

Causality-inspired ML: what can causality do for ML?

Sara Magliacane University of Amsterdam MIT-IBM Watson AI Lab







- Real-world ML needs to deal with:
 - **Biased data** (fairness, selection bias, generalization)
 - Heterogeneous data, small samples, missing/corrupted data, not iid
 - Actionable insights (decisions cannot be made on correlations)

- Real-world ML needs to deal with:
 - Biased data (fairness, selection bias, generalization)
 - Heterogeneous data, small samples, missing/corrupted data, not iid
 - Actionable insights (decisions cannot be made on correlations)
- Causal inference can help with some of these questions:
 - Systematic data fusion and reuse with biased data, heterogenous, not iid data
 - A systematic way to **extract actionable insights**

- **Real-world ML** needs to deal with:
 - **Biased data** (fairness, selection bias, generalization)
 - Heterogeneous data, small samples, missing/corrupted data, not iid
 - Actionable insights (decisions cannot be made on correlations)
- **Causal inference** can help with some of these questions:
 - Systematic data fusion and reuse with biased data, heterogenous, not iid data
 - A systematic way to **extract actionable insights**

Many excellent books and courses, my fav is:

https://stat.ethz.ch/lectures/ss21/causality.php

- **Real-world ML** needs to deal with:
 - **Biased data** (fairness, selection bias, generalization)
 - Heterogeneous data, small samples, missing/corrupted data, not iid
 - Actionable insights (decisions cannot be made on correlations)
- **Causal inference** can help with some of these questions:
 - E.g. too many Systematic data fusion and reuse with biased data experiments to fully identify the graph • A systematic way to extract actionable insights
- "Full" causality can be not necessary or too expensive ->



- **Real-world ML** needs to deal with:
 - **Biased data** (fairness, selection bias, generalization)
 - Heterogeneous data, small samples, missing/corrupted data, not iid
 - Actionable insights (decisions cannot be made on correlations)
- **Causal inference** can help with some of these questions:
 - Systematic data fusion and reuse with biased data, heterogenous, not iid data
 - A systematic way to **extract actionable insights**
- "Full" causality can be not necessary or too expensive -> Gausality Inspired



- **Real-world ML** needs to deal with:
 - **Biased data** (fairness, selection bias, generalization)
 - Heterogeneous data, small samples, missing/corrupted data, not iid
 - Actionable insights (decisions cannot be made on correlations)
- Causal inference can help with some of these questions
 - Systematic data fusion and reuse
 - A systematic way to extract actionable man

In this talk: example in id data domain adaptation, but lots of related work



- Transfer learning:
 - How can I predict what happens when the distribution changes?





- Transfer learning:
 - How can I predict what happens when the distribution changes?



- Transfer learning:
 - How can I predict what happens when the distribution changes?



- Causal inference:
 - How can I predict what happens when the distribution changes after an intervention?
 - Perfect intervention: do-calculus [Pearl, 2009]
 - X is independent of its parents
 - Soft intervention on X:
 - Change of P(X| parents)

- Transfer learning:
 - How can I predict what ha cha when the distribution chan

Very general - can model also changes in distribution that are not from "real" interventions









A term of its parents

- Soft intervention on X:
 - Change of P(X| parents)

- Transfer learning:
 - How can I predict what ha ch when the distribution chan











Very general - can model also changes in distribution that are not from "real" interventions

• X is independent of its parents

- Soft intervention on X:
 - Change of P(X| parents)

Not a new idea!

On Causal and Anticausal Learning ICML 2012

Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun ZhangFIRST.LAST@TUE.MPG.DEMax Planck Institute for Intelligent Systems, Spemannstrasse, 72076 Tübingen, GermanyFIRST.LAST@TUE.MPG.DE

Joris Mooij

Institute for Computing and Information Sciences, Radboud University, Nijmegen, The Netherlands

Abstract

We consider the problem of function estimation in the case where an underlying causal model can be inferred. This has implications for popular scenarios such as covariate shift, concept drift, transfer learning and semi-supervised learning. We argue that causal knowledge may facilitate some approaches for a given problem, and rule out others. In particular, we formulate a hypothesis for when semi-supervised learning can help, and corroborate it with empirical results.

```
J.MOOIJ@CS.RU.NL
```

for causal inference in the machine learning community.

An example illustrating the difference between the statistical and the causal point of view is the correlation between the frequency of storks and the human birth rate (Matthews, 2000). We may be able to train a good predictor of the birth rate which uses the frequency of storks (along with other features) as an input. However, if politicians asked us whether one could boost the birth rate by increasing the number of storks, we would have to tell them that this kind of *intervention* is not covered by the standard i.i.d. assumption of statistical learning. In practice, however, interventions can be relevant, distributions may shift over time, and we might want to combine data recorded under different

Causality allows us to reason systematically about distribution shifts

Causality allows us to reason systematically about distribution shifts

On Causal and Anticausal Learning

| Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang | FIRST.LAST@TUE.MPG.DE | | | | | |
|--|-----------------------|--|--|--|--|--|
| Max Planck Institute for Intelligent Systems, Spemannstrasse, 72076 Tübingen, Germany | | | | | | |
| Joris Mooij | J.MOOIJ@CS.RU.NL | | | | | |
| nstitute for Computing and Information Sciences, Radboud University, Nijmegen, The Netherlands | | | | | | |
| | | | | | | |

Domain Adaptation as a Problem of Inference on Graphical Models

Kun Zhang^{1*}, Mingming Gong^{2*}, Petar Stojanov³ Biwei Huang¹, Qingsong Liu⁴, Clark Glymour¹ ¹ Department of philosophy, Carnegie Mellon University ² School of Mathematics and Statistics, University of Melbourne ³ Computer Science Department, Carnegie Mellon University, ⁴ Unisound AI Lab kunz1@cmu.edu, mingming.gong@unimelb.edu.au, liuqingsong@unisound.com {pstojano, biweih, cg09}@andrew.cmu.edu

Anchor regression: heterogeneous data meet causality

Dominik Rothenhäusler, Nicolai Meinshausen, Peter Bühlmann and Jonas Peters

Invariant Risk Minimization

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, David Lopez-Paz

J. R. Statist. Soc. B (2016) 78, Part 5, pp. 947-1012

Jonas Peters

Invariant Models for Causal Transfer Learning

Mateo Rojas-Carulla Max Planck Institute for Intelligent Systems Tübingen, Germany

Department of Engineering Univ. of Cambridge, United Kingdom

Bernhard Schölkopf Max Planck Institute for Intelligent Systems Tübingen, Germany

Richard Turner Department of Engineering Univ. of Cambridge, United Kingdom

Jonas Peters^{*} Department of Mathematical Sciences Univ. of Copenhagen, Denmark

Invariance, Causality and Robustness

Peter Bühlmann [†] Seminar for Statistics, ETH Zürich

Causal inference by using invariant prediction: identification and confidence intervals

