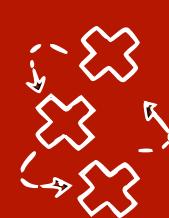


# Causal Data Science

## Lecture 12.2: Causality-inspired ML

Lecturer: Sara Magliacane

UvA - Spring 2024



# Causality + machine learning (non-exhaustive list)

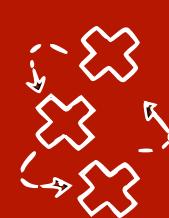
## 1. Machine learning (ML) helps causality

- Causal discovery - learning causal graphs from data
- Causal effect estimation - matching, weighting, double ML
- (Causal) representation learning

## 2. Causality (in the most general definition) helps machine learning

- Robustness, Transfer learning
- Reinforcement Learning
- Bias mitigation, fairness

<https://arxiv.org/pdf/1705.08821.pdf>, <https://arxiv.org/pdf/1802.05664.pdf>, <https://arxiv.org/pdf/1605.03661.pdf>, <https://crl.causalai.net/>, [https://www.youtube.com/watch?v=Obuu3w809CI&ab\\_channel=ConnorJerzak](https://www.youtube.com/watch?v=Obuu3w809CI&ab_channel=ConnorJerzak) and many many others



# Causality + machine learning (non-exhaustive list)

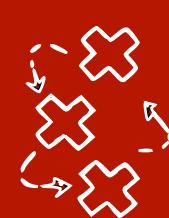
## 1. Machine learning (ML) helps causality

- Causal discovery - learning causal graphs from data
- Causal effect estimation - matching, weighting, double ML
- (Causal) representation learning

## 2. Causality (in the most general definition) helps machine learning

- Robustness, Transfer learning
- Reinforcement Learning
- Bias mitigation, fairness

<https://arxiv.org/pdf/1705.08821.pdf>, <https://arxiv.org/pdf/1802.05664.pdf>, <https://arxiv.org/pdf/1605.03661.pdf>, <https://crl.causalai.net/>, [https://www.youtube.com/watch?v=Obuu3w809CI&ab\\_channel=ConnorJerzak](https://www.youtube.com/watch?v=Obuu3w809CI&ab_channel=ConnorJerzak) and many many others



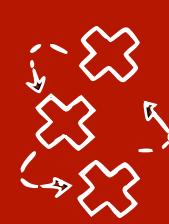
# Causal Hierarchy [Pearl 2009, 2018]



Most ML

Causality

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing $X$ change my belief in $Y$ ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do $X$ ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it $X$ that caused $Y$ ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smoking the past 2 years?



# Causal Hierarchy [Pearl 2009, 2018]

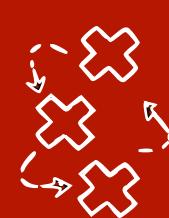


Most ML

Causality

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing $X$ change my belief in $Y$ ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do $X$ ?	What if I take aspirin, will my headache be cured? What if I smoke cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	What if $x$ had been different? What would have happened if $x$ had been different?	E.g. need many experiments or strong assumptions to identify the causal graph or the causal variables

“Full” causality can be **not necessary** or **too expensive** ->



# Causal Hierarchy [Pearl 2009, 2018]

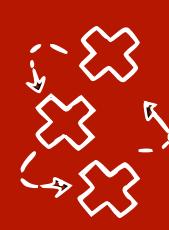


Most ML

Causality

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing $X$ change my belief in $Y$ ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do $X$ ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it $X$ that caused $Y$ ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smoking the past 2 years?

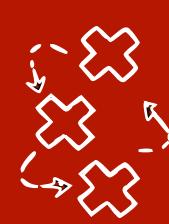
“Full” causality can be **not necessary** or **too expensive** -> *Causality-Inspired*



# Causality vs Transfer learning

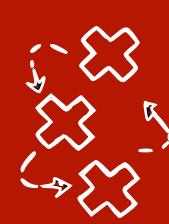
- Transfer learning:
  - How can I predict what happens when the distribution changes?





# Causality vs Transfer learning

- 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  $\text{do}(X)$ :
      - do-calculus [Pearl, 2009]
    - **Soft intervention on X**  $\approx$  change of distribution of  $P(X|\text{parents})$
- 
- 
- 
- 
- 



# Causality vs Transfer learning

- Transfer learning:

- How can I predict what happens when the distribution changes



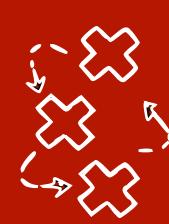
Very general - can model also changes in distribution that are not from “real” interventions



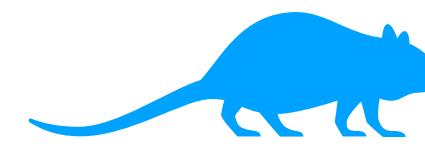
- What does a counterfactual intervention  $\text{do}(X)$ :

- do-calculus [Pearl, 2009]

- Soft intervention on  $X \approx$  change of distribution of  $P(X| \text{parents})$**

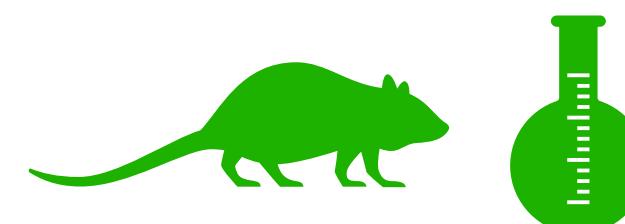


# Unsupervised domain adaptation - toy example



	X1	X2	Y
Wildtype	0,1	2	0
Wildtype	0,2	3	0
Wildtype	1,1	2	1
Wildtype	0,1	3	0

Source domain

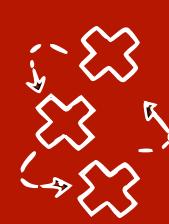


	X1	X2	Y
Gene A	3,1	2	?
Gene A	3,2	3	?
Gene A	4	2	?
Gene A	3,2	3	?

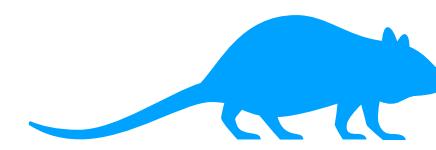
Target domain

No labels in target

- **Task:** Learn a model  $Y = \hat{f}(X_1, X_2)$  on the source domain so that it can reliably estimate  $Y$  in target domain

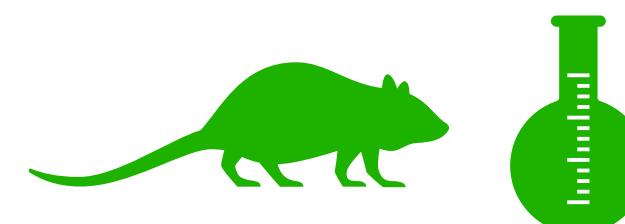


# Domain adaptation from the graphical perspective



C	X1	X2	Y
0	0,1	2	0
0	0,2	3	0
0	1,1	2	1
0	0,1	3	0

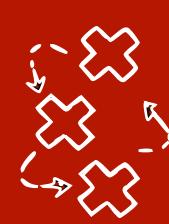
Source domain



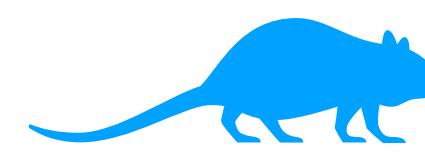
C	X1	X2	Y
1	3,1	2	?
1	3,2	3	?
1	4	2	?
1	3,2	3	?

Target domain

1. We add a context variable C to distinguish the domains



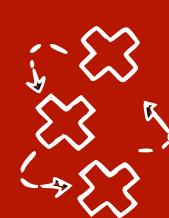
# Domain adaptation from the graphical perspective



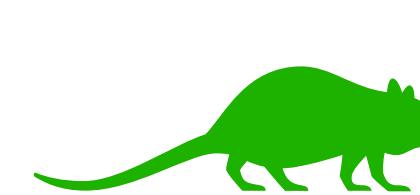
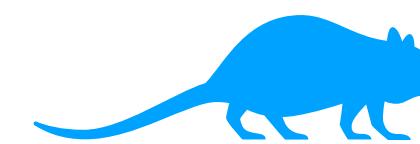
C	X1	X2	Y
0	0,1	2	0
	0,2	3	0
	1,1	2	1
	0,1	3	0
1	3,1	2	?
	3,2	3	?
	4	2	?
	3,2	3	?

Source domain      Target domain

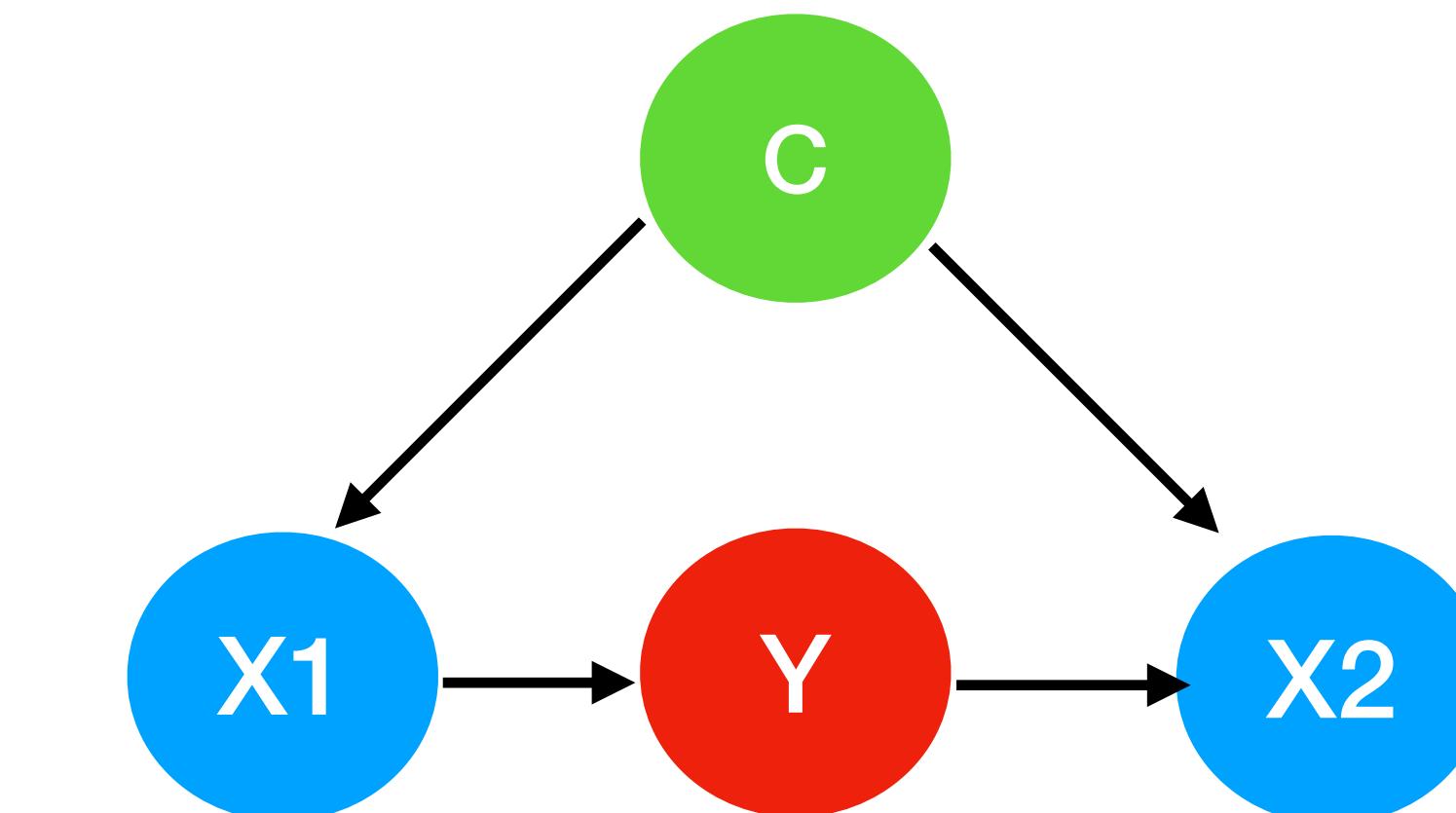
1. We add a context variable C to distinguish the domains
2. We consider the data as coming from a single distribution  $P(X_1, X_2, Y, C)$



# Domain adaptation from the graphical perspective

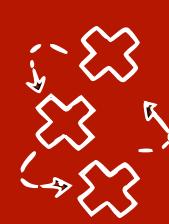


C	X1	X2	Y
0	0,1	2	0
	0,2	3	0
	1,1	2	1
	0,1	3	0
1	3,1	2	?
	3,2	3	?
	4	2	?
	3,2	3	?



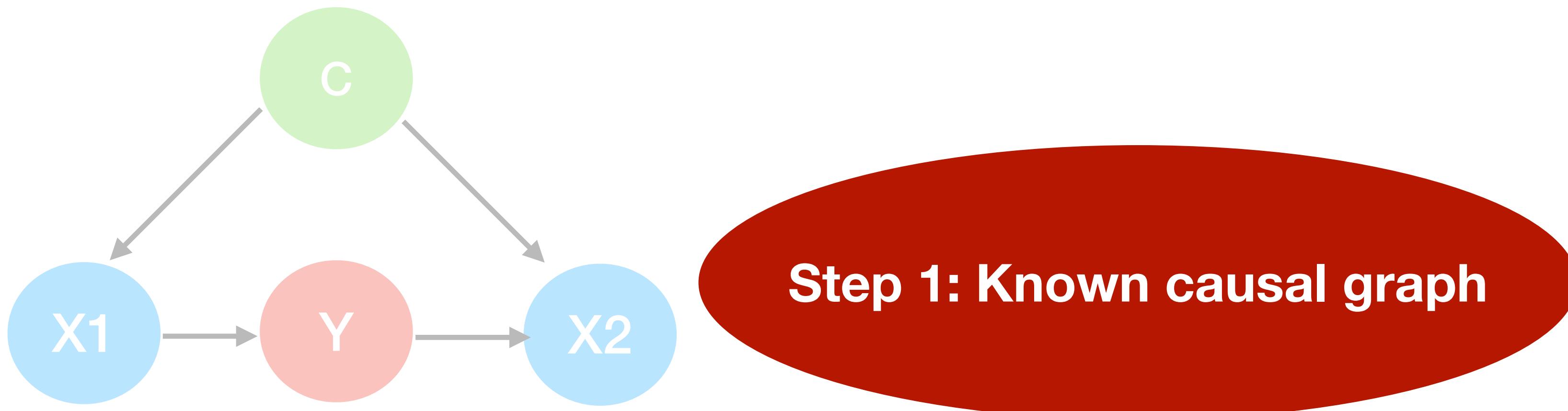
We can represent  $P(X_1, X_2, Y, C)$  with an **(unknown)** causal graph

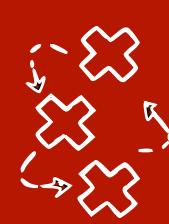
1. We add a context variable C to distinguish the domains
2. We consider the data as coming from a single distribution  $P(X_1, X_2, Y, C)$



# Separating features = safe for (adversarial) domain adaptation

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

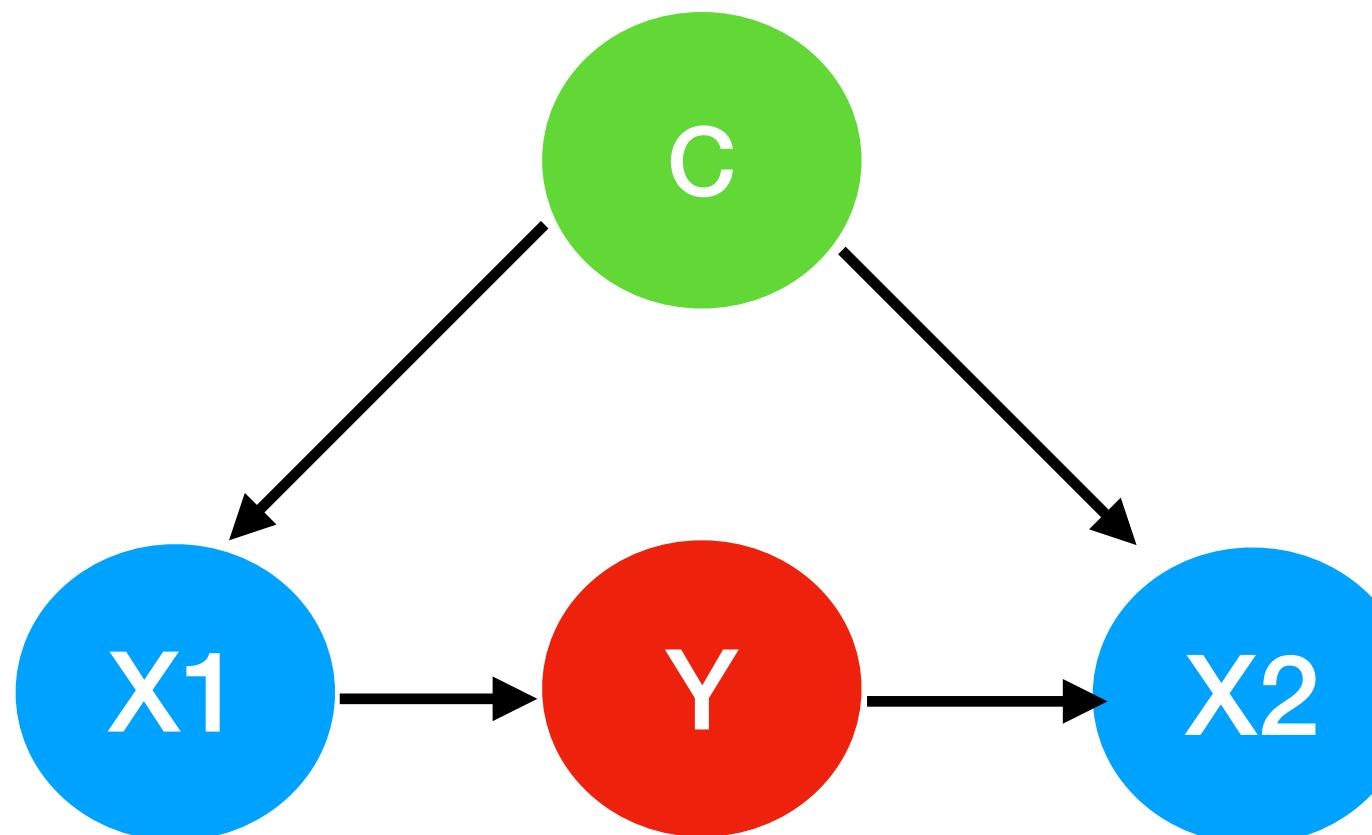


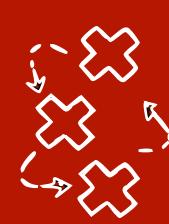


# Separating features = safe for (adversarial) domain adaptation

Aka stable features, invariant features etc.

