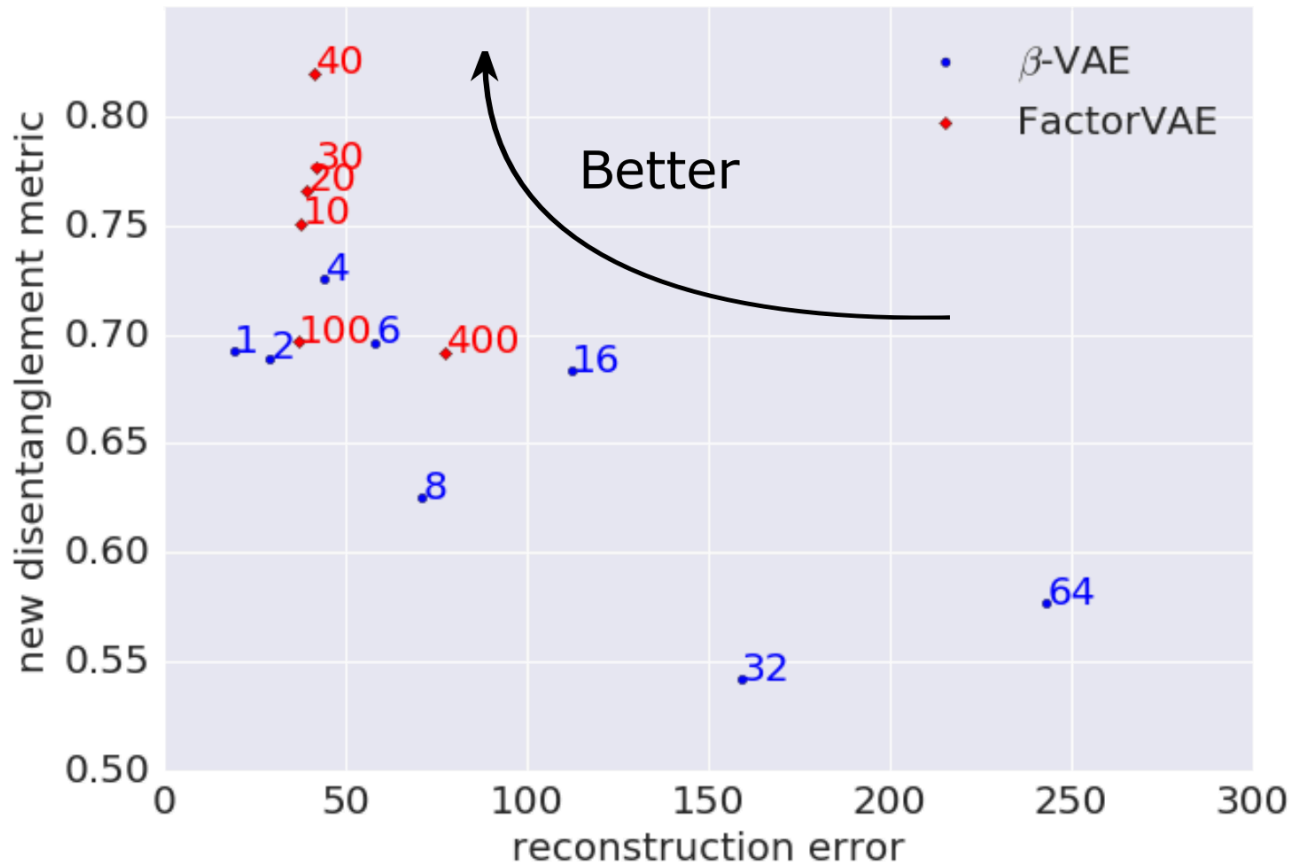
FactorVAE

Idea: β -VAE optimizes the two terms together, FactorVAE separates them



FactorVAE

Idea: β -VAE optimizes the two terms together, FactorVAE separates them



Identifiable VAE (i-VAE)