Max Planck Institute for Intelligent Systems, Tübingen, Germany, and Eidgenössiche Technische Hochschule Zürich, Switzerland

and Peter Bühlmann and Nicolai Meinshausen

Eidgenössiche Technische Hochschule Zürich, Switzerland

BS@TUEBINGEN.MPG.DE

RET26@CAM.AC.UK

JONAS.PETERS@MATH.KU.DK

Counterfactual Invariance to Spurious Correlations: Why and How to Pass Stress Tests

Victor Veitch^{1,2}, Alexander D'Amour¹, Steve Yadlowsky¹, and Jacob Eisenstein¹

> ¹Google Research ²University of Chicago

Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions

Sara Magliacane IBM Research* sara.magliacane@gmail.com

Tom Claassen

Radboud University Nijmegen

tomc@cs.ru.nl

Stephan Bongers University of Amsterdam srbongers@gmail.com

Philip Versteeg University of Amsterdam p.j.j.p.versteeg@uva.nl

Thijs van Ommen

University of Amsterdam thijsvanommen@gmail.com

Joris M. Mooij University of Amsterdam j.m.mooij@uva.nl

A Causal View on Robustness of Neural Networks

2018 Neyman Lecture

Cheng Zhang Microsoft Research Cheng.Zhang@microsoft.com

Kun Zhang Carnegie Mellon University kunz10cmu.edu

Yingzhen Li * Microsoft Research Yingzhen.Li@microsoft.com

and many many more....

MR597@CAM.AC.UK

Causality allows us to reason systematically about distribution shifts, e.g. through graphs





Domain Adaptation as a Problem of Inference on Graphical Models



Anchor regression: heterogeneous data meet causality



J. R. Statist. Soc. B (2016) 78, Part 5, pp. 947–1012



Invariant Models for Causal Transfer Learning



Causal inference by using invariant prediction: identification and confidence intervals

Invariance, Causality and Robustness



Counterfactual Invariance to Spurious Correlations: Why and How to Pass Stress Tests



Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions



A Causal View on Robustness of Neural Networks



and many more....

Causality allows us to reason systematically about distribution shifts, e.g. through graphs



A description of domain adaptation tasks:

Supervised multi-source domain adaptation

| | X1 | X2 | Х3 |
|---|-----------|-----------|------|
| - | 1200 | 1000 | 1500 |
| | 1201 | 800 | 1500 |
| | 1195 | 200 | 1499 |
| | | | |
| | 2000 | 600 | 3000 |
| | 2190 | 450 | 3000 |
| | 2000 | 200 | 2999 |
| | | | |
| | 1200 | 1000 | 1500 |
| | 1201 | 800 | 1500 |
| | 1195 | 200 | 1499 |
| | 1340 | 900 | 1498 |

• Estimate \hat{f} in Y = $\hat{f}(X1, X2, X3, X4)$ from source domains and few labels in target domain



A description of domain adaptation tasks:

Unsupervised multi-source domain adaptation

| | X1 | X2 | Х3 |
|---|-----------|-----------|------|
| - | 1200 | 1000 | 1500 |
| | 1201 | 800 | 1500 |
| | 1195 | 200 | 1499 |
| | | | |
| | 2000 | 600 | 3000 |
| | 2190 | 450 | 3000 |
| | 2000 | 200 | 2999 |
| | | | |
| | 1200 | 1000 | 1500 |
| | 1201 | 800 | 1500 |
| | 1195 | 200 | 1499 |
| | 1340 | 900 | 1498 |

• Estimate \hat{f} in Y = $\hat{f}(X1, X2, X3, X4)$ from source domains and by exploiting the knowledge of the change from the unlabelled data in target





A description of domain adaptation tasks:

Domain generalisation: required to work under any intervention

| | X1 | X2 | Х3 |
|--|------|-----------|------|
| | ? | ? | ? |
| | ? | ? | ? |
| | ? | ? | ? |
| | | | |
| | 2000 | 600 | 3000 |
| | 2190 | 450 | 3000 |
| | 2000 | 200 | 2999 |
| | | | |
| | 1200 | 1000 | 1500 |
| | 1201 | 800 | 1500 |
| | 1195 | 200 | 1499 |
| | 1340 | 900 | 1498 |

• Estimate \hat{f} in Y = $\hat{f}(X1, X2, X3, X4)$ from source domains, no idea about what happens in the target





| | X1 | X2 | Υ |
|--------|-----------|-----------|---|
| Normal | 0.1 | 2 | 0 |
| Normal | 0.2 | 3 | 0 |
| Normal | 1.1 | 2 | 1 |
| Normal | 0.1 | 3 | 0 |







| | X1 | X2 | Y |
|--------|-----------|-----------|---|
| Gene A | 3.1 | 2 | ? |
| Gene A | 3.2 | 3 | ? |
| Gene A | 4 | 2 | ? |
| Gene A | 3.2 | 3 | ? |

| D | | (1 | X2 | Y |
|------|------|-----------|----|---|
| Norm | al O |).1 | 2 | 0 |
| Norm | al C | .2 | 3 | 0 |
| Norm | al 1 | .1 | 2 | 1 |
| Norm | al O | .1 | 3 | 0 |



| | D | X1 | X2 | Y | |
|------|--------|-----------|-----------|---|--|
| | Gene A | 3.1 | 2 | ? | |
| | Gene A | 3.2 | 3 | ? | |
| | Gene A | 4 | 2 | ? | |
| Gene | Gene A | 3.2 | 3 | ? | |

Add a variable D to represent the domain

| D | X1 | X2 | Υ | |
|------------|-----------|-----------|---|--|
| Normal | 0.1 | 2 | 0 | |
| Normal | 0.2 | 3 | 0 | |
| Normal | 1.1 | 2 | 1 | |
| Normal | 0.1 | 3 | 0 | |
| Gene A | 3.1 | 2 | ? | |
| Gene A | 3.2 | 3 | ? | |
| Gene A | 4 | 2 | ? | |
| Gene A | 3.2 | 3 | ? | |
| | | | | |

- Add a variable D to represent the domain
- Consider the data as coming from a single distribution P(X,Y, D)

| D | X1 | X2 | Υ | |
|------------|-----------|-----------|---|--|
| Normal | 0.1 | 2 | 0 | |
| Normal | 0.2 | 3 | 0 | |
| Normal | 1.1 | 2 | 1 | |
| Normal | 0.1 | 3 | 0 | |
| Gene A | 3.1 | 2 | ? | |
| Gene A | 3.2 | 3 | ? | |
| Gene A | 4 | 2 | ? | |
| Gene A | 3.2 | 3 | ? | |
| | | | | |

- Add a variable D to represent the domain
- Consider the data as coming from a single distribution P(X,Y, D)