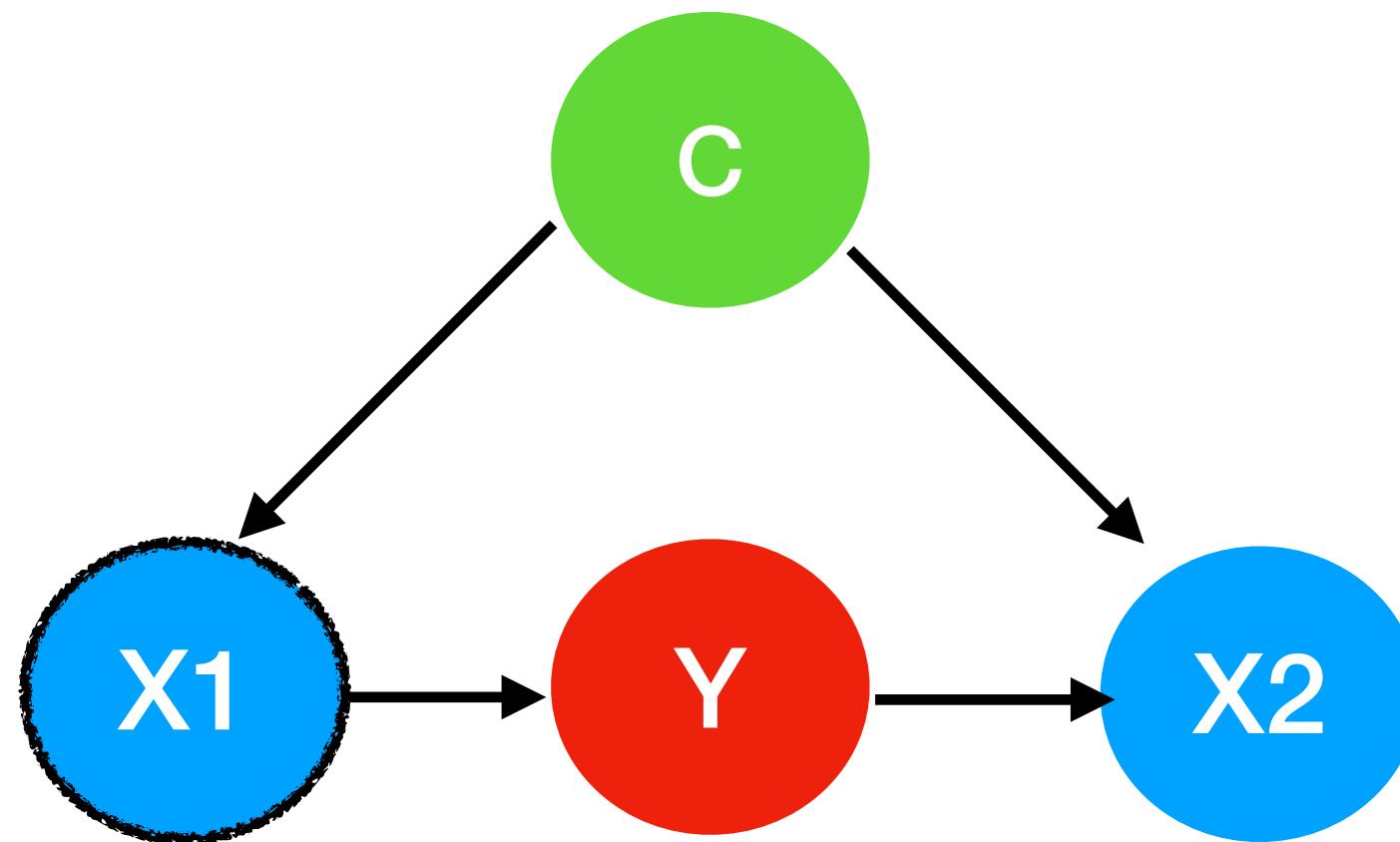
- **Separating features:** sets of features that d-separate Y from the context





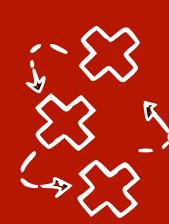
# Separating features = safe for (adversarial) domain adaptation

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



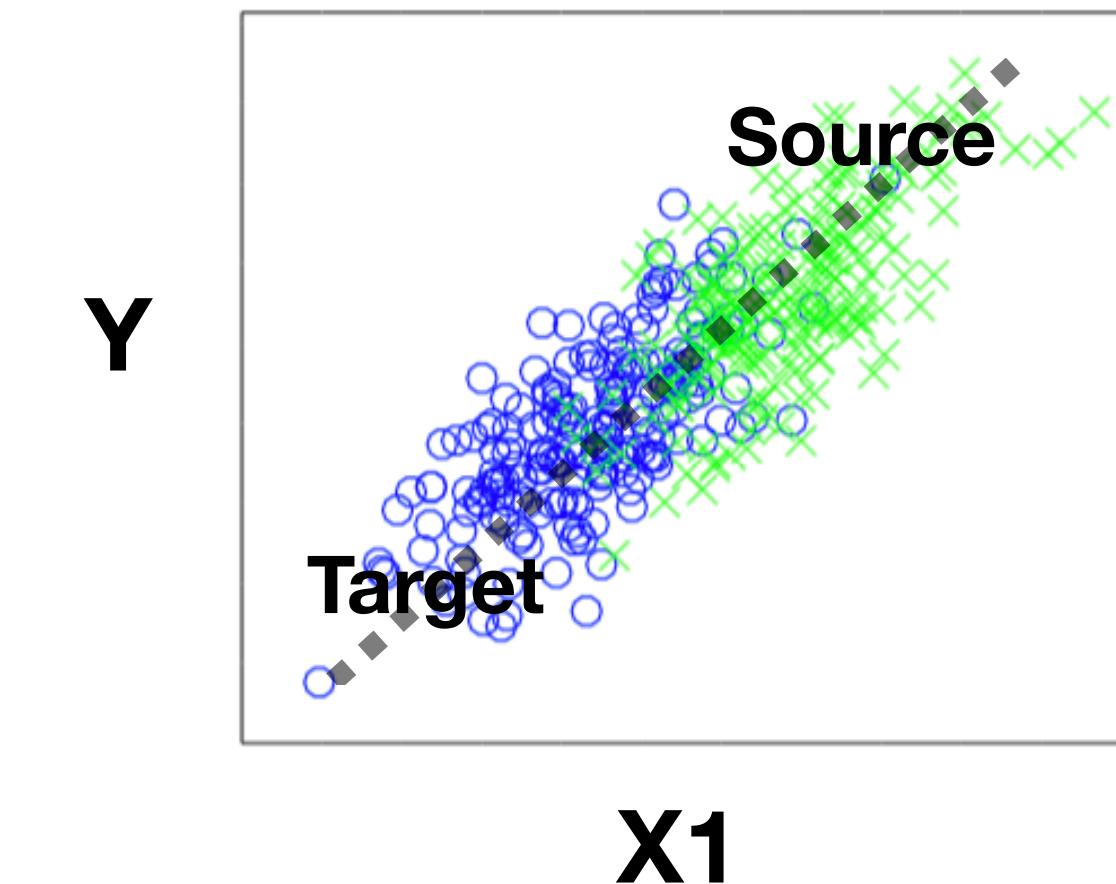
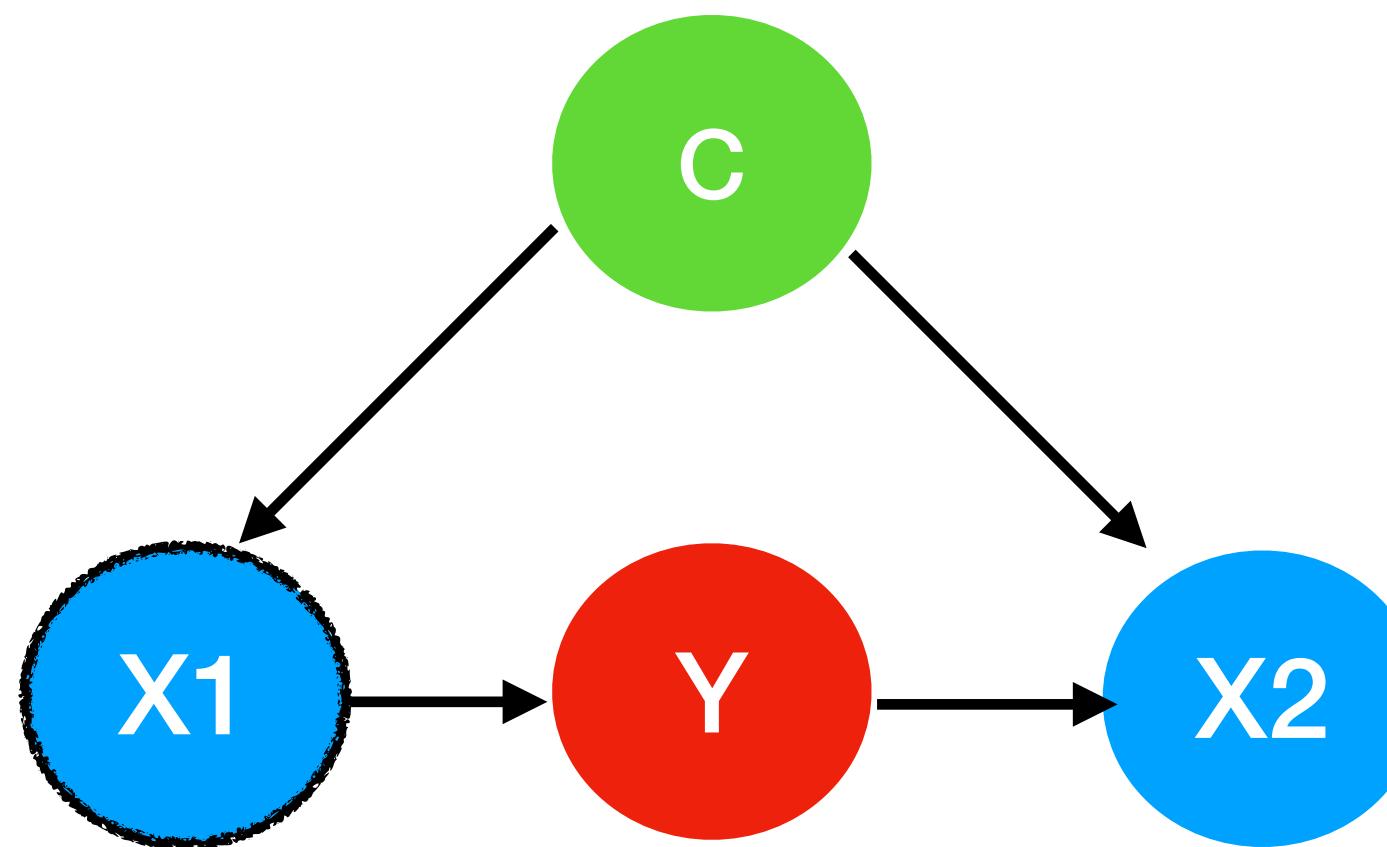
$$Y \perp\!\!\! \perp_d C | X_1 \iff Y \perp\!\!\! \perp C | X_1$$

(under Markov and faithfulness assumptions)



# Separating features = safe for (adversarial) domain adaptation

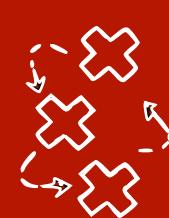
- **Separating features:** sets of features that d-separate Y from the context



$$Y \perp\!\!\! \perp_d C | X_1 \iff Y \perp\!\!\! \perp C | X_1$$

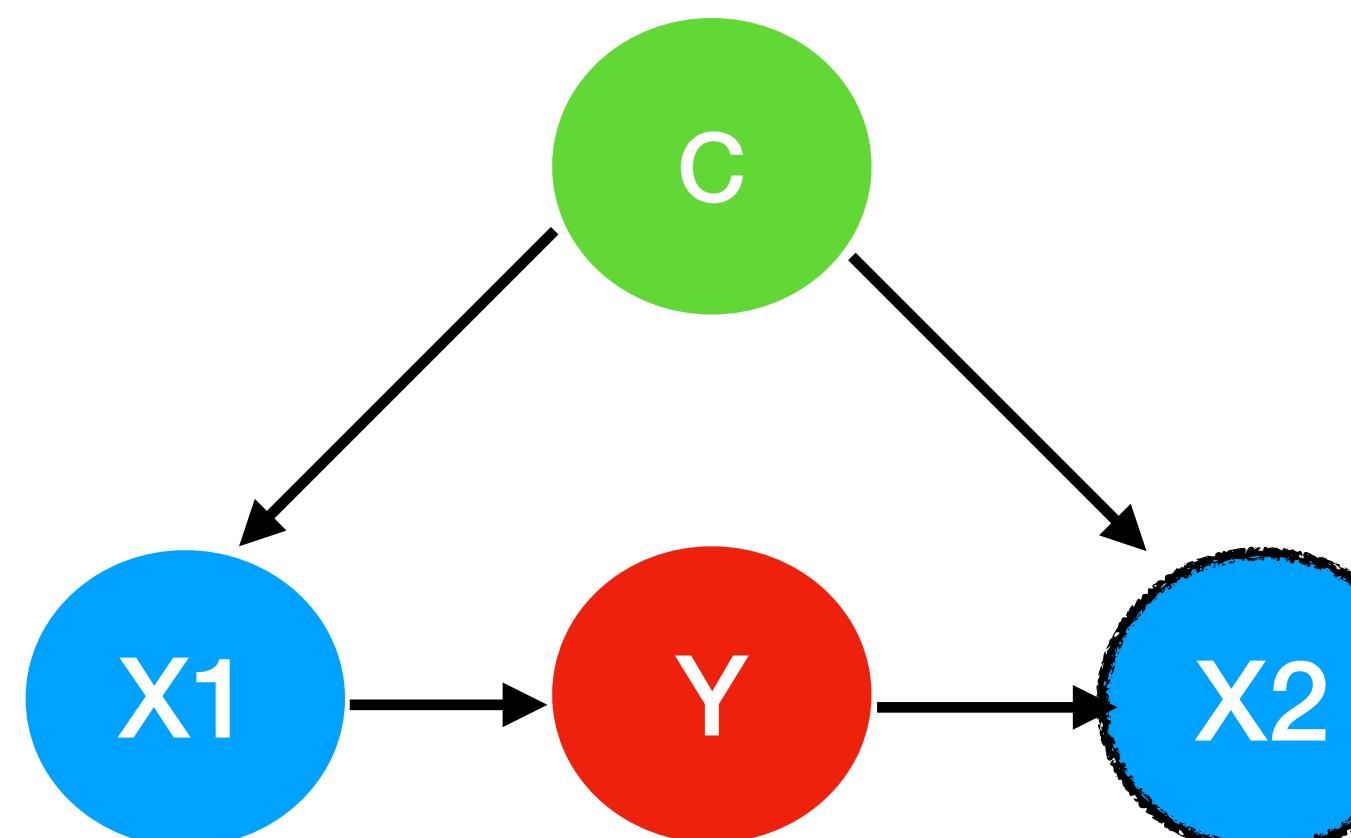
(under Markov and faithfulness assumptions)

$$\begin{aligned} Y \perp\!\!\! \perp C | X_1 &\equiv \\ P(Y|X_1, C=0) &= P(Y|X_1, C=1) \end{aligned}$$

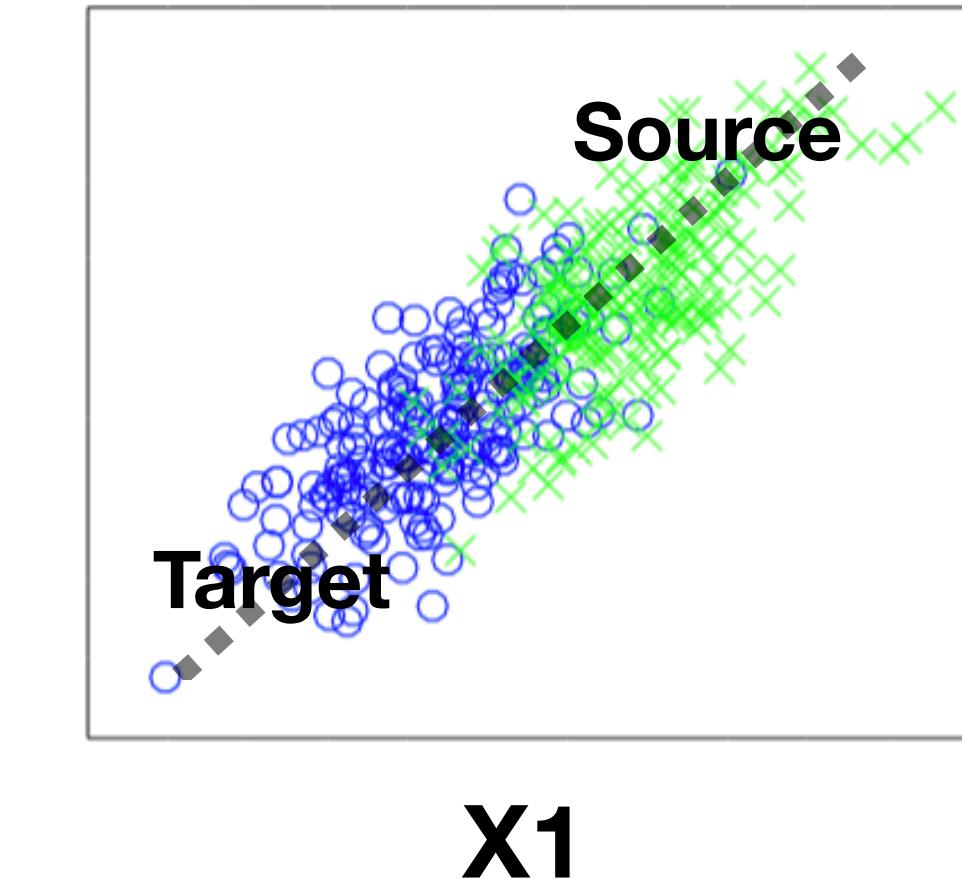


# Separating features = safe for (adversarial) domain adaptation

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



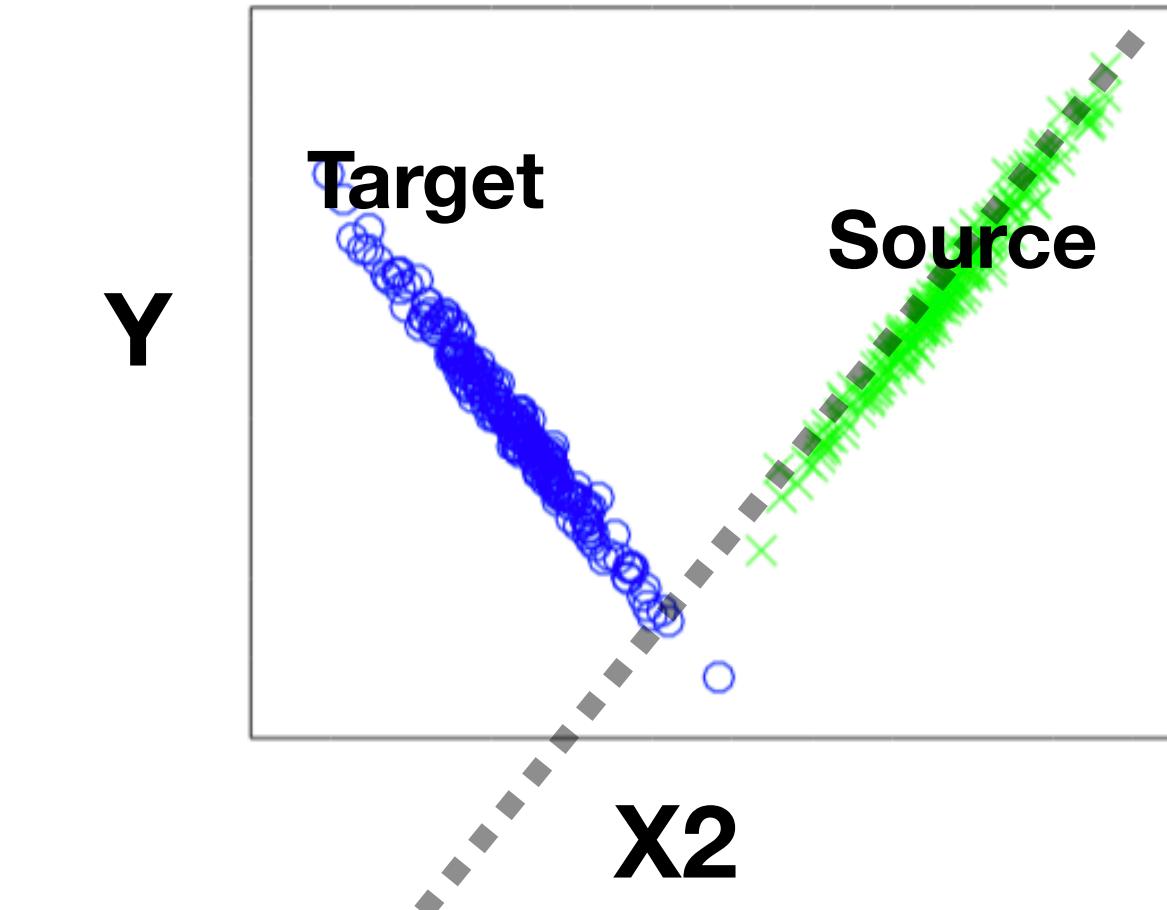
Y



$$Y \perp\!\!\!\perp C | X_2 \iff Y \perp\!\!\!\perp C | X_1$$

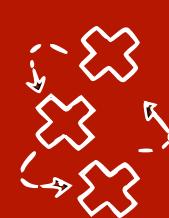
(under Markov and faithfulness assumptions)