$$\text{VAE } p(z) \longrightarrow p(z|u) \quad \text{i-VAE}$$

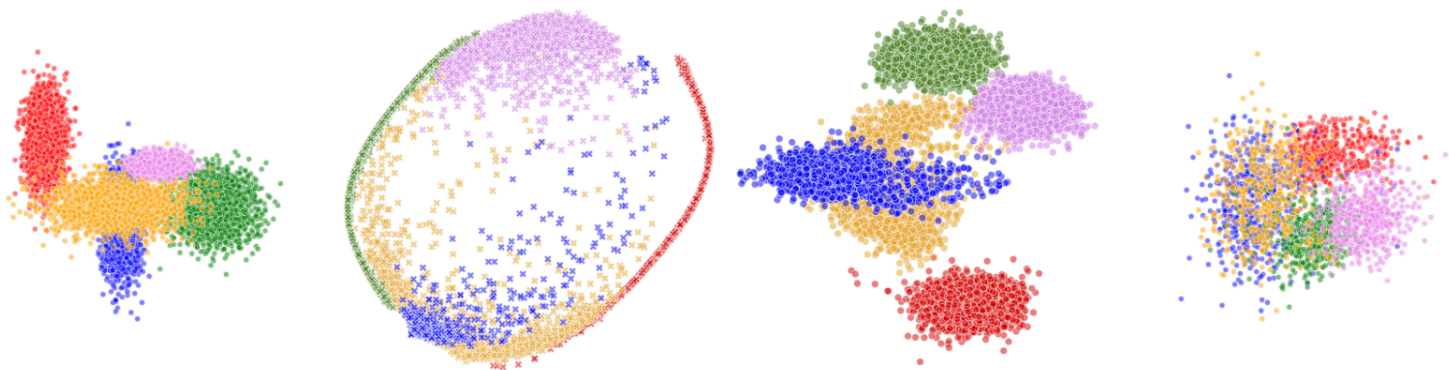
Main Assumption: A conditionally factorized prior distribution over the latent variables $p_\theta(z|u)$, where u is an additionally observed variable
And the data generation stage is a additive noise model $x = f(z) + \epsilon$

$p(z|u)$ is conditionally factorial

$$p(z|u) = \prod_{i=1}^n p(z_i|u),$$

$$\text{Maximize } ELBO = E_D(E_{q_\phi(z|x,u)} \log p_\theta(x|z, u) - KL(q_\phi(z|x, u) || p(z|u)))$$

Identifiable VAE (i-VAE)



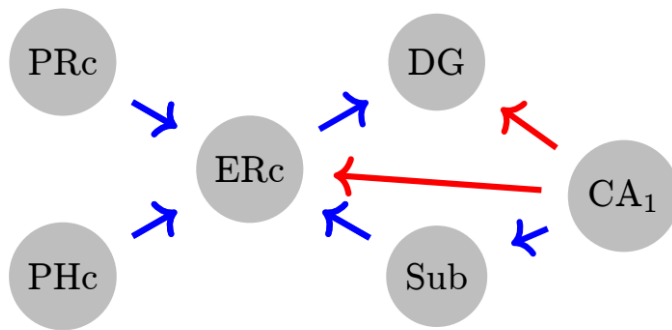
(a) $p_{\theta^*}(\mathbf{z}|\mathbf{u})$ (b) $p_{\theta^*}(\mathbf{x}|\mathbf{u})$ (c) $p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{u})$ (d) $p_{\text{VAE}}(\mathbf{z}|\mathbf{x})$

i-VAE for Causal Discovery

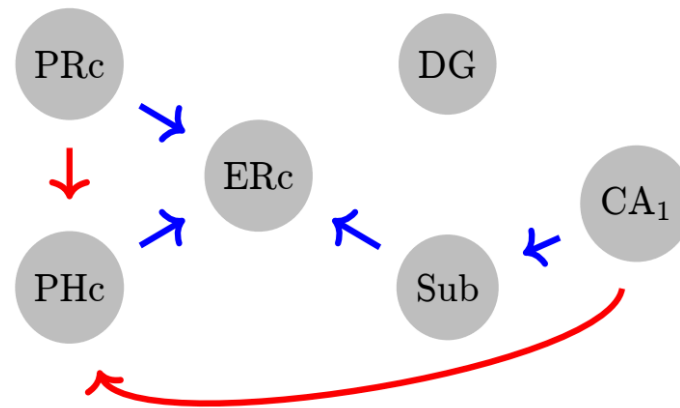
hippocampal fMRI data

Blue: Correct (feasible given anatomical connectivity)

Red: Incorrect (incompatible with anatomical structure)



(a) iVAE 2021



(b) TCL 2016.

Q: Which method is better?

Part III

Summary

Learning Outcomes

- Appreciate how causal learning differs from statistical learning
- Understand the tasks of causal inference and causal discovery
- Be able to describe ICA and its identifiability
- Be able to connect nonlinear ICA and the disentanglement problem in generative models
- Know what β -VAE, FactorVAE, I-VAE are

Motivation

Trustworthy ML