• We can represent P(X,Y, D) with a (possibly unknown) causal graph



| D | X1 | X2 | Υ | |
|------------|-----------|-----------|---|--|
| Normal | 0.1 | 2 | 0 | |
| Normal | 0.2 | 3 | 0 | |
| Normal | 1.1 | 2 | 1 | |
| Normal | 0.1 | 3 | 0 | |
| Gene A | 3.1 | 2 | ? | |
| Gene A | 3.2 | 3 | ? | |
| Gene A | 4 | 2 | ? | |
| Gene A | 3.2 | 3 | ? | |
| | | | | |

- Add a variable D to represent the domain
- Consider the data as coming from a single distribution P(X,Y, D)



• We can represent P(X,Y, D) with a (possibly unknown) causal graph

We can still use d-separations/ conditional independences to reason about invariances



Bayesian networks: d-separation [Pearl 2009]

$$X \longrightarrow Z \longrightarrow Y$$

* Causal Markov assumption: $X \perp A Y \mid Z \implies X \perp Y \mid Z$ Causal faithfulness assumption: $X \perp Y \mid Z \implies X \perp_d Y \mid Z$

 Given a causal graph, d-separation is a graphical criterion that (under standard conditions^{*}) allows us to read conditional independences

• A path between node i and node j is a sequence of distinct nodes (i, \ldots, j) such that each two consecutive nodes are adjacent





 A path between node i and node j is a sequence of distinct nodes (i, \ldots, j) such that each two consecutive nodes are adjacent





- A path between node i and node j is a sequence of distinct nodes
 (i, ..., j) such that each two consecutive nodes are adjacent
- A collider k on a path $\pi = (i, ..., j)$ is a non-endpoint node $(k \neq i, j)$ s.t. the path π contains $\rightarrow k \leftarrow$, the other nodes are **non-colliders**





- A path between node i and node j is a sequence of distinct nodes
 (i, ..., j) such that each two consecutive nodes are adjacent
- A collider k on a path $\pi = (i, ..., j)$ is a non-endpoint node $(k \neq i, j)$ s.t. the path π contains $\rightarrow k \leftarrow$, the other nodes are **non-colliders**





d-separation: blocked paths

- - There is a non-collider on the path that is in A, or
- Otherwise it is **active**

• A path between i and j is blocked by $A \subseteq V$ at least one condition holds: • There is a collider k on the path, but $k \notin A$ and $Desc(k) \cap A = \emptyset$

> Descendants of k (i.e. nodes that can be reached from k by following directed edge)



d-separation: blocked paths - example 1

- A path between i and j is blocked by $A \subseteq V$ at least one condition holds: • There is a non-collider on the path that is in A, or
- - There is a collider k on the path, but $k \notin A$ and $Desc(k) \cap A = \emptyset$
- Otherwise it is **active**





If $3 \in \mathbf{A}$, the path is **blocked**, otherwise it is active

d-separation: blocked paths - example 2

- A path between i and j is blocked by $A \subseteq V$ at least one condition holds: • There is a non-collider on the path that is in A, or
- There is a collider k on the path, but $k \notin A$ and $Desc(k) \cap A = \emptyset$
- Otherwise it is active





collider non-collider

If $1 \in \mathbf{A}$, the path is **blocked** OR

If $3 \notin A$ and $2 \notin A$, the path is blocked



d-separation: definition

- - We denote d-separation as $i \perp j \land A$
- Otherwise we say they are d-connected
 - We denote d-connection as $i \not \perp_d j \mid \mathbf{A}$
- Under standard assumptions: $X \perp A \mid Z \iff X \perp Y \mid Z$
- Demo: http://www.dagitty.net/learn/dsep/index.html

• Nodes i and j is d-separated by $A \subseteq V$ if all paths between i, j are blocked

Note: d-separation is symmetric

Structural causal model - domain/environment variable

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0, 1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0, 1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0, 1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_Y \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$

Structural causal model - domain/environment variable

$$\begin{cases} \epsilon_{1}, \epsilon_{2}, \epsilon_{3}, \epsilon_{Y} \sim \mathcal{N}(0, 1) \\ X_{1} = 10 + \epsilon_{1} \\ Y = 3X_{1} + \epsilon_{Y} \\ X_{3} = 2Y + 0.1\epsilon_{3} \end{cases} \qquad D = 0$$

$$\begin{cases} \epsilon_{1}, \epsilon_{2}, \epsilon_{3}, \epsilon_{Y} \sim \mathcal{N}(0, 1) \\ X_{1} = 10 + \epsilon_{1} \\ Y = 3X_{1} + \epsilon_{Y} \\ X_{2} = 1 \\ X_{3} = 2Y + 0.1\epsilon_{3} \end{cases} \qquad D = 1$$

$$\begin{cases} \epsilon_{1}, \epsilon_{2}, \epsilon_{3}, \epsilon_{Y} \sim \mathcal{N}(0, 1) \\ X_{1} = 10 + \epsilon_{1} \\ Y = 3X_{1} + \epsilon_{Y} \\ X_{2} = \begin{cases} -2Y + \epsilon_{2} \text{ if } D = 0 \\ 1 & \text{if } D = 1 \\ 10Y + \epsilon_{Y} & \text{if } D = 2 \end{cases}$$

$$\begin{cases} \epsilon_{1}, \epsilon_{2}, \epsilon_{3}, \epsilon_{Y} \sim \mathcal{N}(0, 1) \\ X_{1} = 10 + \epsilon_{1} \\ Y = 3X_{1} + \epsilon_{Y} \end{cases}$$

$$D = 2$$

 $X_2 = 10Y + \epsilon_Y$ $X_3 = 2Y + 0.1\epsilon_3$


Domain adaptation example





Domain adaptation example - X1

P(Y|X1) is invariant X1

| d | x1 | У | x2 | x3 |
|---|-----------|-----------|------------|-----------|
| 0 | 8.973763 | 26.130494 | -51.648475 | 52.330948 |
| 0 | 10.428340 | 31.894998 | -64.373356 | 63.802704 |
| 0 | 8.911484 | 25.166962 | -52.313502 | 50.279162 |
| 0 | 9.841798 | 29.783299 | -60.419296 | 59.539914 |
| 0 | 8.969118 | 27.660573 | -55.075839 | 55.327185 |







| х3 | x2 | У | x1 |
|-----------|----|-----------|-----------|
| 57.475345 | 1 | 28.696601 | 9.941015 |
| 51.275390 | 1 | 25.715927 | 8.762380 |
| 56.884332 | 1 | 28.407387 | 9.636201 |
| 62.686789 | 1 | 31.370200 | 10.875069 |
| 62.388444 | 1 | 31.253540 | 10.023968 |
| | | | |

| d | x1 | У | x2 | x3 |
|---|-----------|-----------|------------|-----------|
| 2 | 9.671277 | 26.556214 | 265.034283 | 53.338139 |
| 2 | 9.613139 | 27.120226 | 270.746784 | 54.340341 |
| 2 | 10.718335 | 29.589532 | 295.318526 | 59.291053 |
| 2 | 9.002388 | 26.629254 | 264.942583 | 53.340389 |
| 2 | 9.289340 | 29.030355 | 289.747562 | 58.098312 |