$$Y \perp\!\!\!\perp C | X_1 \equiv P(Y|X_1, C=0) = P(Y|X_1, C=1)$$



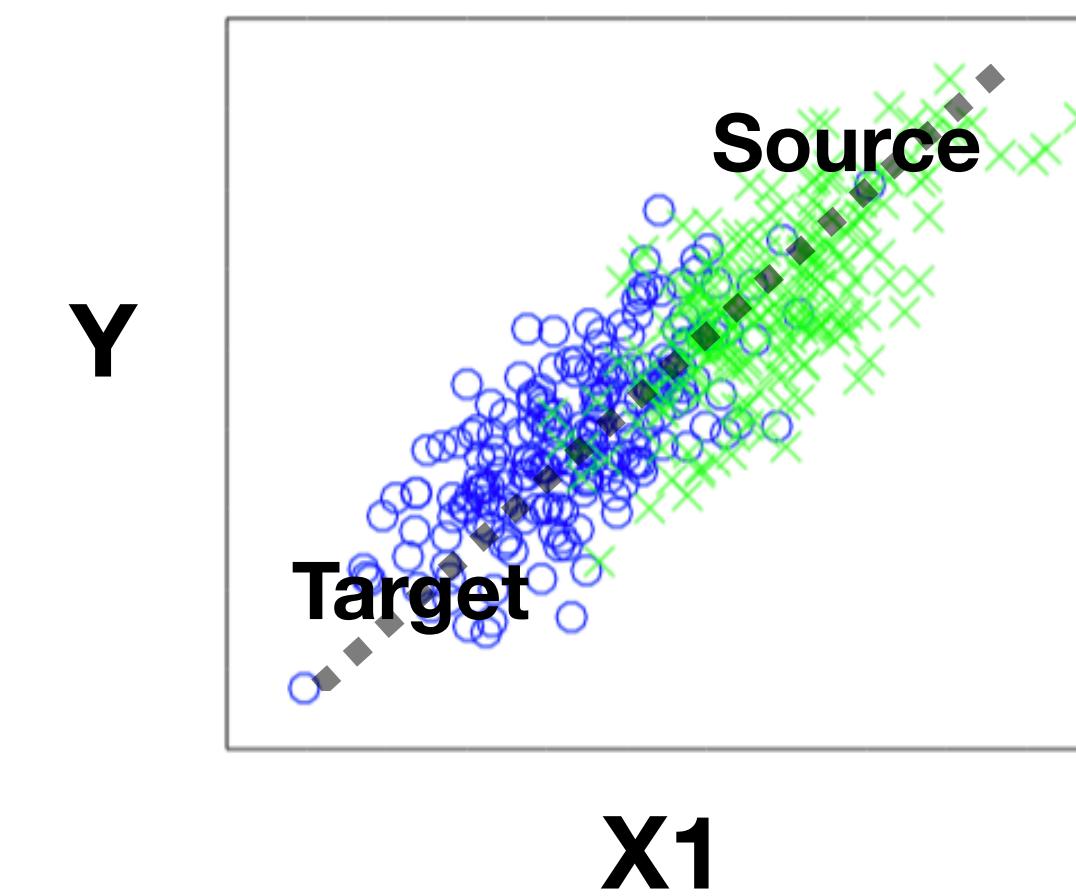
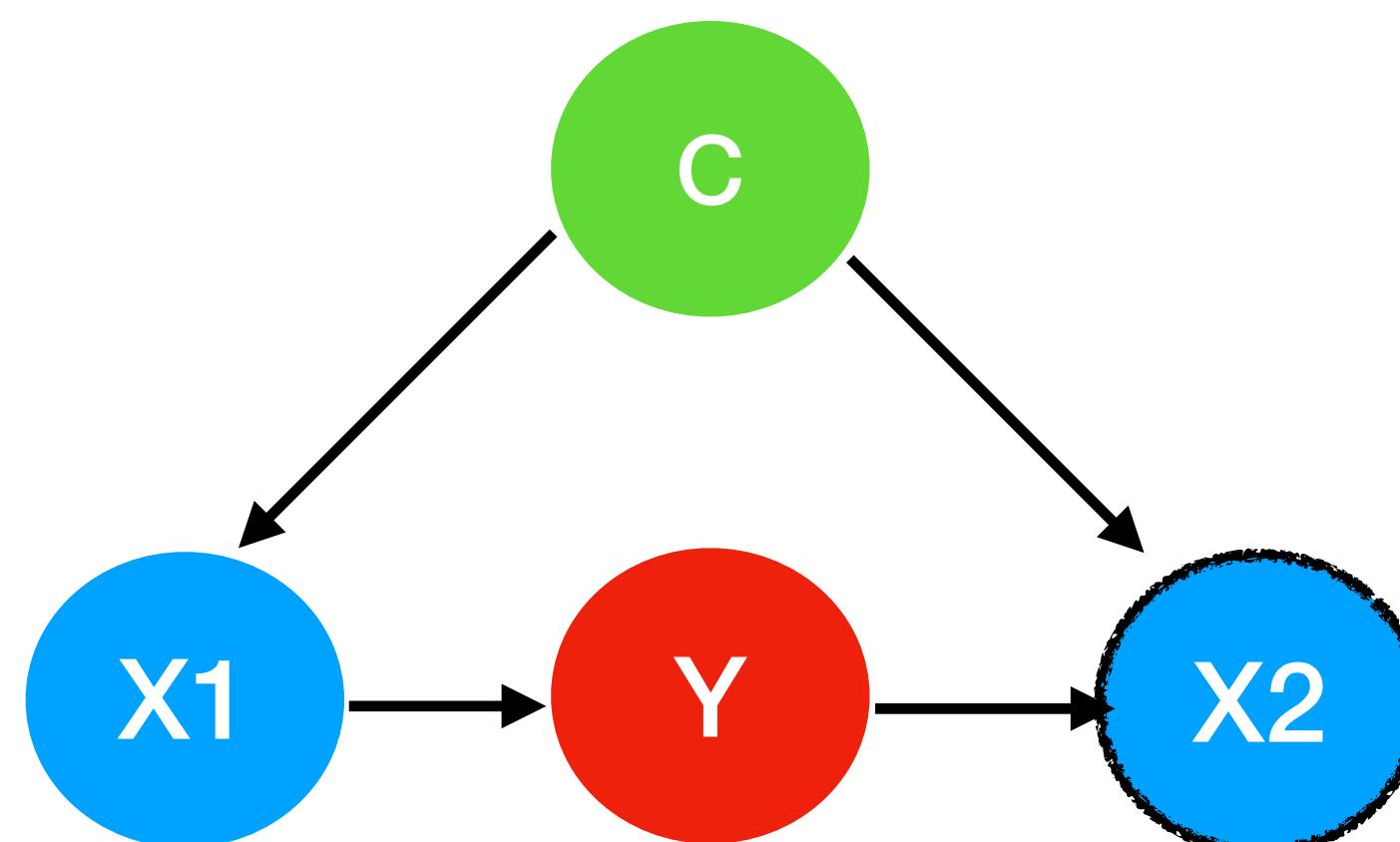
$$Y \perp\!\!\!\perp C | X_2 \equiv$$

$$P(Y|X_2, C=0) \neq P(Y|X_2, C=1)$$



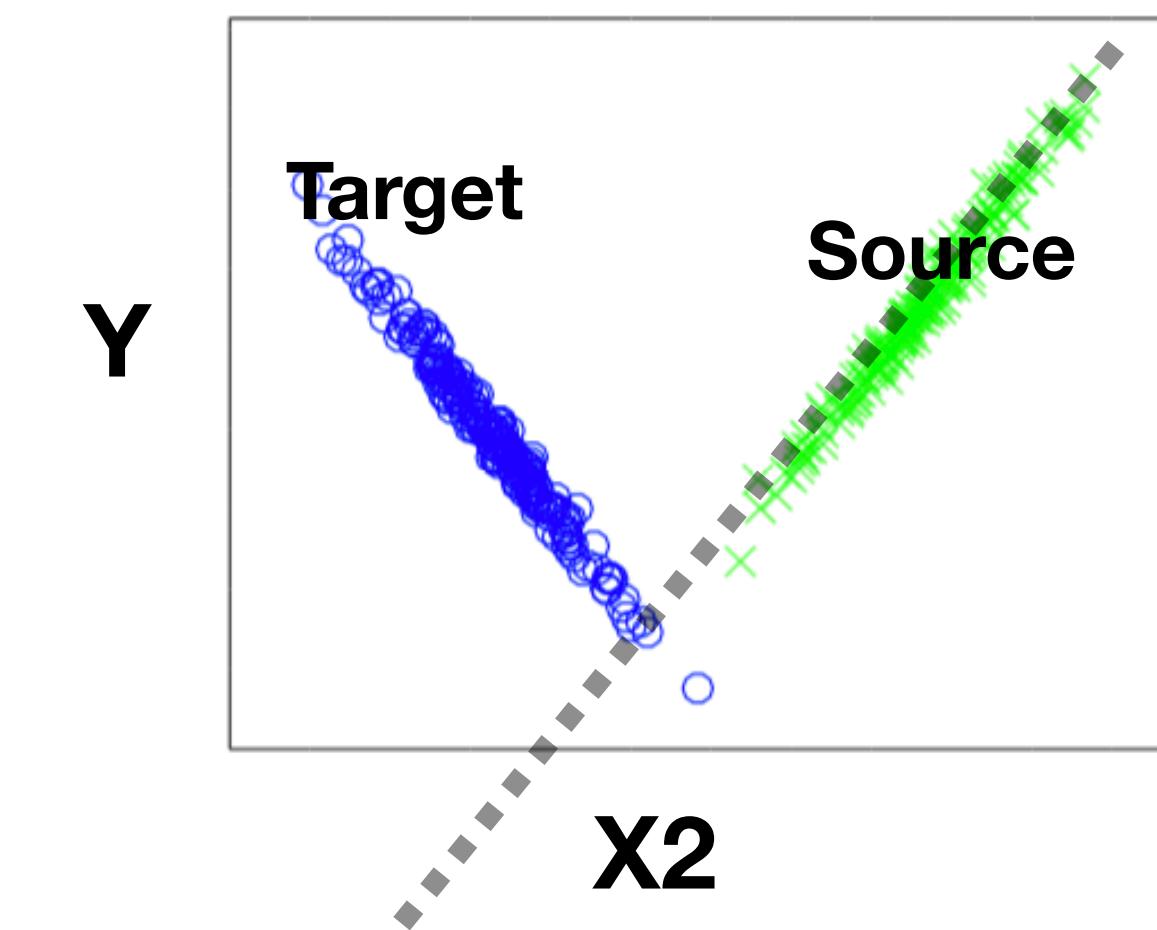
# Separating features = safe for (adversarial) domain adaptation

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



$$Y \perp\!\!\!\perp C | X_1 \equiv$$

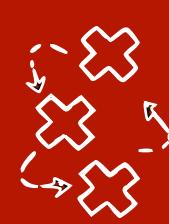
$$P(Y|X_1, C=0) = P(Y|X_1, C=1)$$



$$Y \perp\!\!\!\perp C | X_2 \equiv$$

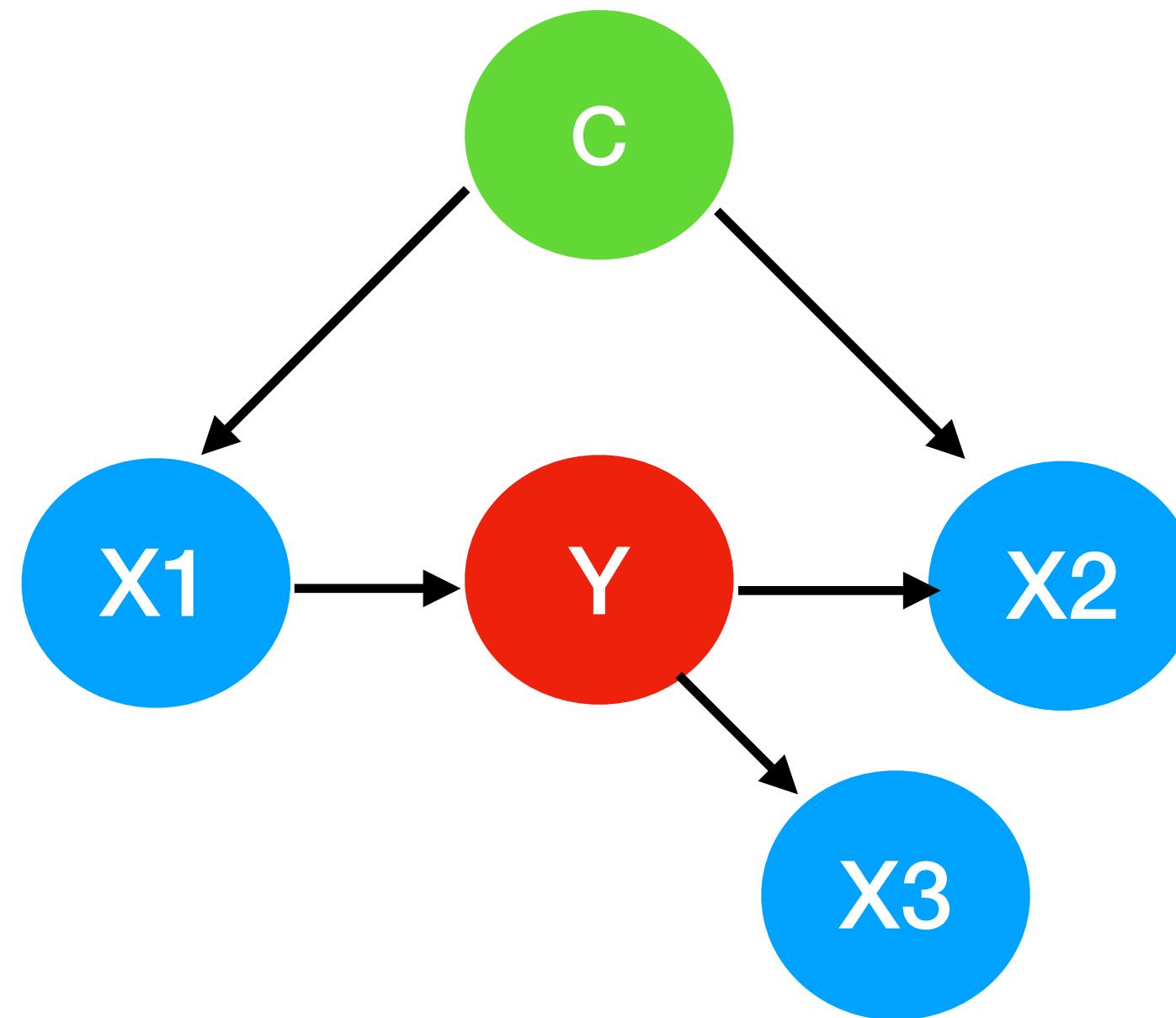
$$P(Y|X_2, C=0) \neq P(Y|X_2, C=1)$$

$\{X_1\}$  is a separating feature,  $\{X_2\}$  and  $\{X_1, X_2\}$  are not  $\rightarrow$  **arbitrarily large error**

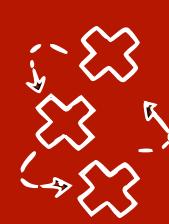


# Separating features = safe for (adversarial) domain adaptation

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

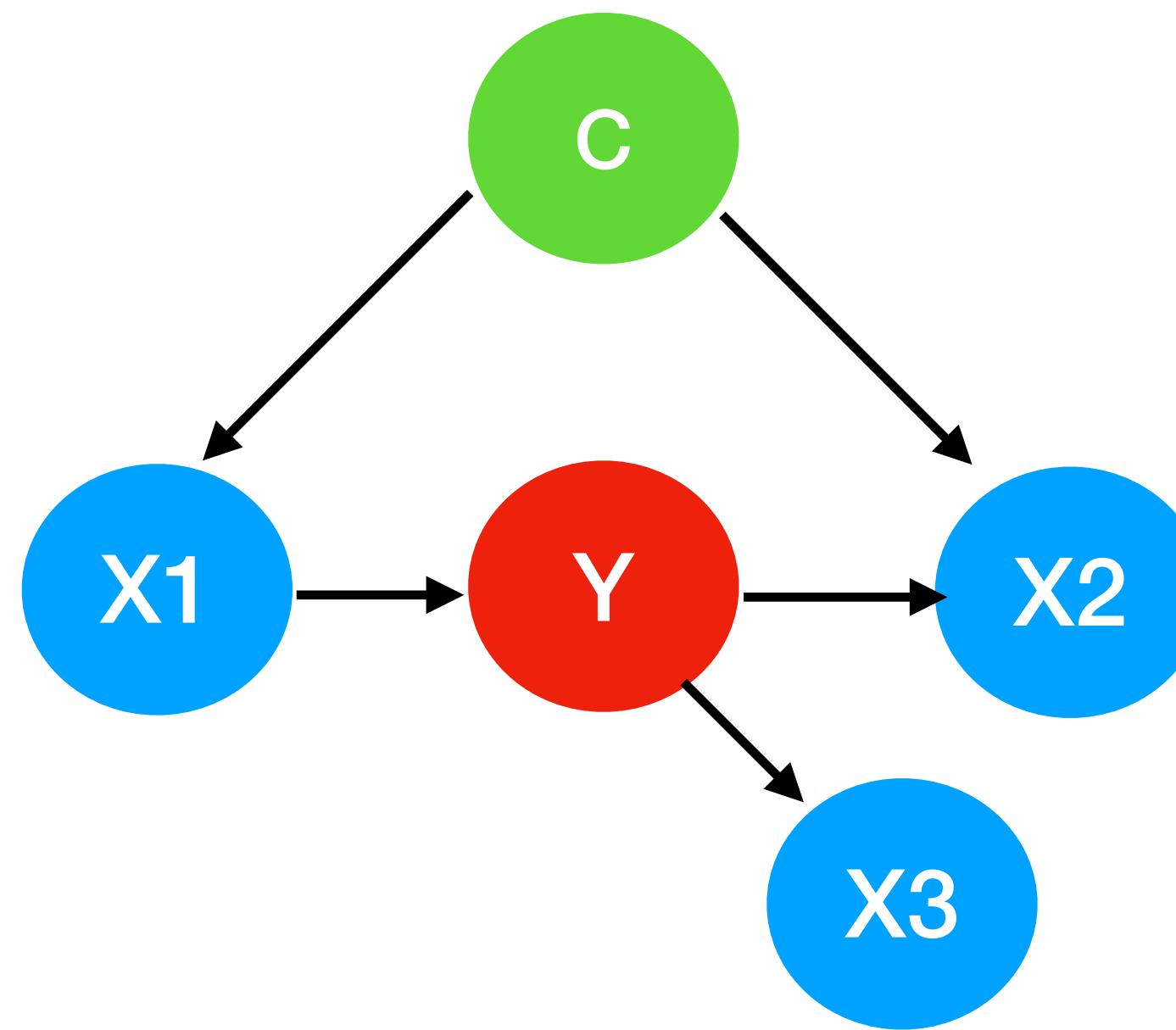


$$Y \perp_d C \mid \{X_1, X_3\} ?$$



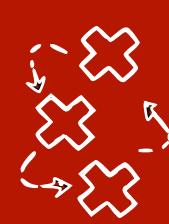
# Separating features = safe for (adversarial) domain adaptation