Domain adaptation example - X1











print("Mean squared error predicting Y in environment 2 based on model learnt in environment 0 from X1", mean_squared_error(Y_2,est_Y_2))

Domain adaptation example - X2



```
sns.scatterplot(data = df, x="x2", y="y", hue="d")
X2_0 = df_0["x2"].values.reshape(-1, 1)
X2_2 = df_2["x2"].values.reshape(-1, 1)
model = LinearRegression().fit(X2_0, Y_0)
est_Y_2 = model.predict(X2_2)
```

Mean squared error predicting Y in environment 2 based on model learnt in environment 0 from X2 30518.374428658524





Separating features intuition - X1



 $P(X_{1}, Y_{1}, X_{2}, X_{3}, D)$

P(Y|X1) is invariant

$P(Y|X_{1}, D=0) = P(Y|X_{1}, D=1) = P(Y|X_{1}, D=2)$ $= P(Y|X_{1})$

Ly this is the if YILD) Xy in true graph d-separation [Pearl 1988]





Separating features intuition - X2



P(4 | Xz) is not invariant $P(Y|X_2, D=0) \neq P(Y|X_2, D=4) \neq P(Y|X_2, D=2)$ 4 this means YKD X2

 $P(X_{1}, Y_{1}, X_{2}, X_{3}, D)$



Separating features intuition - X2



 $P(X_{1}, Y_{1}, X_{2}, X_{3}, D)$

P(4 | Xz) is not invariant $P(Y|X_2, D=0) \neq P(Y|X_2, D=4) \neq P(Y|X_2, D=2)$ 4 this means YXD X2 Look for fratures SSX YLdD **Separating features [Magliacane et al 2018]**





Which variables d-separate Y from D now?



 $P(X_{1}, Y_{1}, X_{2}, X_{3}, D)$

| X1 | X2 | Х3 | Υ |
|------|-----------|------|-------|
| ? | ? | ? | ? |
| ? | ? | ? | ? |
| ? | ? | ? | ? |
| | | | |
| 2000 | 600 | 3000 | -0.21 |
| 2190 | 450 | 3000 | -0.16 |
| 2000 | 200 | 2999 | -0.16 |
| | | | |
| 1200 | 1000 | 1500 | -0.17 |
| 1201 | 800 | 1500 | -0.14 |
| 1195 | 200 | 1499 | -0.07 |
| 1340 | 900 | 1498 | -0.14 |
| | | | |

Intervention on every variable except Y = domain generalisation



Common misconceptions: 1. An invariant feature need not be causal



• Y|X1,X2 is invariant \implies invariant features are not necessarily parents of Y



 $Y \perp D \mid X_1$

Common misconceptions: 1. An invariant feature need not be causal



Invariant feature across "many different datasets" is not enough in general to find causal parents, need more assumptions

• Y|X1,X2 is invariant \implies invariant features are not necessarily parents of Y

Common misconceptions: 1. An invariant feature need not be causal $Y \perp D \mid X_1$ X2 $Y \perp D \mid X_1, X_2$

- YX1,X2 is invariant \implies invariant features are not necessarily parents of Y
- Invariant Causal Prediction [Peters et al. 2016] under causal sufficiency: $\{X_1, X_2\} \cap \{X_1\} = \{X_1\}$ Pa(Y)

$$\mathbf{S}^* = \bigcap_{Y \perp L D \mid \mathbf{S}} \mathbf{S} \subseteq$$

Common misconceptions: 1. An invariant feature need not be causal $Y \perp D \mid X_1$ X2 $Y \perp D \mid X_1, X_2$

$$\begin{array}{c} \hline D \end{array} \longrightarrow \begin{array}{c} \hline X1 \end{array} \longrightarrow \begin{array}{c} \hline Y \end{array} \\ \hline \end{array} \end{array}$$

- Y|X1,X2 is invariant \implies invariant features are not necessarily parents of Y
- Invariant Causal Prediction [Peters et al. 2016] under causal sufficiency: $\mathbf{S}^* = \bigcap \ \mathbf{S} \subseteq Pa(Y) \qquad \{X_1, X_2\} \cap \{X_1\} = \{X_1\}$





→ $\{X_1\} \cap \{X_2\} \cap \{X_1, X_2\} = \emptyset$

Common misconception 2: Parents are not enough under latent confounding

Y



• Y|X1 is invariant, Y|X2 is not

 $\begin{array}{c} Y \perp D \mid X_1 \\ Y \perp D \mid X_2 \\ Y \perp D \mid X_2 \end{array}$

Common misconception 2: Parents are not enough under latent confounding



• Y|X1 is invariant, Y|X2 is not

 $Y \perp D \mid X_1$ $Y \perp D \mid X_{2}$ $Y \perp D \mid X_1, X_2$

Even if we knew all the parents, under latent confounding this wouldn't necessarily help transfer

Common misconception 2: Parents are not enough under latent confounding $Y \perp D \mid X_1$ Н



• Y|X1 is invariant, Y|X2 is not

• Conclusion: causality (e.g. using the causal parents, learning the

Even if we knew all the parents, under latent confounding this wouldn't necessarily help transfer

 $Y \perp D \mid X_2$

 $Y \perp D \mid X_1, X_2$

complete causal graph) is neither necessary or sufficient* for transfer

Desiderata for a causality inspired domain adaptation method

- X, Y and changes can be represented by an unknown causal graph
- Allow for latent confounders
- Avoid parametric assumptions, allow for heterogeneous effects across domains

Desiderata for a causality inspired domain adaptation method

- X, Y and changes can be represented by an unknown causal graph
- Allow for latent confounders
- Avoid parametric assumptions, allow for heterogeneous effects across domains
- Instead of modeling changes between each domain, distinguish the change between the mixture of sources and the target

Desiderata for a causality inspired domain adaptation method

- X, Y and changes can be represented by an unknown causal graph
- Allow for latent confounders
- Avoid parametric assumptions, allow for heterogeneous effects across domains
- Instead of modeling changes between each domain, distinguish the change between the mixture of sources and the target
- Avoid common assumption that if Y T(X) is invariant across multiple source domains, then Y T(X) is invariant also in the target domain
 - This also (implicitly) assumed by methods based on the idea that invariance => causality

Causal domain adaptation problem [Magliacane et al. 2018]

Unsupervised multi-source domain adaptation

X2

2

3

2

3

• We interpret the change in the target domain as a soft intervention

Y

0

0

0

• We assume Y cannot be intervened upon directly - P(Y) can still change



| | X1 | X2 | Υ |
|--------|-----|----|---|
| Gene A | 3.1 | 2 | 1 |
| Gene A | 3.2 | 3 | 1 |
| Gene A | 4 | 1 | 1 |
| Gene A | 3.2 | 3 | 0 |

X1

0.1

0.2

1.1

0.1

Normal

Normal

Normal

Normal



| T | | X1 | X2 | Υ |
|---|--------|-----------|-----------|---|
| | Gene B | 0.2 | 1 | ? |
| | Gene B | 0.3 | 1 | ? |
| | Gene B | 0.3 | 2 | ? |
| | Gene B | 04 | 1 | ? |

Causal domain adaptation Multiple context variable [Magliacane et al. 2018 C1, C2...

- Unsupervised multi-source domain adaptation
- We interpret the change in the target domain as a **soft intervention**

Normal

Normal

Normal

Normal

X1

0.1

0.2

1.1

0.1



| | X1 | X2 | Υ |
|--------|-----|----|---|
| Gene A | 3.1 | 2 | 1 |
| Gene A | 3.2 | 3 | 1 |
| Gene A | 4 | 1 | 1 |
| Gene A | 3.2 | 3 | 0 |



Y

0

0

0

X2

2

3

2

3

• We assume Y cannot be intervened upon directly - P(Y) can still change

| T | | X1 | X2 | Y |
|---|--------|-----------|-----------|---|
| | Gene B | 0.2 | 1 | ? |
| | Gene B | 0.3 | 1 | ? |
| | Gene B | 0.3 | 2 | ? |
| | Gene B | 04 | 1 | 2 |



Causal domain adaptation problem [Magliacane et al. 2018]

- Unsupervised multi-source domain adaption
- We interpret the change in the target domain as a soft intervention

Normal

Normal

Normal

X1

0.1

0.2

1.1

X2

2

3

2

Y

0

0



| Normal | 0.1 | 3 | 0 |
|--------|-----------|----|---|
| | X1 | X2 | Y |
| Gene A | 3.1 | 2 | 1 |
| Gene A | 3.2 | 3 | 1 |
| Gene A | 4 | 1 | 1 |
| Gene A | 3.2 | 3 | 0 |



C1 = 1

• We assume Y cannot be intervened upon directly - P(Y) can still change

| T | | X1 | X2 | Υ |
|----------|--------|-----------|-----------|---|
| | Gene B | 0.2 | 1 | ? |
| | Gene B | 0.3 | 1 | ? |
| | Gene B | 0.3 | 2 | ? |
| | Gene B | 04 | 1 | 2 |

Causal domain adaptation problem [Magliacane et al. 2018]

Unsupervised multi-source domain adaptation

X2

2

• We interpret the change in the target domain as a soft intervention

0

• We assume Y cannot be intervened upon directly - P(Y) can still change





X1

0.1

Normal





Joint Causal Inference [Mooij et al. 2020]

- \bullet
- **disentangle** changes in distribution across the datasets



We represent jointly different distributions as an **unknown single causal graph**

Instead of a single domain variable, we add several context variables so we can

• If we know nothing about the changes in the datasets, we use indicator variables

| C 1 | C 2 | X1 | X2 | Υ |
|------------|------------|-----------|-----------|---|
| 0 | 0 | 0.1 | 2 | 0 |
| 0 | 0 | 0.2 | 3 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 0 | 0 | 0.1 | 3 | 0 |
| 1 | 0 | 3.1 | 2 | 1 |
| 1 | 0 | 3.2 | 3 | 1 |
| 1 | 0 | 4 | 1 | 1 |
| 1 | 0 | 3.2 | 3 | 0 |
| 0 | 1 | 0.2 | 1 | ? |
| 0 | 1 | 0.3 | 1 | ? |
| 0 | 1 | 0.3 | 2 | ? |
| 0 | 1 | 0.4 | 1 | ? |

Joint Causal Inference [Mooij et al. 2020]

| C1 | C 2 | X1 | X2 | Y | |
|----|------------|-----------|----|---|-------|
| 0 | 0 | 0.1 | 1 | 0 | C_1 |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | X_1 |
| 1 | 0 | 3.1 | 2 | 1 | 1 |
| 1 | 0 | 3.2 | 3 | 1 | X_1 |
| 1 | 0 | 4 | 3 | 1 | L |
| 0 | 1 | 0.2 | 0 | 0 | |
| 0 | 1 | 0.3 | 0 | 1 | |
| 0 | 1 | 0.3 | 1 | 0 | |

 We can learn an equivalence class of the unknown single causal graph using conditional independence tests on systematically pooled data

We treat context variables as normal variables that we know are uncaused

 $\parallel Y \mid X_1$ $\perp X_2 \mid Y, C_2$





Causal domain adaptation: separating features

variable C1 representing the target domain



{X1} is a separating feature set, {X1, X2} could lead to arbitrary large error

• Separating features: sets of features that d-separate Y from the context

Aka stable features, invariant features etc.

• Idea: we could test the conditional independence in the data $Y \perp C_1 | X_1? \qquad Y \perp C_1 | X_2?$

- Idea: we could test the conditional independence in the data $Y \perp C_1 \mid X_1? \qquad Y \perp C_1 \mid X_2?$
- **Problem:** Y is always missing when C1=1, so we cannot test these

| C1 | C2 | X1 | X2 | Υ |
|----|-----------|-----------|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 1 | 0 | 3.1 | 2 | ? |
| 1 | 0 | 3.2 | 3 | ? |
| 1 | 0 | 4 | 3 | ? |
| 0 | 1 | 0.2 | 0 | 0 |
| 0 | 1 | 0.3 | 0 | 1 |
| 0 | 1 | 0.3 | 1 | 0 |

- Idea: we could test the conditional independence in the data $Y \perp C_1$ $Y \perp C$
- **Problem:** Y is always missing when C1=1, so we cannot test these \bullet

| C1 | C2 | X1 | X2 | Υ |
|----|----|-----------|----|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 1 | 0 | 3.1 | 2 | ? |
| 1 | 0 | 3.2 | 3 | ? |
| 1 | 0 | 4 | 3 | ? |
| 0 | 1 | 0.2 | 0 | 0 |
| 0 | 1 | 0.3 | 0 | 1 |
| 0 | 1 | 0.3 | 1 | 0 |

Invariant Models for Causal Transfer Learning

Mateo Rojas-Carulla Max Planck Institute for Intelligent Systems Tübingen, Germany

Department of Engineering Univ. of Cambridge, United Kingdom

Bernhard Schölkopf Max Planck Institute for Intelligent Systems Tübingen, Germany

Richard Turner Department of Engineering Univ. of Cambridge, United Kingdom

Jonas Peters^{*} Department of Mathematical Sciences Univ. of Copenhagen, Denmark