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



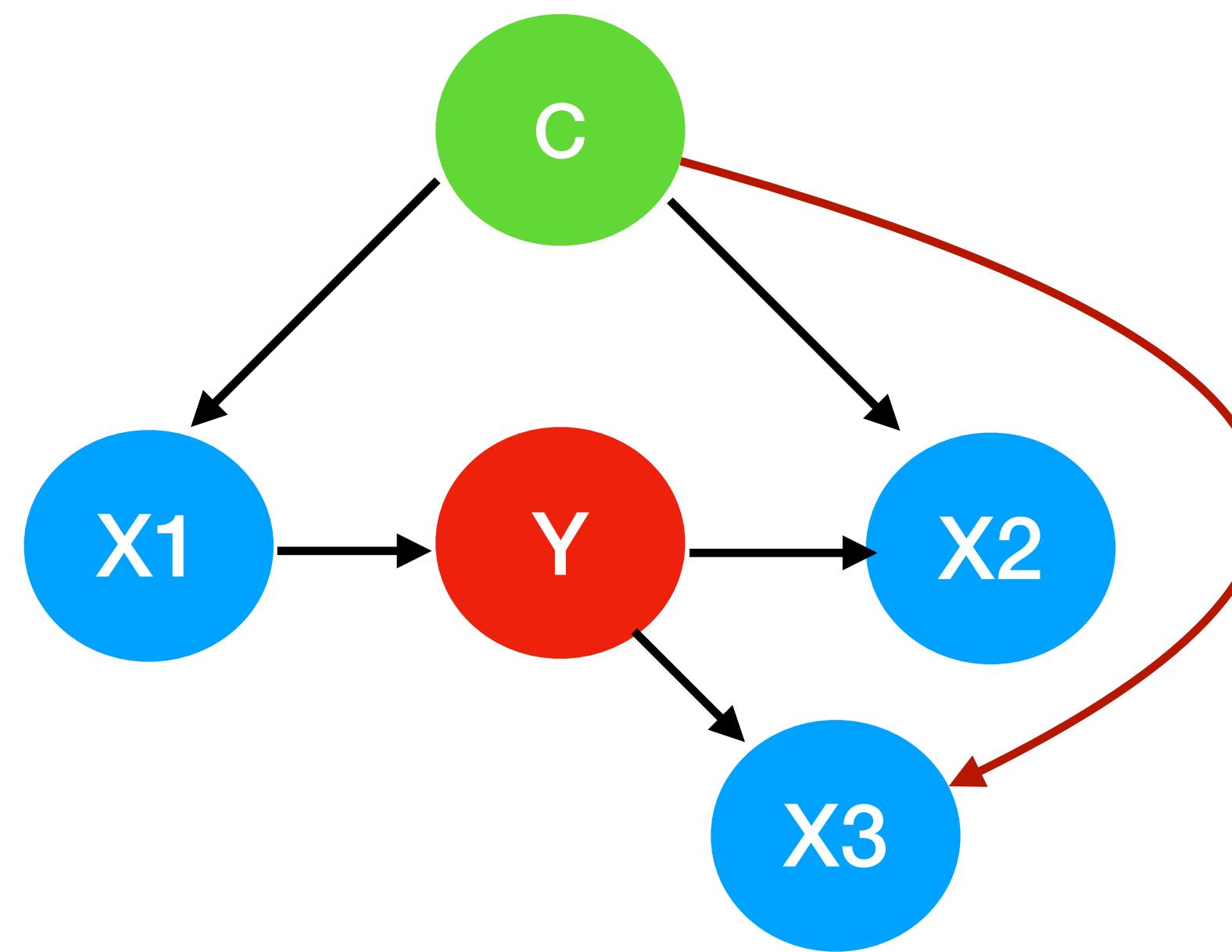
$$Y \perp_d C | \{X_1, X_3\} ?$$

$Y \perp_d C | \{X_1, X_3\}$  (separating features are not necessarily parents of Y)



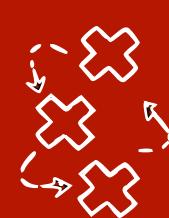
# Separating features = safe for (adversarial) domain adaptation

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



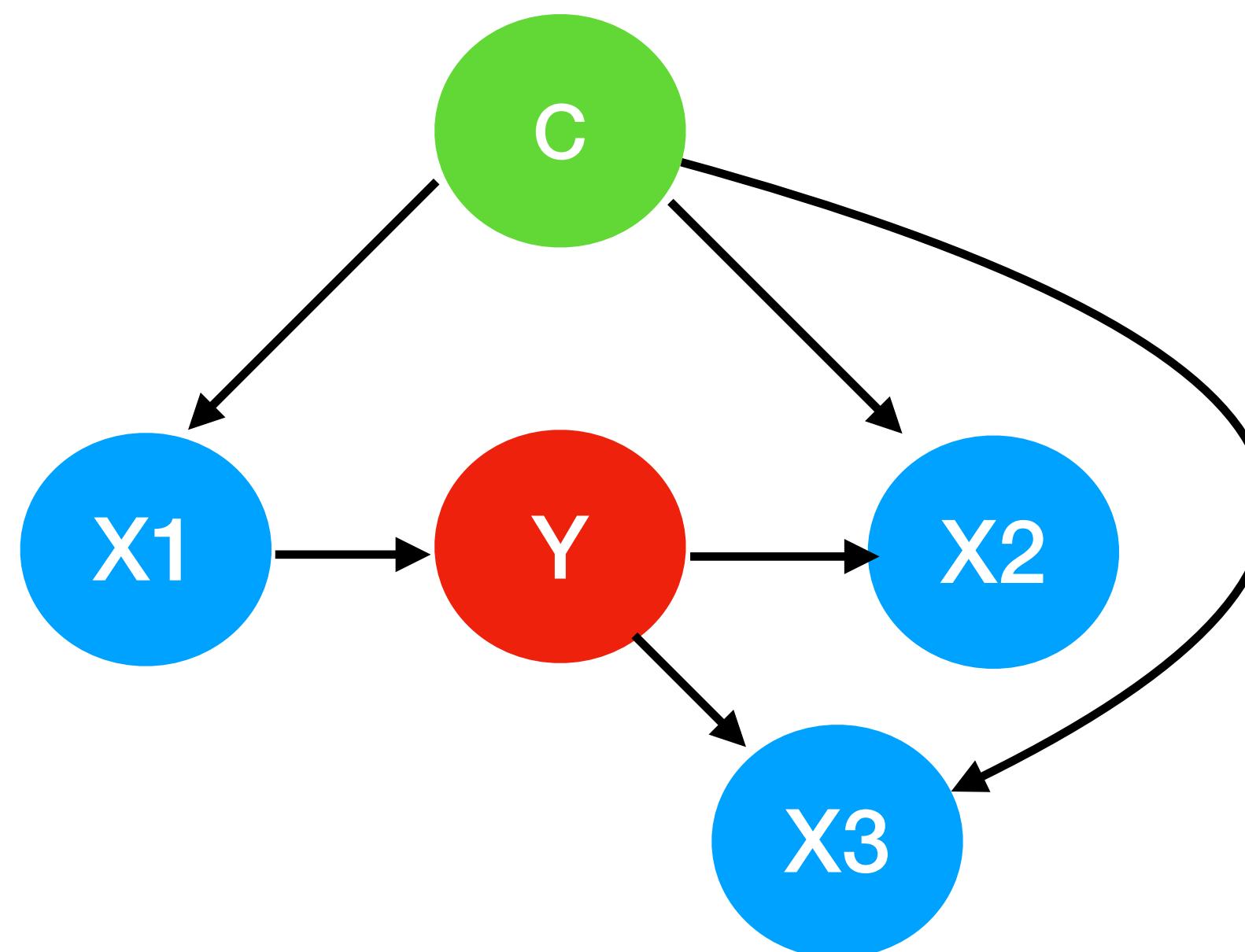
Intervention on every variable except Y =  
domain generalisation

$$Y \perp_d C \mid \{X_1, X_3\} ?$$



# Separating features = safe for (adversarial) domain adaptation

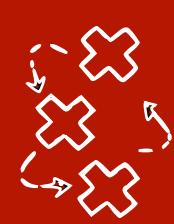
- **Separating features:** sets of features that d-separate Y from the context



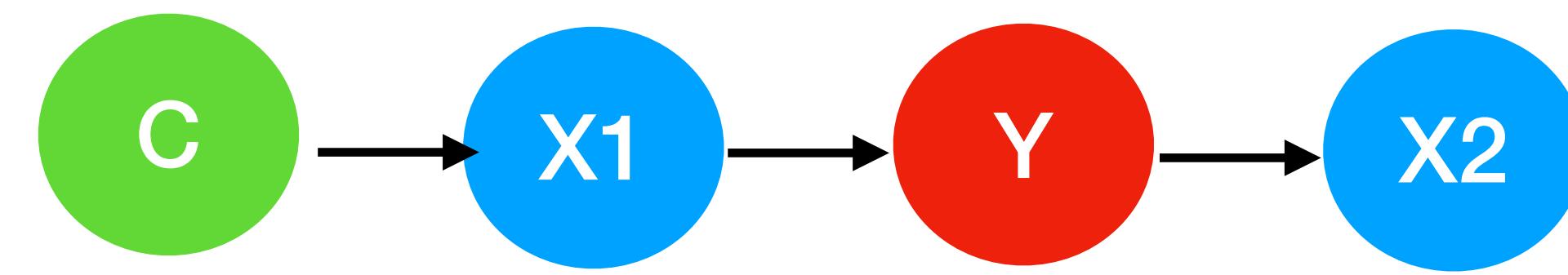
Intervention on every variable except Y =  
domain generalisation

$$Y \perp_d C \mid \{X_1, X_3\} ?$$

Under every possible distribution shift (except directly on the label Y), the separating set of feature are the **causal parents**

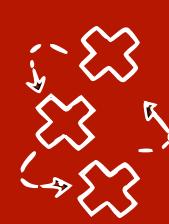


# Misconceptions 1: an invariant feature need not be causal

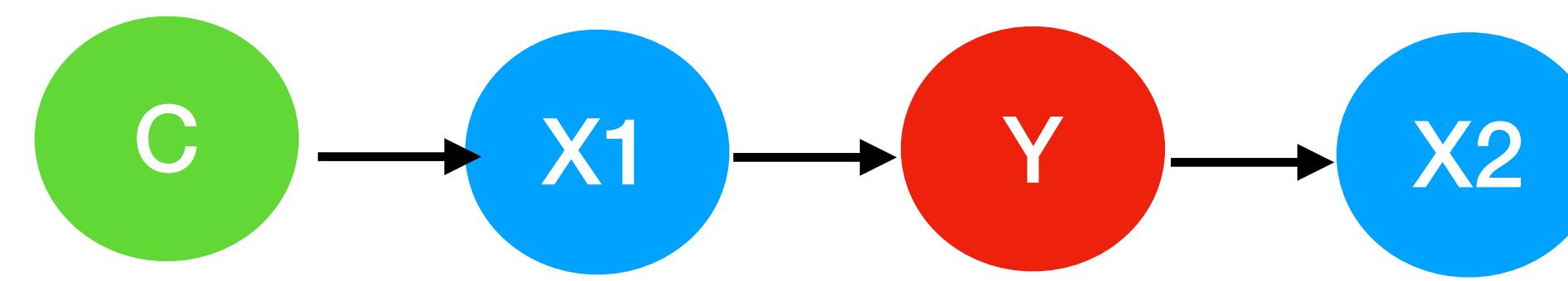


$$\begin{aligned} Y \perp\!\!\!\perp C | X_1 \\ Y \perp\!\!\!\perp C | X_1, X_2 \end{aligned}$$

- $P(Y | X_1, X_2)$  is invariant  $\implies$  invariant features are not necessarily parents of  $Y$



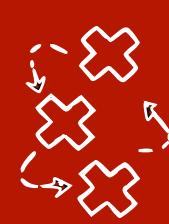
# Misconceptions 1: an invariant feature need not be causal



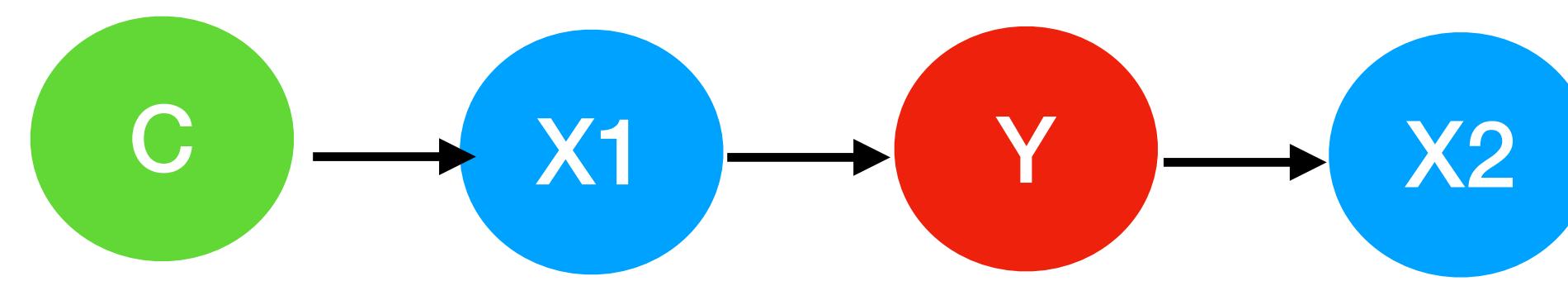
$$\begin{aligned} Y \perp\!\!\!\perp C | X_1 \\ Y \perp\!\!\!\perp C | X_1, X_2 \end{aligned}$$

- $P(Y | X_1, X_2)$  is invariant  $\implies$  invariant features are not necessarily parents of  $Y$

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



# Misconceptions 1: an invariant feature need not be causal



$$\begin{aligned} Y \perp\!\!\!\perp C | X_1 \\ Y \perp\!\!\!\perp C | X_1, X_2 \end{aligned}$$

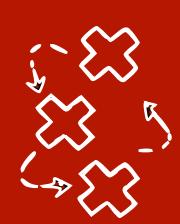
- $P(Y | X_1, X_2)$  is invariant  $\implies$  invariant features are not necessarily parents of  $Y$

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

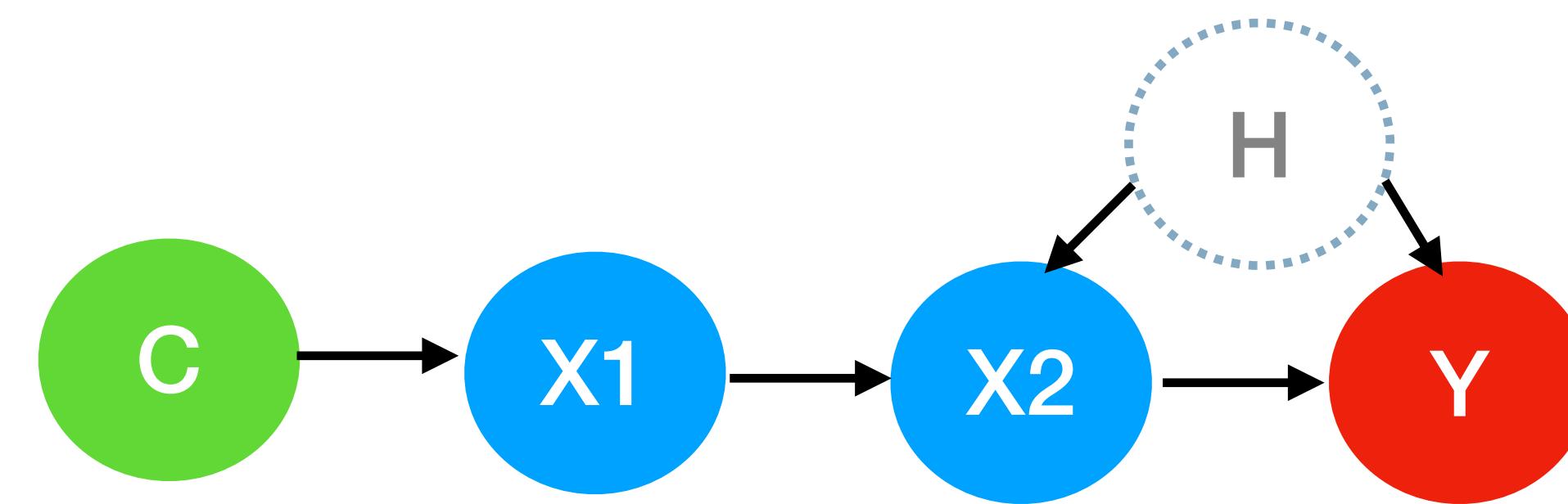
- Invariant Causal Prediction [Peters et al. 2016] under causal sufficiency:

$$S^* = \bigcap_{Y \perp\!\!\!\perp C | S} S \subseteq Pa(Y)$$

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



## Misconceptions 2: observed parents are not enough under latent confounding

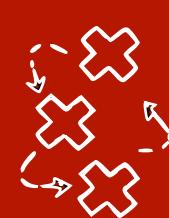


$$Y \perp\!\!\!\perp C | X_1$$

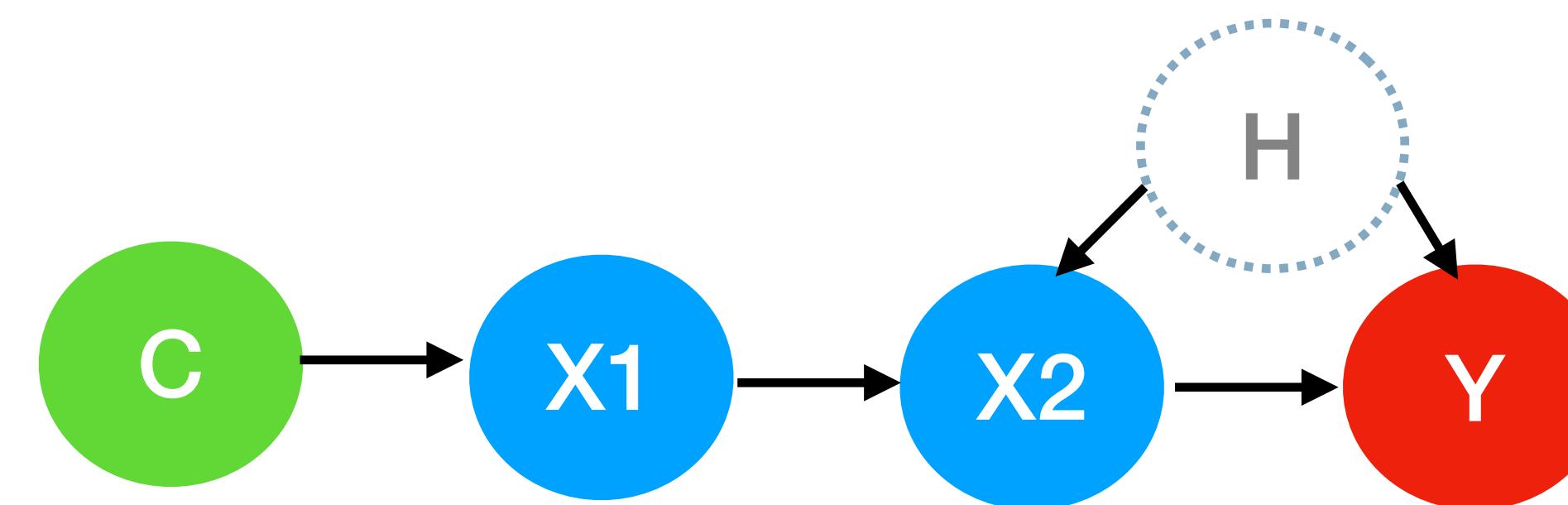
$$Y \not\perp\!\!\!\perp C | X_2$$

$$Y \perp\!\!\!\perp C | X_1, X_2$$

- $P(Y|X_1)$  is invariant,  $P(Y|X_2)$  is not



## Misconceptions 2: observed parents are not enough under latent confounding

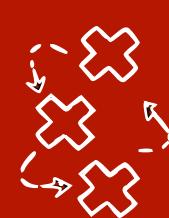


$$Y \perp\!\!\!\perp C | X_1$$

$$Y \not\perp\!\!\!\perp C | X_2$$

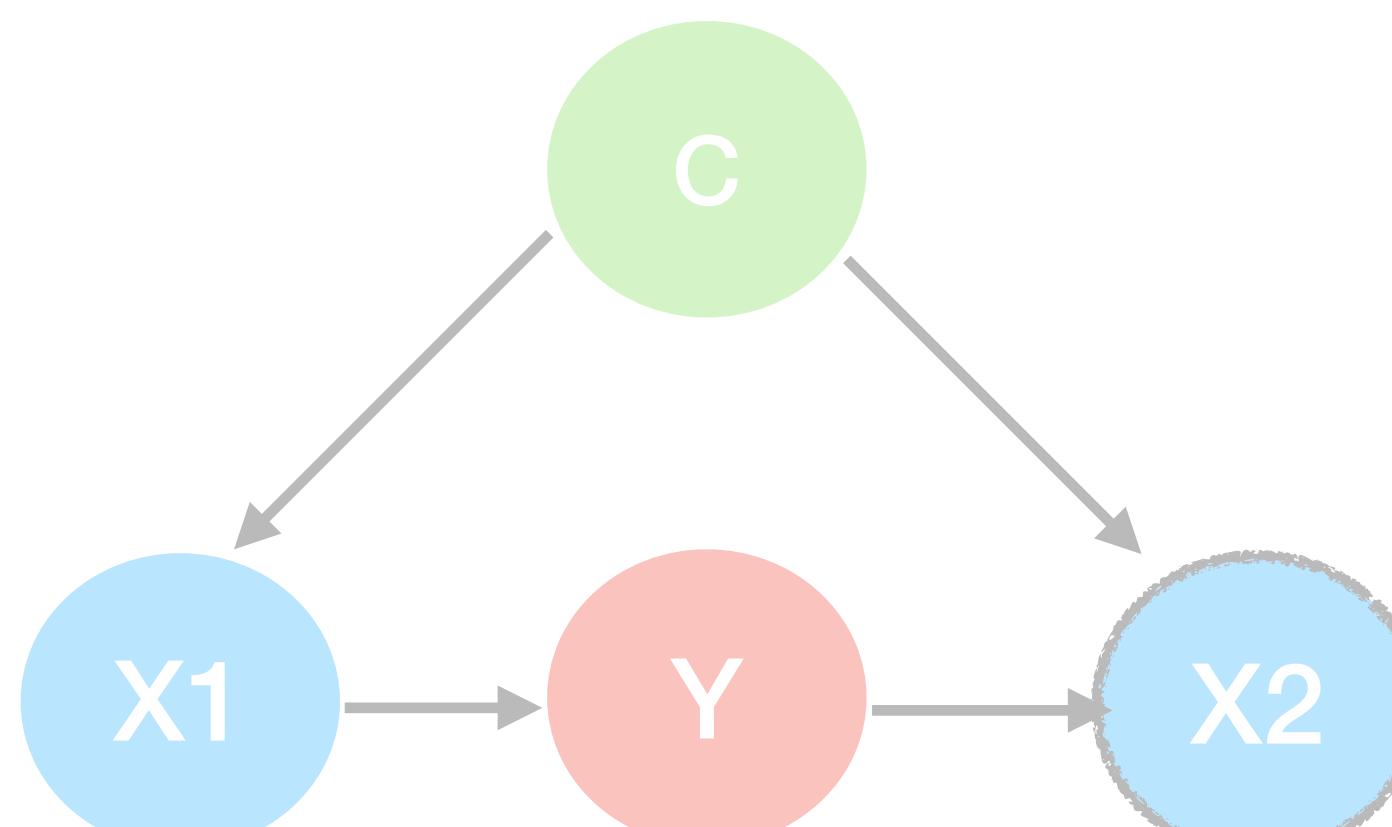
$$Y \perp\!\!\!\perp C | X_1, X_2$$

- $P(Y|X_1)$  is invariant,  $P(Y|X_2)$  is not
- **Conclusion:** causality (e.g. using the causal parents, learning the complete causal graph) is **neither necessary or sufficient\*** for transfer, what we care about are **conditional independences/d-separations**



# Separating features = safe for (adversarial) domain adaptation

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



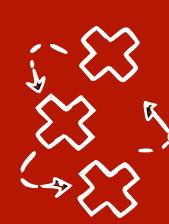
$$Y \perp\!\!\! \perp_d C | X_2$$

$$\begin{aligned} Y \perp\!\!\! \perp C | X_1 \equiv \\ P(Y|X_1, C=0) = P(Y|X_1, C=1) \end{aligned}$$



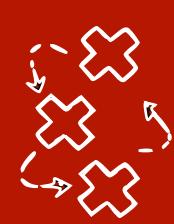
$$\begin{aligned} Y \perp\!\!\! \perp C | X_2 \equiv \\ P(Y|X_2, C=0) \neq P(Y|X_2, C=1) \end{aligned}$$

{X1} is a separating feature, {X2} and {X1, X2} are not -> **arbitrarily large error**



# Desiderata for a causal domain adaptation method

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

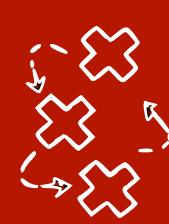


# Unsupervised multi-source domain adaptation - toy example



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	?

Source domains      Target domain

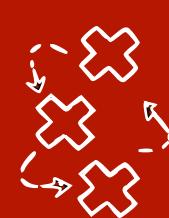


# Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions NeurIPS 2018

Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, Joris M. Mooij

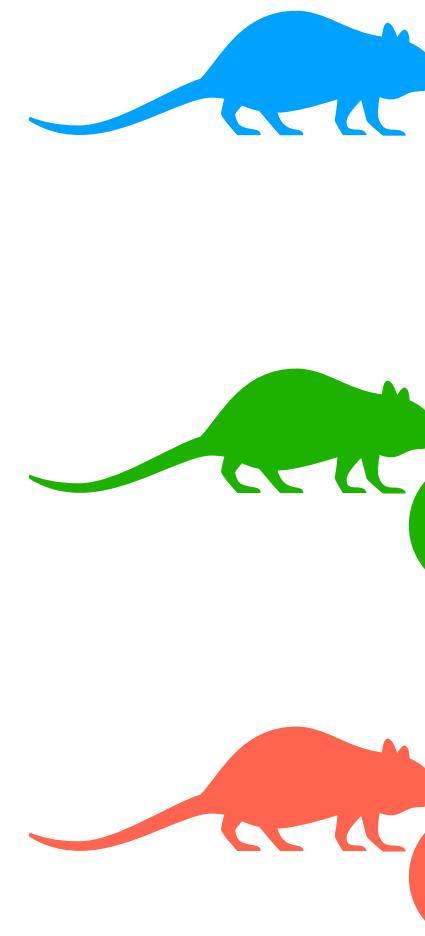
C1	C2	X1	X2	Y
0	0	0,1	1	0
0	0	0,2	1	0
0	0	1,1	2	1
0	1	3,1	2	1
0	1	3,2	3	1
0	1	4	3	1
1	0	0,2	0	?
1	0	0,3	0	?
1	0	0,3	1	?

- We assume we can model all the domains in with a **single unknown acyclic causal graph** with **multiple context variables** [Mooij et al. 2020]



# Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions NeurIPS 2018

Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, Joris M. Mooij

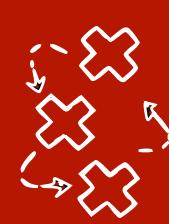


C1	C2	X1	X2	Y
0	0	0,1	1	0
0	0	0,2	1	0
0	0	1,1	2	1
0	1	3,1	2	1
0	1	3,2	3	1
0	1	4	3	1
1	0	0,2	0	?
1	0	0,3	0	?
1	0	0,3	1	?

Step 1: Known causal graph

- Find features  $S \subseteq X : Y \perp_d C_1 | S$

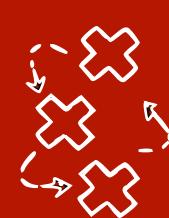
- We assume we can model all the domains in with a **single unknown acyclic causal graph** with **multiple context variables** [Mooij et al. 2020]



# What if the causal graph is unknown?

- **Idea:** we could test the conditional independence in the data

$$Y \perp\!\!\!\perp C_1 | X_1? \quad Y \perp\!\!\!\perp C_1 | X_2?$$



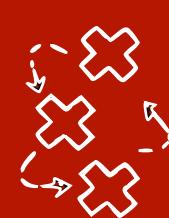
# What if the causal graph is unknown?

- **Idea:** we could test the conditional independence in the data

$$Y \perp\!\!\!\perp C_1 | X_1 ? \quad Y \perp\!\!\!\perp C_1 | X_2 ?$$

- **Problem:** Y is always missing when  $C_1=1$ , so we cannot test these

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	0,2	0	0
0	1	0,3	0	1
0	1	0,3	1	0
1	0	3,1	2	?
1	0	3,2	3	?
1	0	4	3	?



# What if the causal graph is unknown?

- **Idea:** we could test the conditional independence in the data

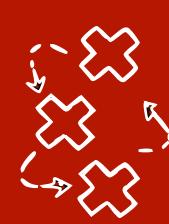
$$\begin{array}{l} Y \perp\!\!\!\perp C_1 | X_1 ? \\ Y \perp\!\!\!\perp C_1 | X_2 ? \end{array}$$

- **Problem:** Y is always missing when  $C_1=1$ , so we cannot test these

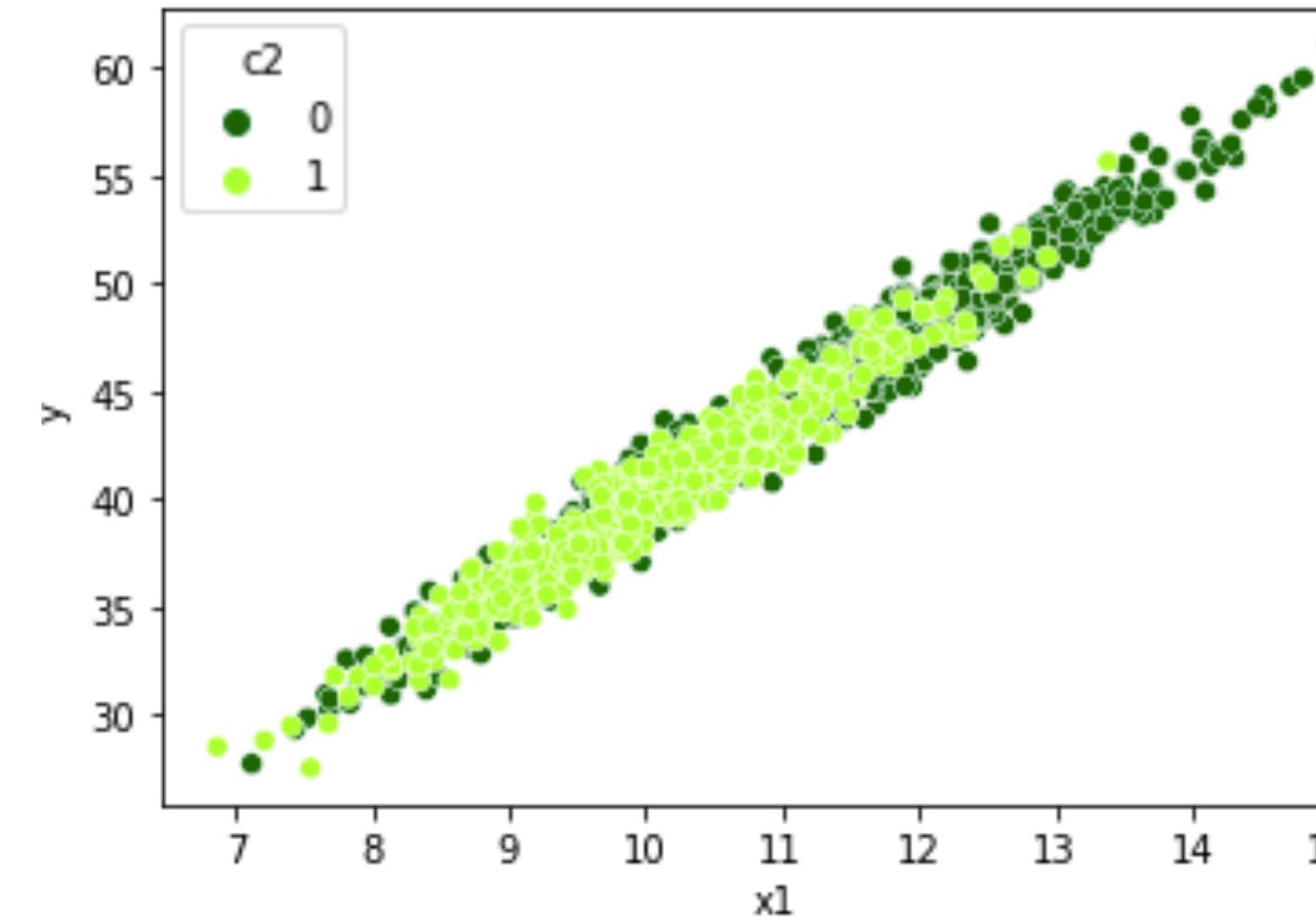
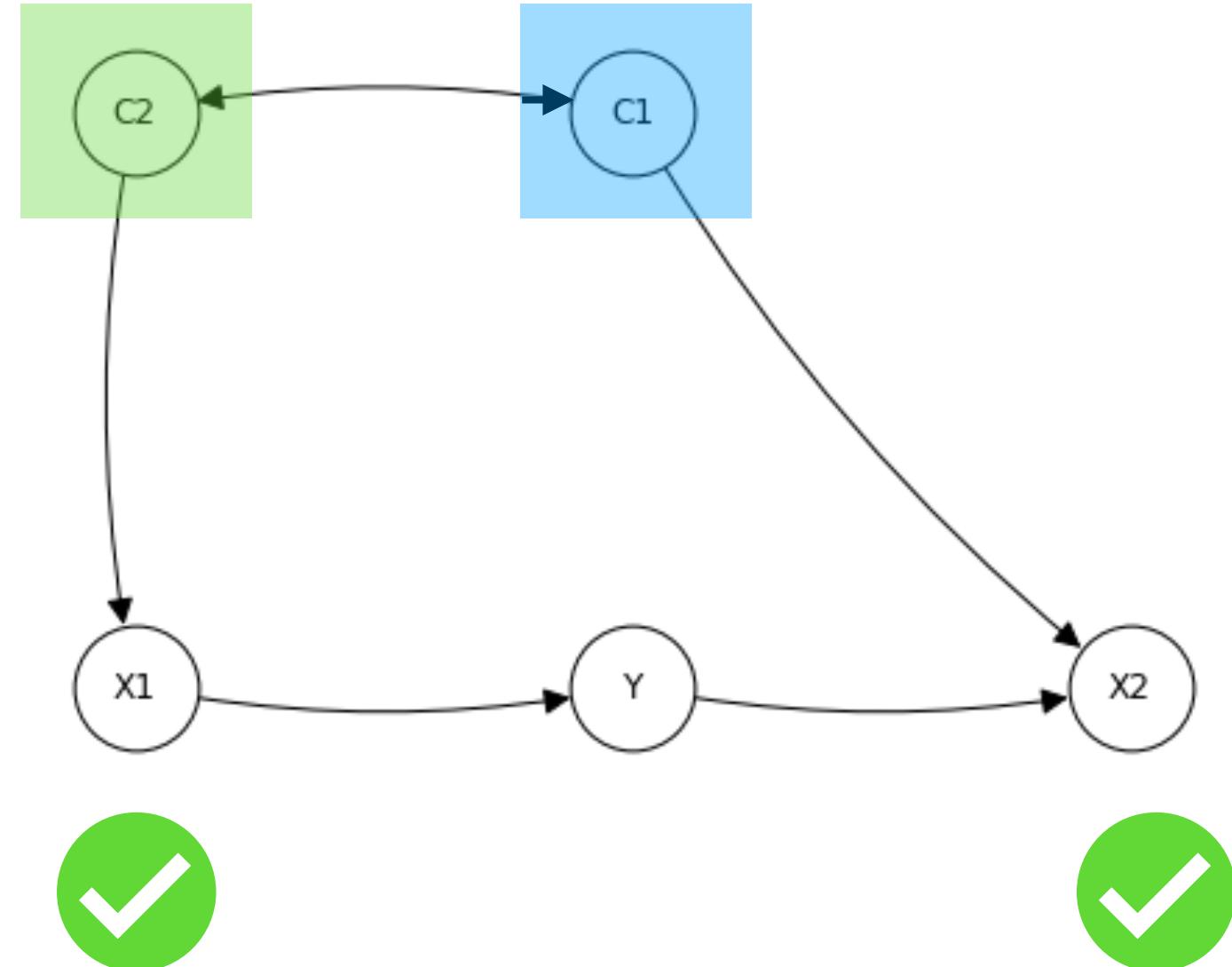
C1	C2	X1	X2	Y
0	0	0,1	1	0
0	0	0,2	1	0
0	0	1,1	2	1
0	1	0,2	0	0
0	1	0,3	0	1
0	1	0,3	1	0

**Idea:** Invariant features in source domains are also separating in the target domain

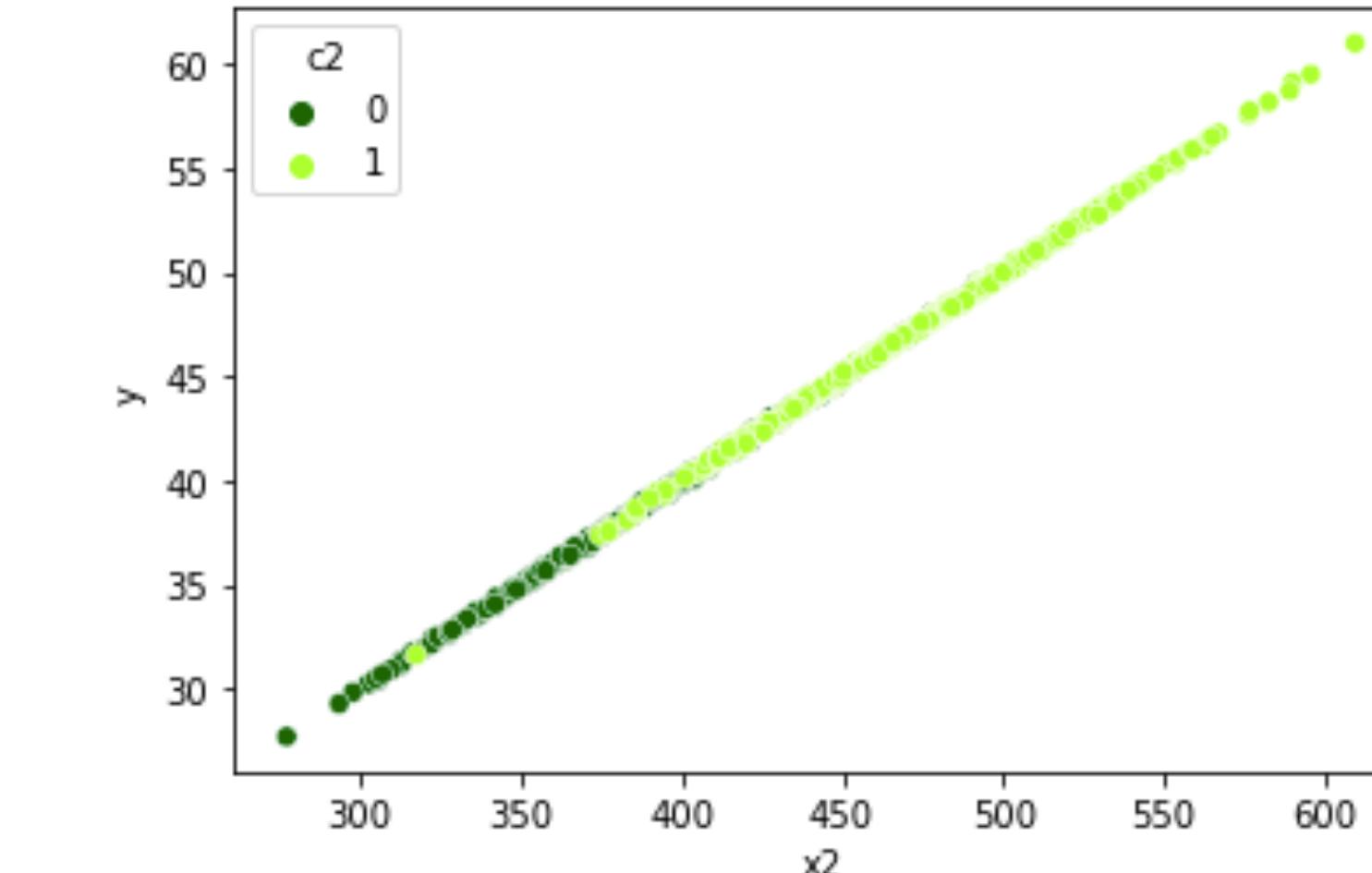
$$Y \perp\!\!\!\perp C_2 | \{X_1, C_1 = 0\} \implies Y \perp\!\!\!\perp C_1 | X_1$$



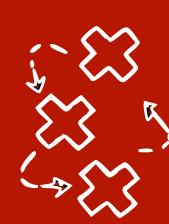
# Separating features in sources are also separating in target - counterexample



$$Y \perp\!\!\!\perp C_2 | \{X_1, C_1 = 0\}$$

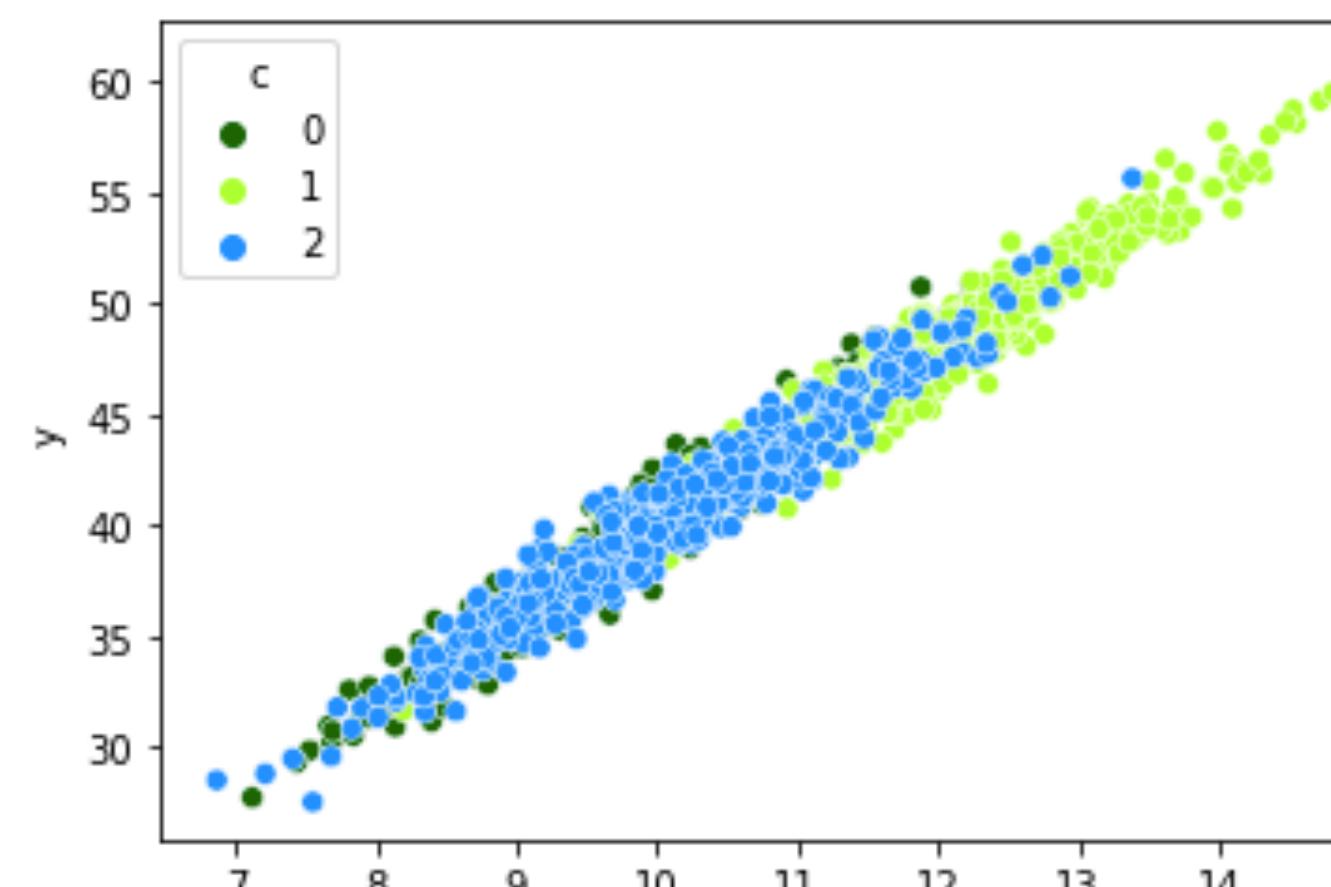
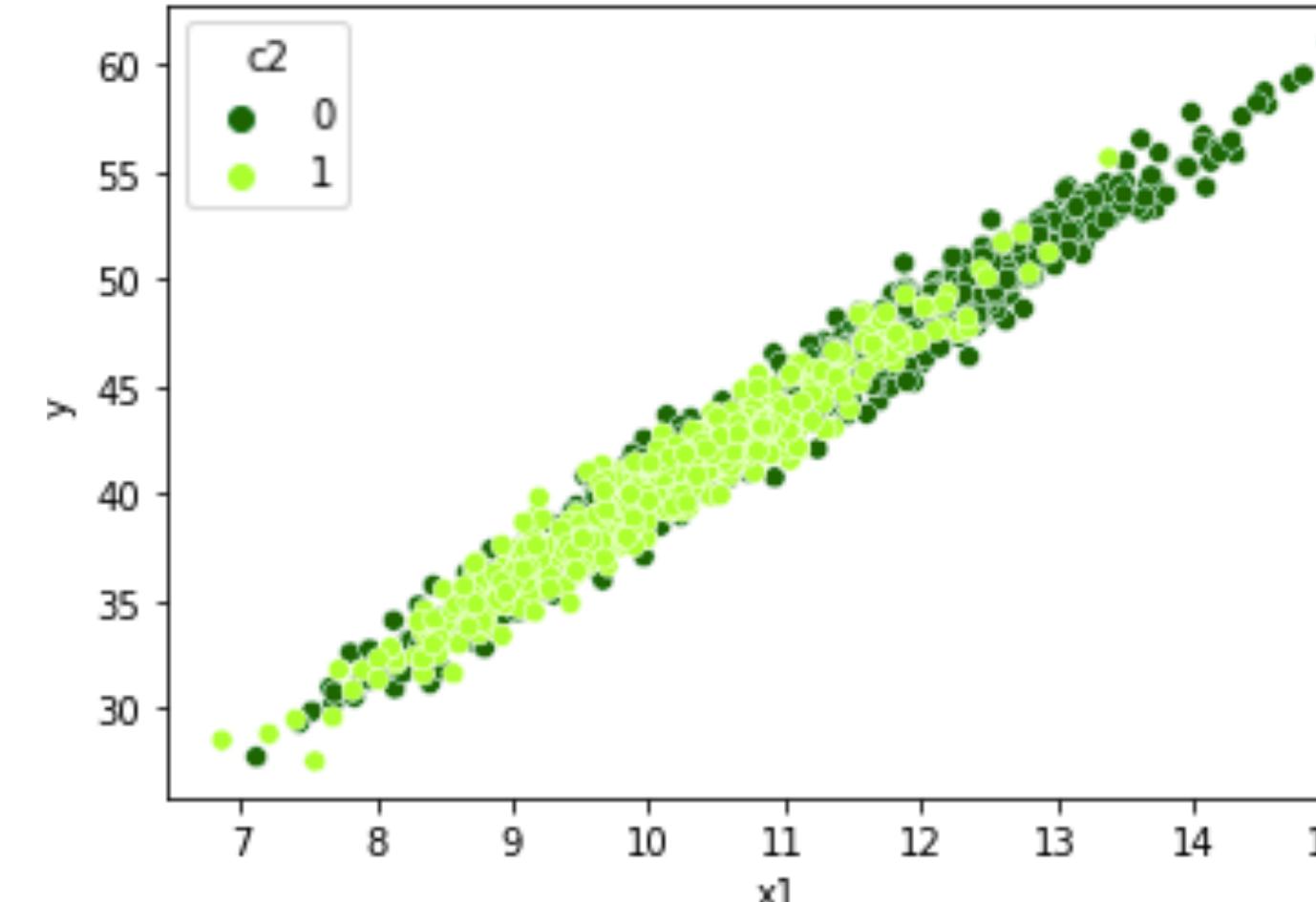
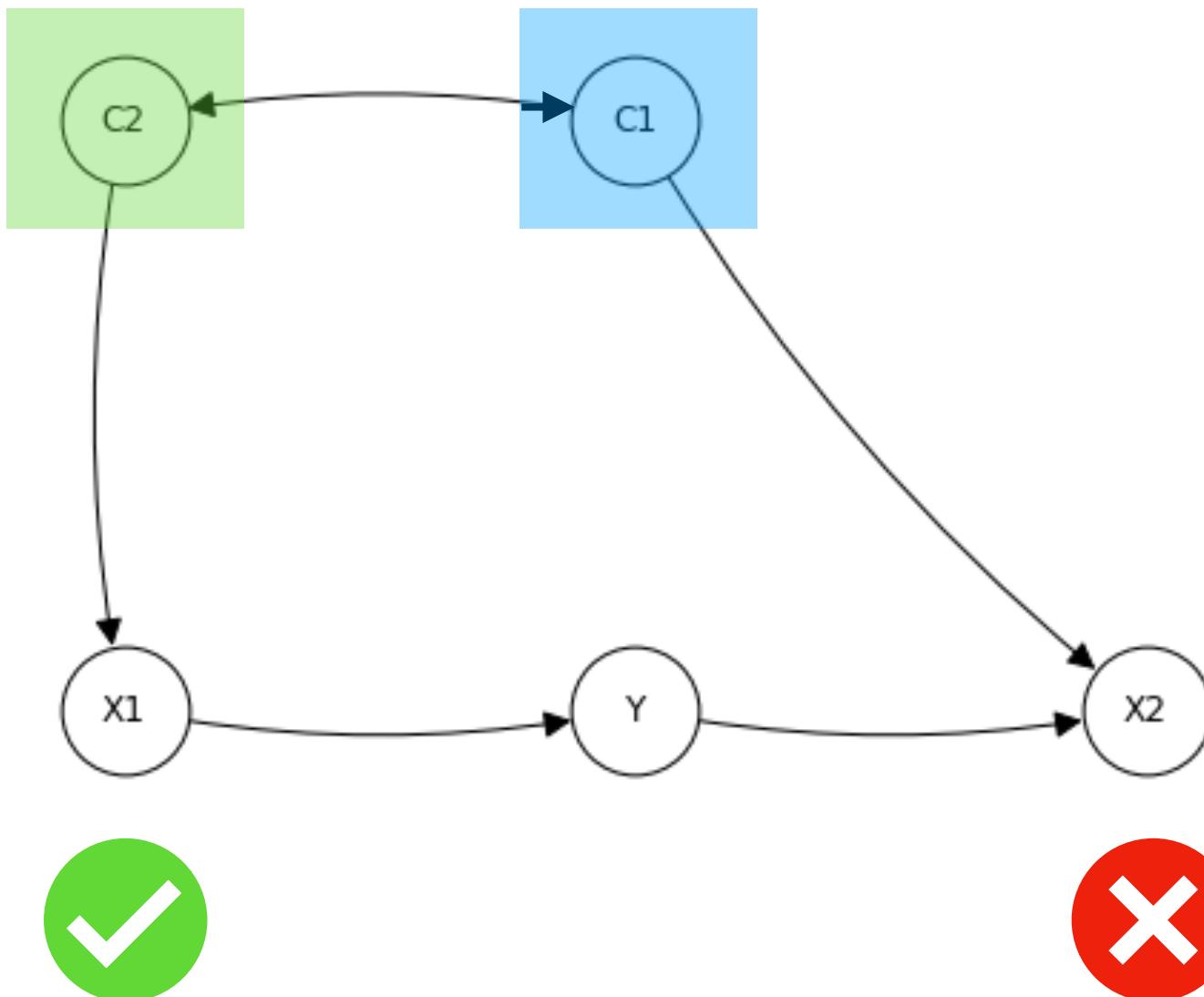


$$Y \perp\!\!\!\perp C_2 | \{X_2, C_1 = 0\}$$



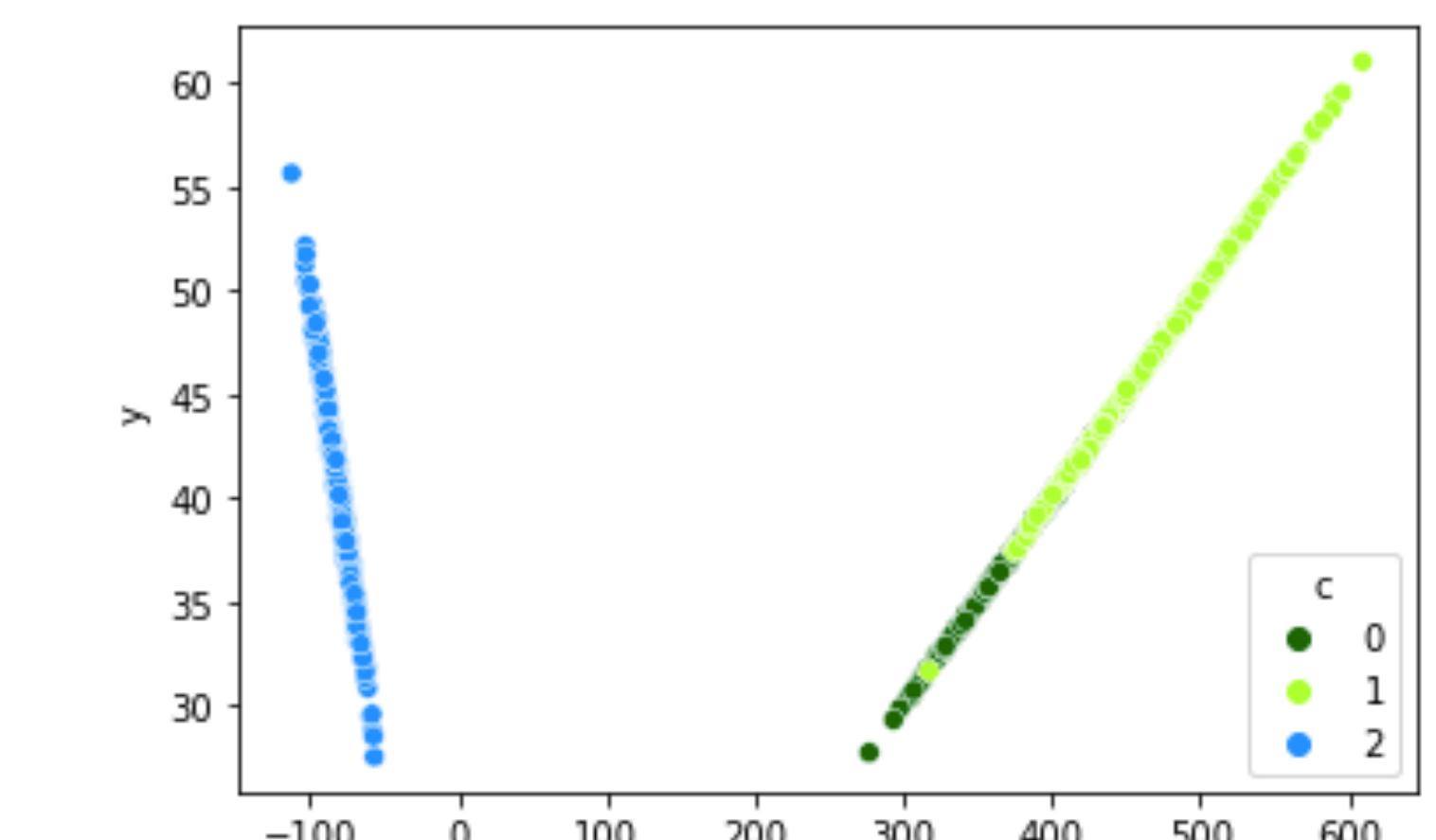
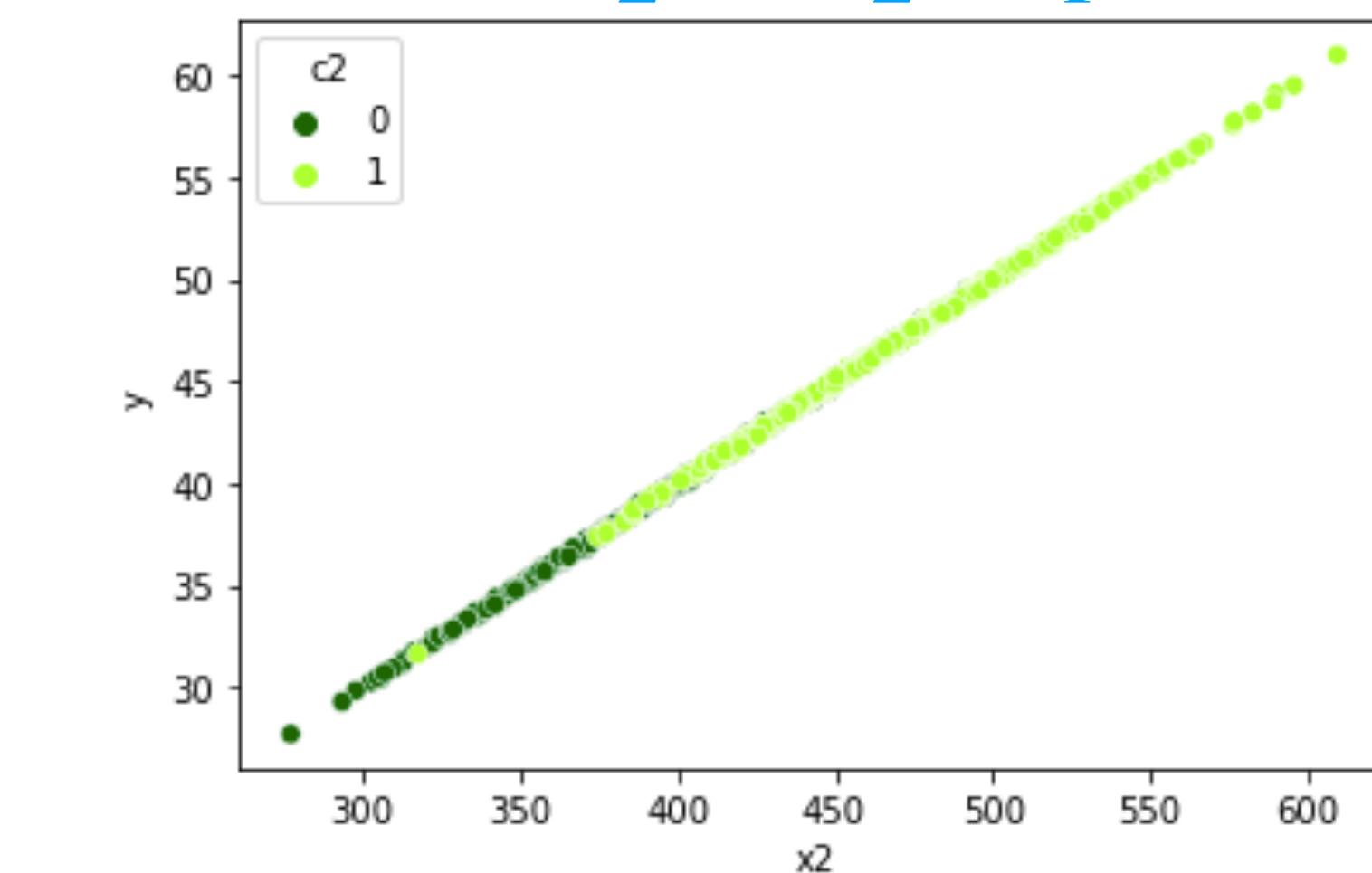
# Separating features in sources are also separating in target - counterexample