MR597@CAM.AC.UK

BS@TUEBINGEN.MPG.DE

RET26@CAM.AC.UK

JONAS.PETERS@MATH.KU.DK

Idea: Separating features in sources are also separating in target

 $Y \perp \!\!\!\perp C_2 \mid X_1 \implies Y \perp \!\!\!\perp C_1 \mid X_1$

Separating features in sources are also separating in target - counterexample



c1 = epsilon_c1 c2 = epsilon_c2 x1 = epsilon_x1 + 10 + 2 * c2 y = 4 * x1 + epsilon_y x2 = - 2 * y * c1 + epsilon_x2 + 10 * y *(1-c1)



• Idea: we could test the conditional independence in the data $Y \perp C_1 \mid X_1? \qquad Y \perp C_1 \mid X_2?$

• **Problem:** Y is always missing when C1=1, so we cannot test these

| C1 | C2 | X1 | X2 | Υ |
|----|-----------|-----|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 1 | 0 | 3.1 | 2 | ? |
| 1 | 0 | 3.2 | 3 | ? |
| 1 | 0 | 4 | 3 | ? |
| 0 | 1 | 0.2 | 0 | 0 |
| 0 | 1 | 0.3 | 0 | 1 |
| 0 | 1 | 0.3 | 1 | 0 |

 $X_1 \perp X_2$ $X_1 \perp C_1$ $X_1 \perp X_2 \mid C_1$ $X_1 \perp X_2 \mid Y, C_1 = 0$

Idea: Can we use all other in/dependences?

• • •



Assumptions [Magliacane et al. 2018]

- (Joint Causal Inference)
- We assume Y cannot be intervened upon directly

• We assume that there exists an acyclic causal graph that fits all the data

Assumptions [Magliacane et al. 2018]

- We assume that there exists an acyclic causal graph that fits all the data (Joint Causal Inference)
- We assume Y cannot be intervened upon directly
- We assume no extra dependences involving Y in target domain C1=1 $A, D, \mathbf{B} \in \mathbf{V} \setminus \{Y, C_1\}$ $Y \perp \!\!\!\perp A \mid \mathbf{B}, C_1 = 0 \implies Y \perp \!\!\!\perp A \mid \mathbf{B}, C_1 = 1$ $A \perp \!\!\!\perp D \mid \mathbf{B}, Y, C_1 = 0 \implies A \perp \!\!\!\perp D \mid \mathbf{B}, Y, C_1 = 1$
- Note that this does not assume anything about the separating set test :



Inferring separating sets without enumerating all possible causal graphs

Query $Y \perp C_1 | X_1 ?$

Assumptions

All testable conditional independences from data $X_1 \perp \!\!\!\perp X_3 \mid X_4$ $Y \perp \!\!\!\perp C_2 \mid X_1, C_1 = 0$ $X_2 \perp \!\!\!\perp C_2 \mid Y, C_1 = 0$

Logic encoding of d-separation [Hyttinen et al. 2014]



A simple causal feature selection algorithm

Source domains data

| C1 | C2 | X1 | X2 | Y | |
|----|----|-----------|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

List of combinations of features ordered by source domain loss in predicting Y

✤ L=({X1, C2}, {X1, X2, C2}, {X1, X2}, …)

A simple causal feature selection algorithm

Sel

Source domains data

| C1 | C2 | X1 | X2 | Υ | |
|----|----|-----------|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

List of combinations of features ordered by source domain loss in predicting Y

A simple causal feature selection algorithm

Source domains data

| C1 | C2 | X1 | X2 | Y | |
|----|-----------|-----|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

All data (including target)

| C1 | C2 | X1 | X2 | Υ | |
|----|-----------|-----|-----------|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |
| 1 | 0 | 0.2 | 0 | ? | |
| 1 | 0 | 0.3 | 0 | ? | |
| 1 | 0 | 0.3 | 1 | ? | |
| | | | | | |

Query $Y \perp C_1 | S?$

Assumptions

All testable conditional independences from data

 $X_{1} \perp X_{3} \mid X_{4}$ $Y \perp C_{2} \mid X_{1}, C_{1} = 0$ $X_{2} \perp C_{2} \mid Y, C_{1} = 0$

Logic encoding of d-separation [Hyttinen et al. 2014] List of combinations of features ordered by source domain loss in predicting Y

S= {X1, C2}



Select new set S
Source domains data

| C1 | C2 | X1 | X2 | Y | |
|----|-----------|-----------|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

All data (including target)

| C1 | C2 | X1 | X2 | Υ | |
|----|-----------|-----|-----------|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |
| 1 | 0 | 0.2 | 0 | ? | |
| 1 | 0 | 0.3 | 0 | ? | |
| 1 | 0 | 0.3 | 1 | ? | |
| | | | | | |

Query $Y \perp C_1 | S?$

Assumptions

All testable conditional independences from data

 $X_{1} \perp X_{3} \mid X_{4}$ $Y \perp C_{2} \mid X_{1}, C_{1} = 0$ $X_{2} \perp C_{2} \mid Y, C_{1} = 0$

Logic encoding of d-separation [Hyttinen et al. 2014] List of combinations of features ordered by source domain loss in predicting Y





Source domains data

| C1 | C2 | X1 | X2 | Y | |
|----|-----------|-----------|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

All data (including target)

| C1 | C2 | X1 | X2 | Υ | |
|----|-----------|-----|-----------|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |
| 1 | 0 | 0.2 | 0 | ? | |
| 1 | 0 | 0.3 | 0 | ? | |
| 1 | 0 | 0.3 | 1 | ? | |
| | | | | | |

Query $Y \perp C_1 | S?$

Assumptions

All testable conditional independences from data

 $X_{1} \perp X_{3} \mid X_{4}$ $Y \perp C_{2} \mid X_{1}, C_{1} = 0$ $X_{2} \perp C_{2} \mid Y, C_{1} = 0$

Logic encoding of d-separation [Hyttinen et al. 2014] List of combinations of features ordered by source domain loss in predicting Y



Source domains data

| C1 | C2 | X1 | X2 | Y | |
|----|-----------|-----------|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

All data (including target)

| C1 | C2 | X1 | X2 | Υ | |
|----|-----------|-----|-----------|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |
| 1 | 0 | 0.2 | 0 | ? | |
| 1 | 0 | 0.3 | 0 | ? | |
| 1 | 0 | 0.3 | 1 | ? | |
| | | | | | |

Query $Y \perp C_1 | S?$

Assumptions

All testable conditional independences from data

 $X_{1} \perp X_{3} \mid X_{4}$ $Y \perp C_{2} \mid X_{1}, C_{1} = 0$ $X_{2} \perp C_{2} \mid Y, C_{1} = 0$

Logic encoding of d-separation [Hyttinen et al. 2014] List of combinations of features ordered by source domain loss in predicting Y



Select new set S

S= {X1, X2, C2}



Source domains data

| C1 | C2 | X1 | X2 | Y | |
|----|-----------|-----------|----|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |

Standard feature selection

All data (including target)

| C1 | C2 | X1 | X2 | Υ | |
|----|-----------|-----|-----------|---|--|
| 0 | 0 | 0.1 | 1 | 0 | |
| 0 | 0 | 0.2 | 1 | 0 | |
| 0 | 0 | 1.1 | 2 | 1 | |
| 0 | 1 | 3.1 | 2 | 1 | |
| 0 | 1 | 3.2 | 3 | 1 | |
| 0 | 1 | 4 | 3 | 1 | |
| 1 | 0 | 0.2 | 0 | ? | |
| 1 | 0 | 0.3 | 0 | ? | |
| 1 | 0 | 0.3 | 1 | ? | |
| | | | | | |

Query $Y \perp C_1 | S?$

Assumptions

All testable conditional independences from data

 $X_{1} \perp X_{3} \mid X_{4}$ $Y \perp C_{2} \mid X_{1}, C_{1} = 0$ $X_{2} \perp C_{2} \mid Y, C_{1} = 0$

Logic encoding of d-separation [Hyttinen et al. 2014] List of combinations of features ordered by source domain loss in predicting Y



Source domains data

| C2 | X1 | X2 | Υ | |
|-----------|-----------|-----------|---|-----------------|
| | 0.1 | 1 | 0 | |
| | 0.2 | 1 | 0 | a second second |
| | 1.1 | 2 | 1 | |
| | 3.1 | 2 | 1 | |
| | 3.2 | 3 | 1 | |
| | 4 | 3 | 1 | |

Standard feature selection

All data (including target)

| C2 | X1 | X2 | Y |
|-----------|-----|----|-----|
| | 0.1 | 1 | 0 |
| | 0.2 | 1 | 0 |
| | 1.1 | 2 | 1 |
| | 3.1 | 2 | 1 # |
| | 3.2 | 3 | 1 |
| | 4 | 3 | 1 |
| | 0.2 | 0 | ? |
| | 0.3 | 0 | ? |
| | 0.3 | 1 | ? |

Query $Y \perp C_1 | S?$

Assumptions

All testable conditional ndependences from data $X_1 \perp \!\!\!\perp X_3 \mid X_4$ $Y \perp \!\!\!\perp C_2 \mid X_1, C_1 = 0$

ogic encoding of d-separation. [Hyttinen et al. 2014] List of combinations of features ordered by source domain loss in predicting Y



Desiderata for a causality inspired domain adaptation method

- X, Y and changes can be represented by an unknown causal graph
 - Allow for latent confounders

 - Instead of modeling changes between each domain, distinguish the change between the mixture of sources and the target
 - domains, then Y|T(X) is invariant also in the target domain
 - Only search for invariant features with respect to current target task

Avoid parametric assumptions, allow for heterogeneous effects across domains

Avoid common assumption that if Y T(X) is invariant across multiple source



Desiderata for a causality inspired domain adaptation method

- X, Y and changes can be represented by an unknown causal graph
 - Allow for latent confounders
 - Avoid parametric assumptions, allow for heterogeneous effects across domains
 - Instead of modeling changes between each domain, distinguish the change between the mixture of sources and the target
 - Avoid common assumption that if Y T(X) is invariant across multiple source domains, then Y|T(X) is invariant also in the target domain
 - Only search for invariant features with respect to current target task

No need to find causal graph or equivalence class



Limitations and future work

- separating
 - most informative interventions?
- input scales to tens of vars (including context variables)
 - Extension: use approximate algorithms, combine with low dim representations
- Can we apply this to multi-task RL (e.g. in factored MDPs)?

Potentially too conservative: Separating sets may exist that are not provably

• Extension: can we use active learning/intervention design to decide

• Scalability: using (error-correcting) logic-based encoding with all CI tests as

A sneak peak in applications of causality-inspired ML:

An Approach to Data-Driven Domain Adaptation



- Only relevant features needed to predict *Y*
- Augmented graph learned by CD-NOD
 - Independently changing modules θ_i
 - Special case: invariant modules
- Domain adaption: inference on this graphical model
 - Infer the posterior of *Y* in target domain ٠
 - Nonparametric methods to model conditional distributions •

Zhang, Gong, Stojanov, Huang, Liu, and Glymour, "Domain Adaptation As a Problem of Inference on Graphical Models," NeurIPS 2020. (Huang et al., ICML'19 for time series data)

<u>https://www.youtube.com/watch?v=_MVi6XzOdD0&ab_channel=OnlineCausalInferenceSeminar</u>



https://arxiv.org/abs/1903.01672



AdaRL: What, Where, and How to Adapt in Transfer RL

Biwei Huang, Fan Feng, Chaochao Lu, Sara Magliacane, Kun Zhang

- We have n source domains with random trajectories
- Learn a factored MDP (symbolic inputs) or POMDP (images) with latent change
 - Identify the minimal dimensions of the state and change factors that are **necessary and sufficient** for policy optimisation

factors that are constant in each domain, but vary across domains over sources

AdaRL: What, Where, and How to Adapt in Transfer RL

Biwei Huang, Fan Feng, Chaochao Lu, Sara Magliacane, Kun Zhang

- We have n source domains with random trajectories
- Learn a factored MDP (symbolic inputs) or POMDP (images) with latent change
 - Identify the minimal dimensions of the state and change factors that are **necessary and sufficient** for policy optimisation
- In the target domain learn the value of the change factor and apply this policy

https://arxiv.org/abs/2107.02729

factors that are constant in each domain, but vary across domains over sources

Learn a policy over all source domains, parametrised in the minimal change factors

AdaRL: What, Where, and How to Adapt in Transfer RL

Biwei Huang, Fan Feng, Chaochao Lu, Sara Magliacane, Kun Zhang



https://arxiv.org/abs/2107.02729

Target domain **Identify compact** domaingeneralisable representations $(\mathbf{s}_t^{min}, \theta_k^{min})$ Estimate Policy learning on domain-specific source domains parameters θ_{target}^{min} $\pi^*(\theta_k^{min})$ Optimal parametrised policy Simplifying assumption: no $\pi^*(\theta_{target}^{min})$ new edges in Optimal target policy target domain

ICLR 2022



FansRL: Factored Adaptation for Non-Stationary Reinforcement Learning **NeurIPS 2022** Fan Feng, Biwei Huang, Kun Zhang, Sara Magliacane

 The latent change factors are not constant anymore and they model nonstationarity





Change factors follow a Markov process: Non-stationary environments (wind changes)

- **Discrete/abrupt** changes
- Continuous/smooth changes



Non-stationary rewards (target changes)

https://arxiv.org/abs/2203.16582



Conclusions

- distribution shift

 - This is true even with:
 - Unknown causal graph
 - Missing data/CI (so unknown MEC)
- separations even with missing data, even without reconstructing MEC

D-separation [Pearl 1988] is a principled way to reason about invariances and

• Not a new observation, known since [Schoelkopf et al 2012, Zhang et al. 2013]

D-separation logic encodings [Hyttinen et al 2014] allow us to reason about d-

Conclusions

- distribution shift

 - This is true even with:
 - Unknown causal graph
 - Missing data/CI (so unknown MEC)
- separations even with missing data, even without reconstructing MEC

D-separation [Pearl 1988] is a principled way to reason about invariances and

• Not a new observation, known since [Schoelkopf et al 2012, Zhang et al. 2013]

D-separation logic encodings [Hyttinen et al 2014] allow us to reason about d-

Thanks! Questions?