$$Y \perp\!\!\!\perp C_2 | \{X_1, C_1 = 0\}$$

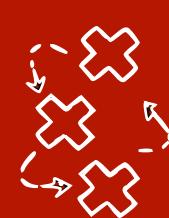


$$Y \perp\!\!\!\perp C_1 | X_1$$

$$Y \perp\!\!\!\perp C_2 | \{X_2, C_1 = 0\}$$



$$Y \perp\!\!\!\perp C_1 | X_2$$



# What if the causal graph is unknown?

- **Idea:** we could test the conditional independence in the data

$$\begin{array}{l} Y \perp\!\!\!\perp C_1 | X_1 ? \\ Y \perp\!\!\!\perp C_1 | X_2 ? \end{array}$$

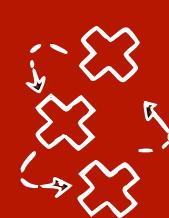
- **Problem:** Y is always missing when  $C_1=1$ , so we cannot test these

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	0,2	0	0
0	1	0,3	0	1
0	1	0,3	1	0

**Idea:** Invariant features in source domains are also separating in the target domain

$$Y \perp\!\!\!\perp C_2 | \{X_1, C_1 = 0\} \implies Y \perp\!\!\!\perp C_1 | X_1$$

This is a strong assumption



# What if the causal graph is unknown?

- **Idea:** we could test the conditional independence in the data

$$Y \perp\!\!\!\perp C_1 | X_1 ? \quad Y \perp\!\!\!\perp C_1 | X_2 ?$$

- **Problem:** Y is always missing when  $C_1=1$ , so we cannot test these

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	0,2	0	0
0	1	0,3	0	1
0	1	0,3	1	0
1	0	3,1	2	?
1	0	3,2	3	?
1	0	4	3	?

$$X_1 \perp\!\!\!\perp X_2$$

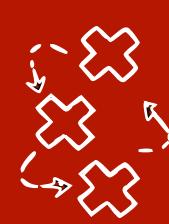
$$X_1 \perp\!\!\!\perp C_1$$

$$X_1 \perp\!\!\!\perp X_2 | C_1$$

$$X_1 \perp\!\!\!\perp X_2 | Y, C_1 = 0$$

...

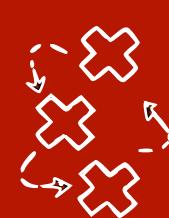
- **Idea:** Can we use all other in/dependences?



# Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions NeurIPS 2018

Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, Joris M. Mooij

- We search for **separating features** that d-separate  $Y$  from  $C_1$  (target)
- We assume **no extra dependences involving  $Y$**  in target domain  $C_1=1$



# Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions NeurIPS 2018

Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, Joris M. Mooij

- We search for **separating features** that d-separate Y from  $C_1$  (target)
- We assume **no extra dependences involving Y** in target domain  $C_1=1$

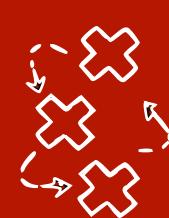
C1	C2	X1	X2	Y
0	0	0,1	1	0
0	0	0,2	1	0
0	0	1,1	2	1
0	1	3,1	2	1
0	1	3,2	3	1
0	1	4	3	1
1	0	0,2	0	?
1	0	0,3	0	?
1	0	0,3	1	?

$$Y \perp\!\!\!\perp C_2 | C_1 = 0$$

$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

Perform allowed CI tests



# Domain Adaptation by Using Causal Inference to Predict Invariant Conditional Distributions NeurIPS 2018

Sara Magliacane, Thijs van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, Joris M. Mooij

- We search for **separating features** that d-separate  $Y$  from  $C_1$  (target)
- We assume **no extra dependences involving  $Y$**  in target domain  $C_1=1$

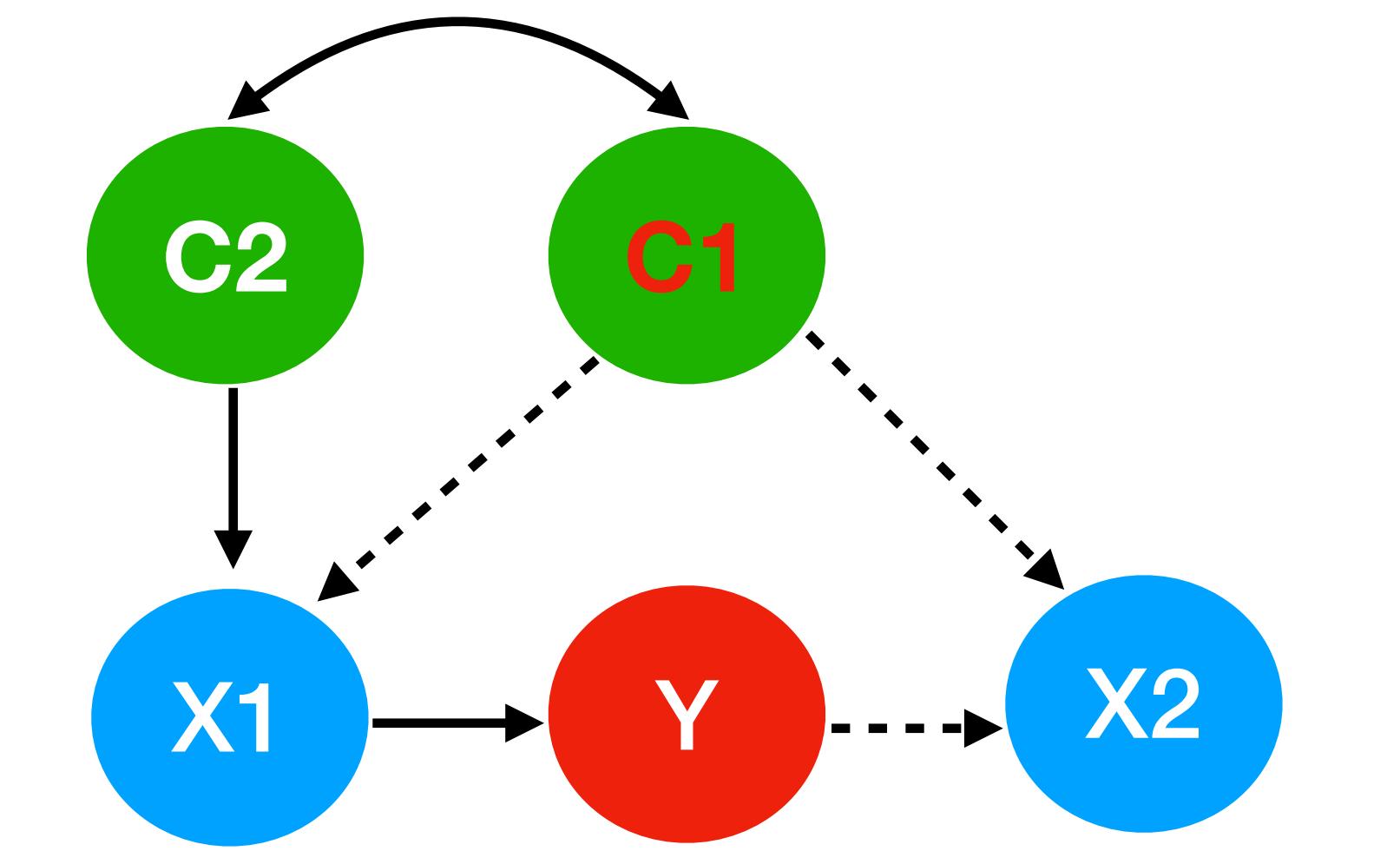
C1	C2	X1	X2	Y
0	0	0,1	1	0
0	0	0,2	1	0
0	0	1,1	2	1
0	1	3,1	2	1
0	1	3,2	3	1
0	1	4	3	1
1	0	0,2	0	?
1	0	0,3	0	?
1	0	0,3	1	?

$$Y \perp\!\!\!\perp C_2 | C_1 = 0$$

$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

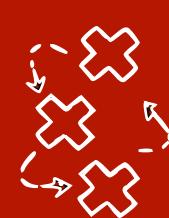
$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

Perform allowed CI tests

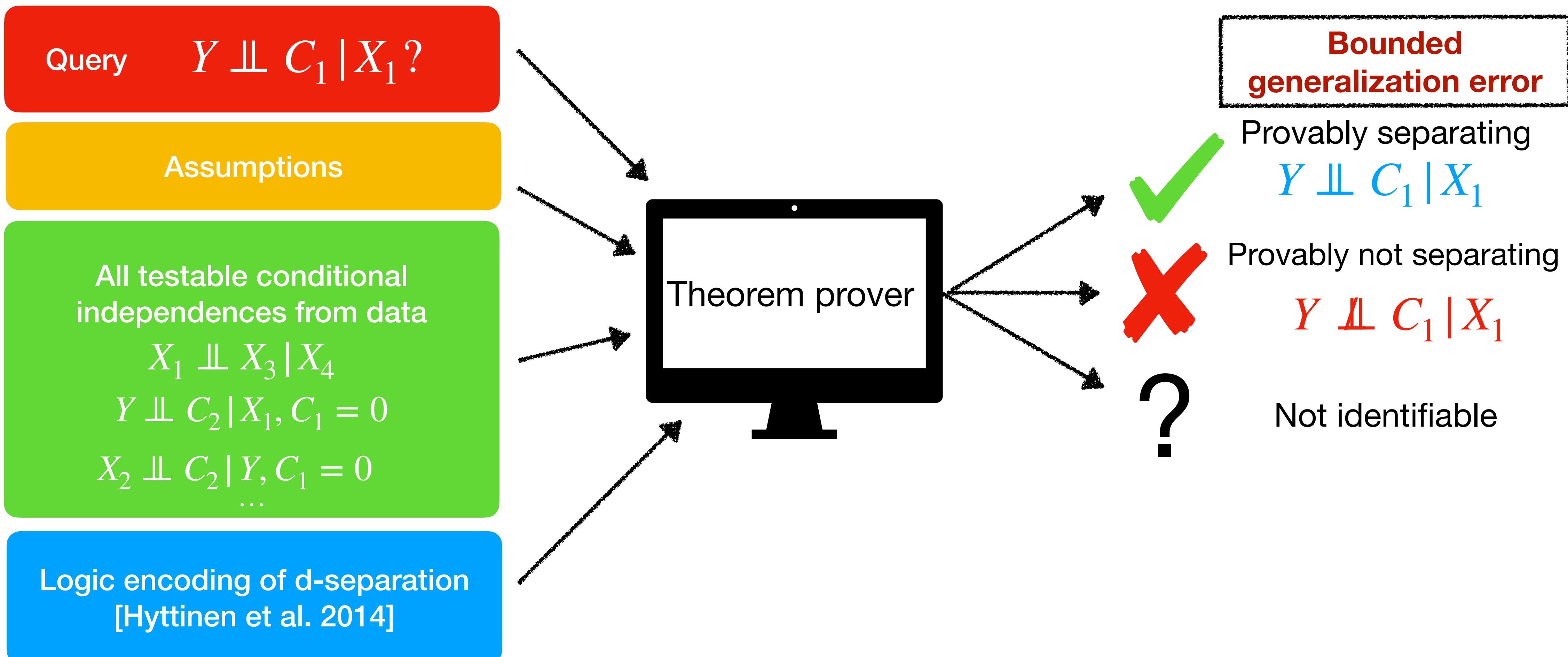


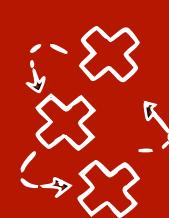
All possible compatible graphs

$$Y \perp\!\!\!\perp C_1 | X_1 ?$$

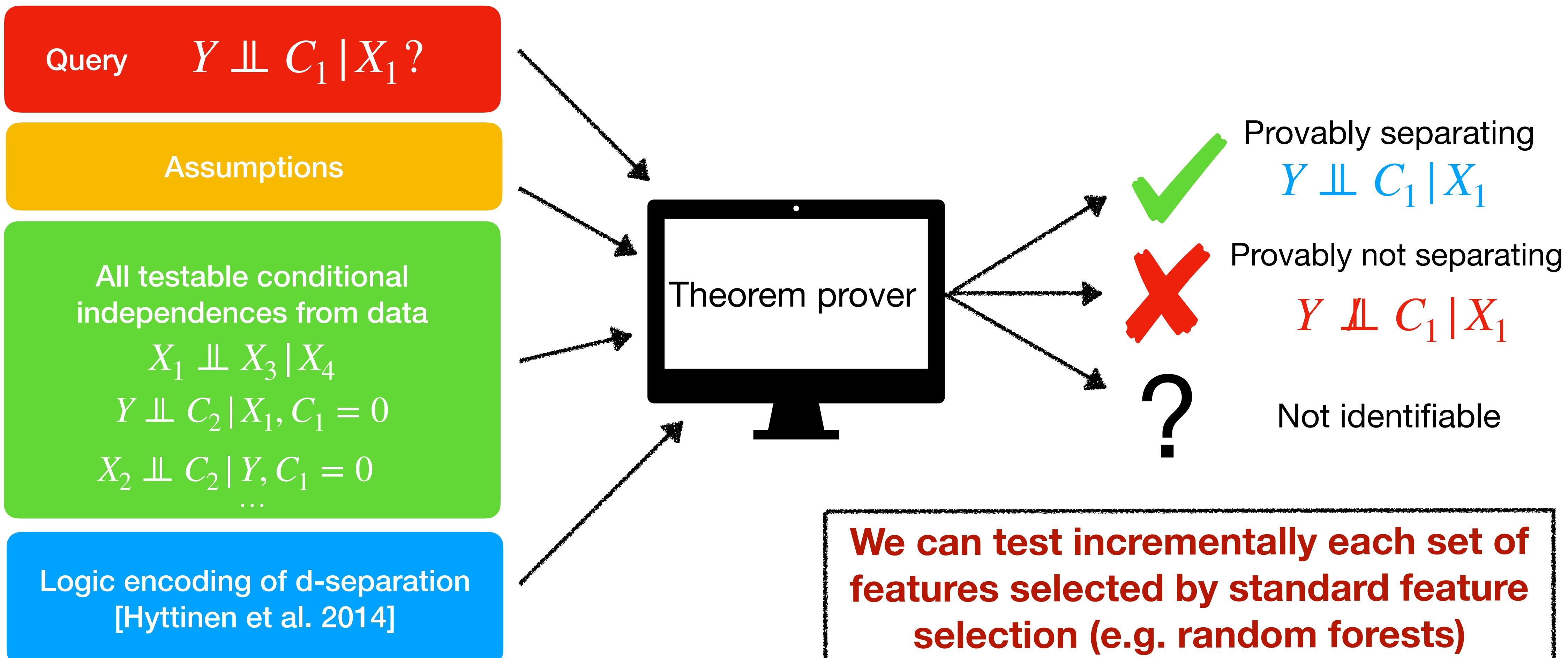


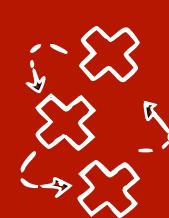
# Inferring separating sets of features





# Inferring separating sets of features





# A simple causal feature selection algorithm

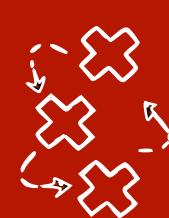
Source domains data

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



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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$



# 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

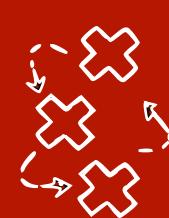
Standard feature selection

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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

$$S = \{X_1, C_2\}$$



# 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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

$$S = \{X_1, C_2\}$$

All data (including target)

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	?

Query  $Y \perp\!\!\!\perp C_1 | S$ ?