(joint work with Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, Joris Mooij, Biwei Huang, Fan Feng, Chaochao Lu and Kun Zhang)

- (Joint Causal Inference)
- We assume Y cannot be intervened upon directly

• We assume that there exists an acyclic causal graph that fits all the data

- (Joint Causal Inference)
- We assume Y cannot be intervened upon directly
- We assume no extra dependences involving Y in target domain C1=1 $A, D, \mathbf{B} \subset \mathbf{V} \setminus \{Y, C_1\}$

• We assume that there exists an acyclic causal graph that fits all the data



- We assume that there exists an acyclic causal graph that fits all the data (Joint Causal Inference)
- We assume Y cannot be intervened upon directly
- We assume no extra dependences involving Y in target domain C1=1 $A, D, \mathbf{B} \in \mathbf{V} \setminus \{Y, C_1\}$ $Y \perp A \mid \mathbf{B}, C_1 = 0 \implies Y \perp A \mid \mathbf{B}, C_1 = 1$ $A \perp D \mid \mathbf{B}, Y, C_1 = 0 \implies A \perp D \mid \mathbf{B}, Y, C_1 = 1$

There can be extra independences in the target



- We assume that there exists an acyclic causal graph that fits all the data (Joint Causal Inference)
- We assume Y cannot be intervened upon directly
- We assume no extra dependences involving Y in target domain C1=1 $Y \perp A \mid \mathbf{B}, C_1 = 0 \implies Y \perp A \mid \mathbf{B}, C_1 = 1$ $A, D, \mathbf{B} \subset \mathbf{V} \setminus \{Y, C_1\}$ $A \perp D \mid \mathbf{B}, Y, C_1 = 0 \implies A \perp D \mid \mathbf{B}, Y, C_1 = 1$

Some CIs are still missing: $A \perp D \mid \mathbf{B}, Y$ $Y \perp A \mid \mathbf{B}$



- We assume that there exists an **acyclic** causal graph that fits all the data (Joint Causal Inference)
- We assume Y cannot be intervened upon directly
- We assume no extra dependences involving Y in target domain C1=1 $A, D, \mathbf{B} \in \mathbf{V} \setminus \{Y, C_1\}$ $Y \perp A \mid \mathbf{B}, C_1 = 0 \implies Y \perp A \mid \mathbf{B}, C_1 = 1$ $A \perp D \mid \mathbf{B}(Y, C_1 = 0) \implies A \perp D \mid \mathbf{B}, Y, C_1 = 1$
- Note that this does not assume anything about the separating set test :



d-separation: complete example



Nodes i and j is d-separated by A if all paths between them are blocked

d-separation: complete example



Nodes i and j is d-separated by A if all paths between them are blocked

non-collider collider

- If $1 \in \mathbf{A}$, the path is **blocked**, **OR**
- If $3 \notin A$ and $2 \notin A$, the path is blocked

d-separation: complete example



Nodes i and j is d-separated by A if all paths between them are blocked

A simple example

| C1 | C2 | X1 | X2 | Υ |
|-----------|-----------|-----------|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 0 | 1 | 3.1 | 2 | 1 |
| 0 | 1 | 3.2 | 3 | 1 |
| 0 | 1 | 4 | 3 | 1 |
| 1 | 0 | 0.2 | 0 | ? |
| 1 | 0 | 0.3 | 0 | ? |
| 1 | 0 | 0.3 | 1 | ? |

A simple example

| C1 | C2 | X1 | X2 | Y |
|-----------|-----------|-----------|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 0 | 1 | 3.1 | 2 | 1 |
| 0 | 1 | 3.2 | 3 | 1 |
| 0 | 1 | 4 | 3 | 1 |
| 1 | 0 | 0.2 | 0 | ? |
| 1 | 0 | 0.3 | 0 | ? |
| 1 | 0 | 0.3 | 1 | ? |

 $X_2 \perp C_2 \mid Y, C_1 = 0$

- $Y \perp C_2 \mid C_1 = 0$ $Y \perp C_2 | X_1, C_1 = 0$
- Perform allowed CI tests

| C1 | C2 | X1 | X2 | Y |
|-----------|-----------|-----------|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 0 | 1 | 3.1 | 2 | 1 |
| 0 | 1 | 3.2 | 3 | 1 |
| 0 | 1 | 4 | 3 | 1 |
| 1 | 0 | 0.2 | 0 | ? |
| 1 | 0 | 0.3 | 0 | ? |
| 1 | 0 | 0.3 | 1 | ? |

C2

X1



A simple example

| C1 | C2 | X1 | X2 | Y |
|-----------|-----------|-----------|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 0 | 1 | 3.1 | 2 | 1 |
| 0 | 1 | 3.2 | 3 | 1 |
| 0 | 1 | 4 | 3 | 1 |
| 1 | 0 | 0.2 | 0 | ? |
| 1 | 0 | 0.3 | 0 | ? |
| 1 | 0 | 0.3 | 1 | ? |

- $Y \perp C_2 \mid C_1 = 0$ $Y \perp C_2 | X_1, C_1 = 0$ $X_2 \perp C_2 \mid Y, C_1 = 0$
- Perform allowed CI tests



All possible compatible graphs



A simple example

| C1 | C2 | X1 | X2 | Υ |
|-----------|-----------|-----------|-----------|---|
| 0 | 0 | 0.1 | 1 | 0 |
| 0 | 0 | 0.2 | 1 | 0 |
| 0 | 0 | 1.1 | 2 | 1 |
| 0 | 1 | 3.1 | 2 | 1 |
| 0 | 1 | 3.2 | 3 | 1 |
| 0 | 1 | 4 | 3 | 1 |
| 1 | 0 | 0.2 | 0 | ? |
| 1 | 0 | 0.3 | 0 | ? |
| 1 | 0 | 0.3 | 1 | ? |

 $Y \perp C_1 \mid X_1$

- $Y \perp C_2 \mid C_1 = 0$ $Y \perp C_2 | X_1, C_1 = 0$
- $X_2 \perp C_2 \mid Y, C_1 = 0$
- Perform allowed CI tests



All possible compatible graphs

• We can prove untestable separating test without reconstructing the graph:

True in all possible compatible graphs