Assumptions

All testable conditional independences from data

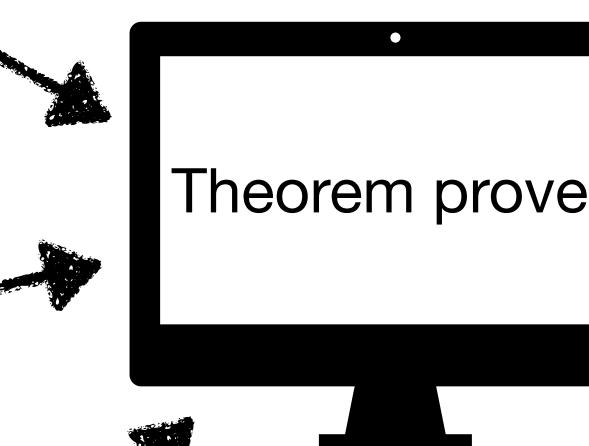
$$X_1 \perp\!\!\!\perp X_3 | X_4$$

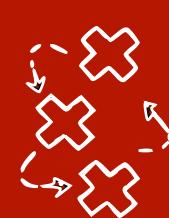
$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

...

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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

$$S = \{X_1, C_2\}$$

All data (including target)

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	?

Query  $Y \perp\!\!\!\perp C_1 | S$ ?

Assumptions

All testable conditional independences from data

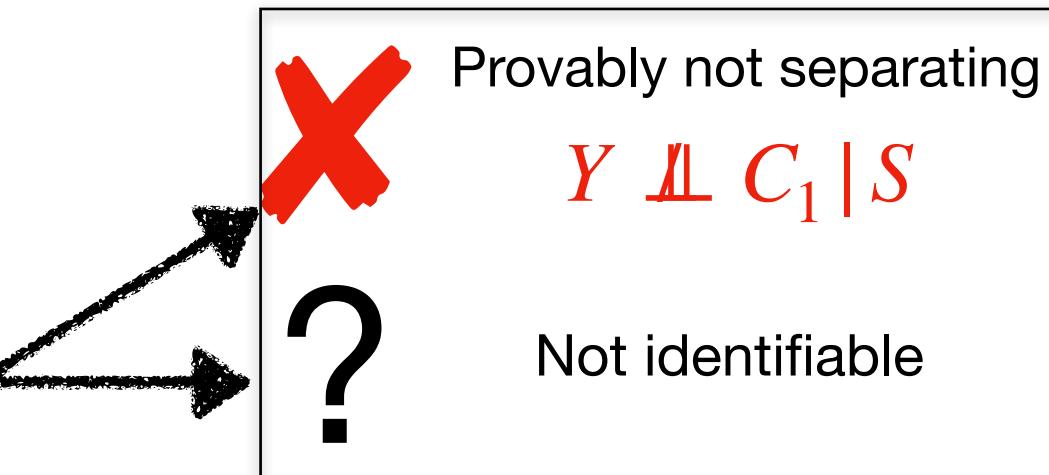
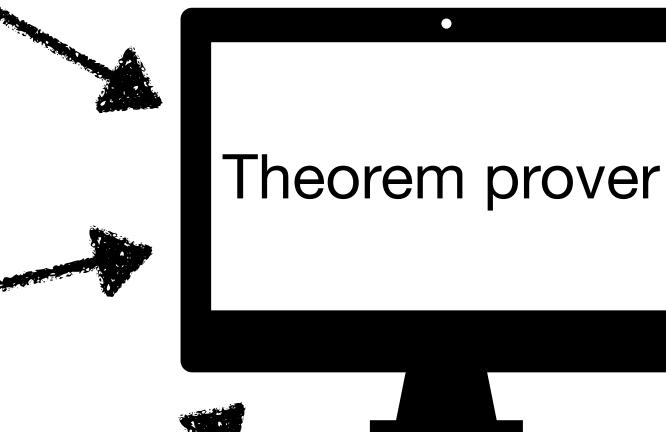
$$X_1 \perp\!\!\!\perp X_3 | X_4$$

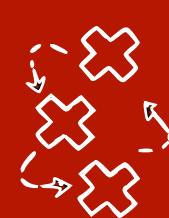
$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

...

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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

$$S = \{X_1, X_2, C_2\}$$

All data (including target)

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	?

Query  $Y \perp\!\!\!\perp C_1 | S$ ?

Assumptions

All testable conditional independences from data

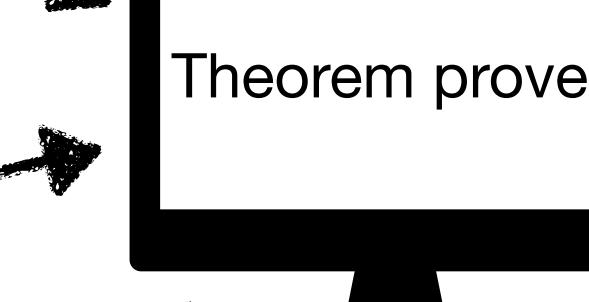
$$X_1 \perp\!\!\!\perp X_3 | X_4$$

$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

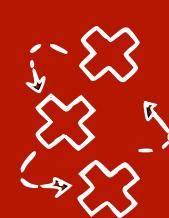
$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

...

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



Provably not separating  
 $Y \perp\!\!\!\perp C_1 | S$   
Not identifiable



# 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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

$$S = \{X_1, X_2, C_2\}$$

All data (including target)

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	?

Query  $Y \perp\!\!\!\perp C_1 | S$ ?

Assumptions

All testable conditional independences from data

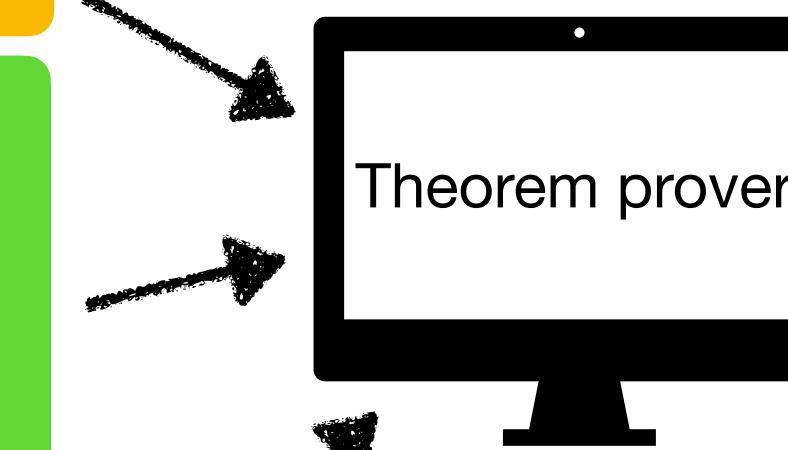
$$X_1 \perp\!\!\!\perp X_3 | X_4$$

$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

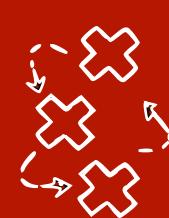
...

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



Provably separating  
 $Y \perp\!\!\!\perp C_1 | S$

Learn  $\hat{f}(S)$   
on source domains



# 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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

All data (including target)

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	?

Query  $Y \perp\!\!\!\perp C_1 | S$ ?

Assumptions

All testable conditional independences from data

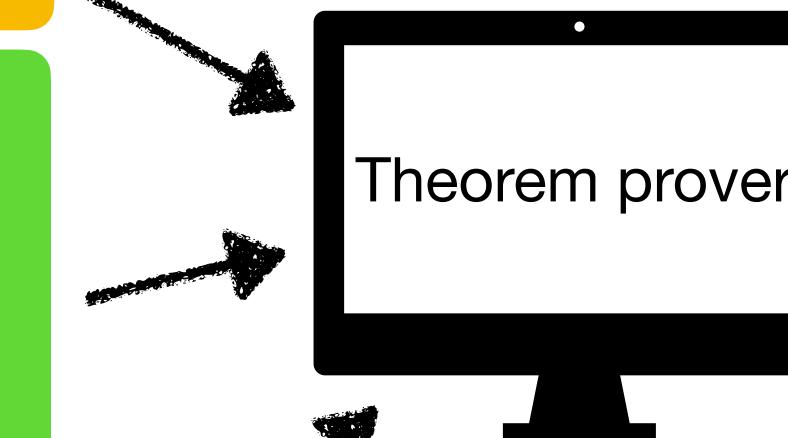
$$X_1 \perp\!\!\!\perp X_3 | X_4$$

$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$

$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$

...

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

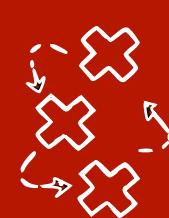


Iterate until empty

Provably not separating  
 $Y \perp\!\!\!\perp C_1 | S$   
Not identifiable

Provably separating  
 $Y \perp\!\!\!\perp C_1 | S$

Learn  $\hat{f}(S)$   
on source domains



# 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

$$L = (\{X_1, C_2\}, \{X_1, X_2, C_2\}, \{X_1, X_2\}, \dots)$$

Select new set S

**Bounded generalisation error**

All data (including target)

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	?

Query  $Y \perp\!\!\!\perp C_1 | S$ ?

Assumptions

All testable conditional independences from data

$$X_1 \perp\!\!\!\perp X_3 | X_4$$

$$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$$
  
$$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$$
  
...

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



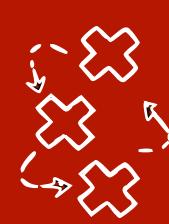
C1	C2	X1	X2	Y
0	1	0,2	0	?
0	1	0,3	0	?
0	1	0,3	1	?

Provably not separating  
 $Y \perp\!\!\!\perp C_1 | S$

Not identifiable

Provably separating  
 $Y \perp\!\!\!\perp C_1 | S$

Learn  $\hat{f}(S)$  on source domains



# Inferring separating sets of features

Query  $Y \perp\!\!\!\perp C_1 |$

Assumptions

All testable conditional independences from data

$X_1 \perp\!\!\!\perp X_3 | X_4$

$Y \perp\!\!\!\perp C_2 | X_1, C_1 = 0$

$X_2 \perp\!\!\!\perp C_2 | Y, C_1 = 0$   
...

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

No need to find causal graph or equivalence class, we only care about conditional independences/d-separations

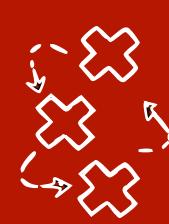
A big (current) limitation:  
Scalability

Provably separating  
 $Y \perp\!\!\!\perp C_1 | X_1$

$Y \perp\!\!\!\perp C_1 | X_1$

Not identifiable

We can test incrementally each set of features selected by standard feature selection (e.g. random forests)



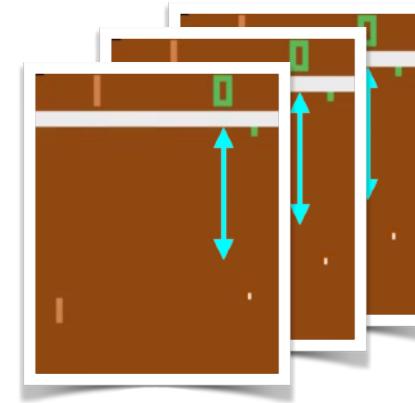
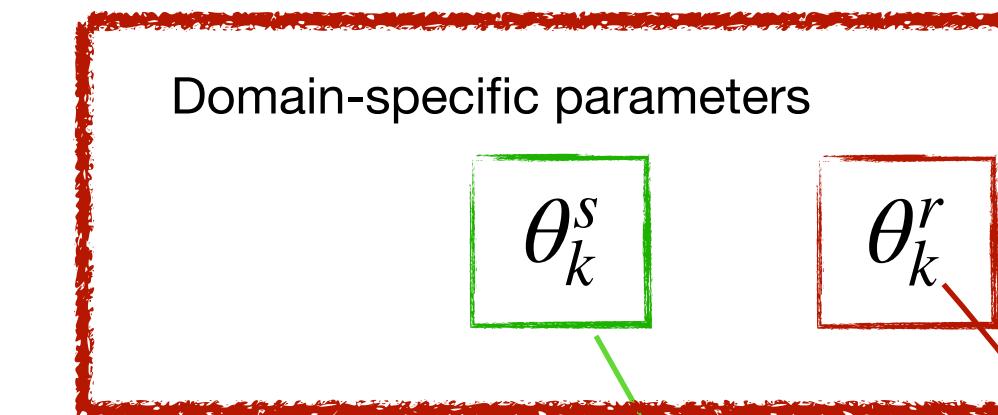
# AdaRL: What, Where, and How to Adapt in Transfer RL

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

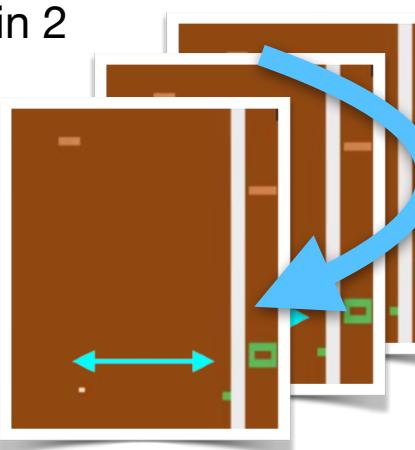
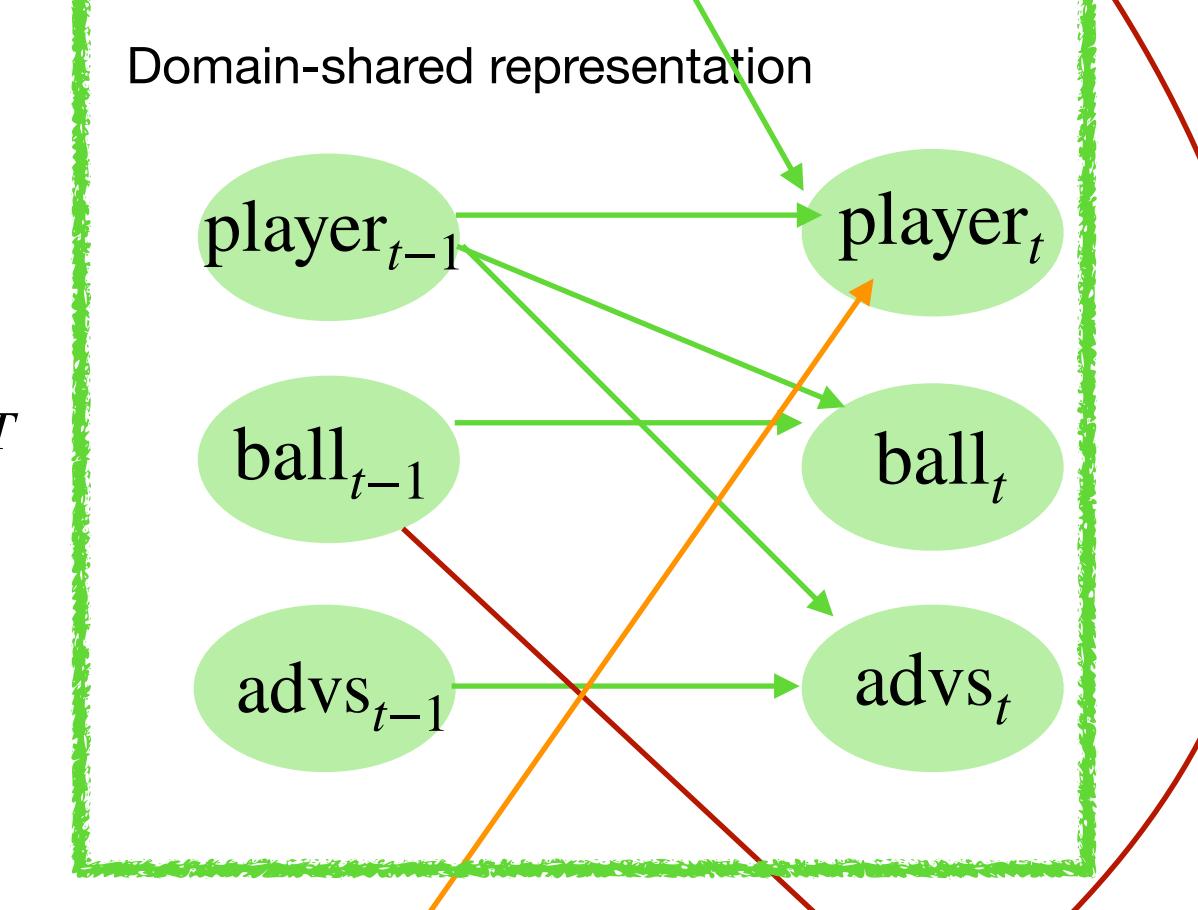
ICLR 2022

Source domains

Domain 1

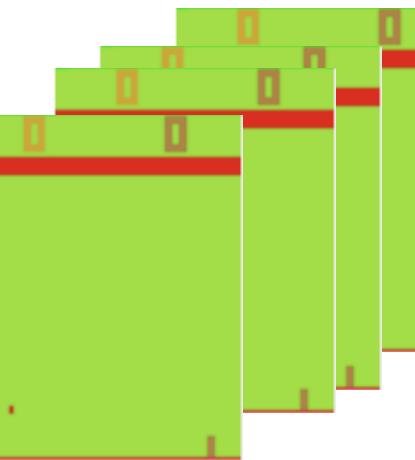
 $\{\text{player}_t, \text{ball}_t, \text{advs}_t, a_t, r_t\}_{t=0,\dots,T}$ 

Domain 2

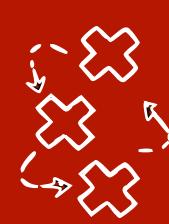
 $\{\text{player}_t, \text{ball}_t, \text{advs}_t, a_t, r_t\}_{t=0,\dots,T}$ 

**When we learn from symbolic inputs, the causal graph can be identified, but we don't have guarantees on what the latent change factors are**

Domain n

 $\{\text{player}_t, \text{ball}_t, \text{advs}_t, a_t, r_t\}_{t=0,\dots,T}$ 

timeslice t-1      timeslice t



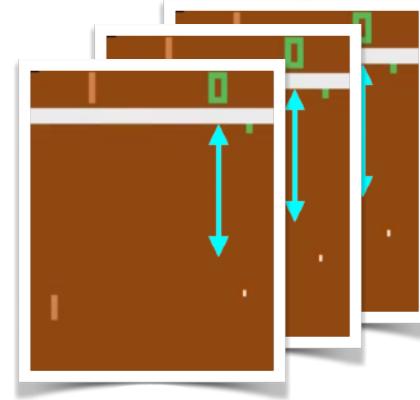
# AdaRL: What, Where, and How to Adapt in Transfer RL

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

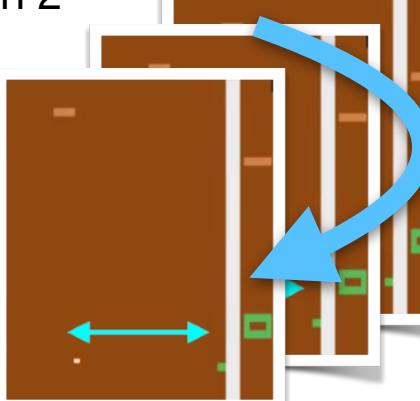
ICLR 2022

Source domains

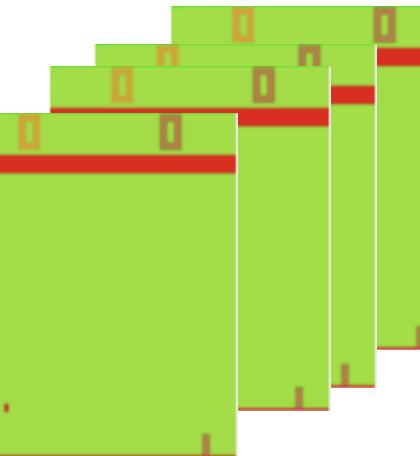
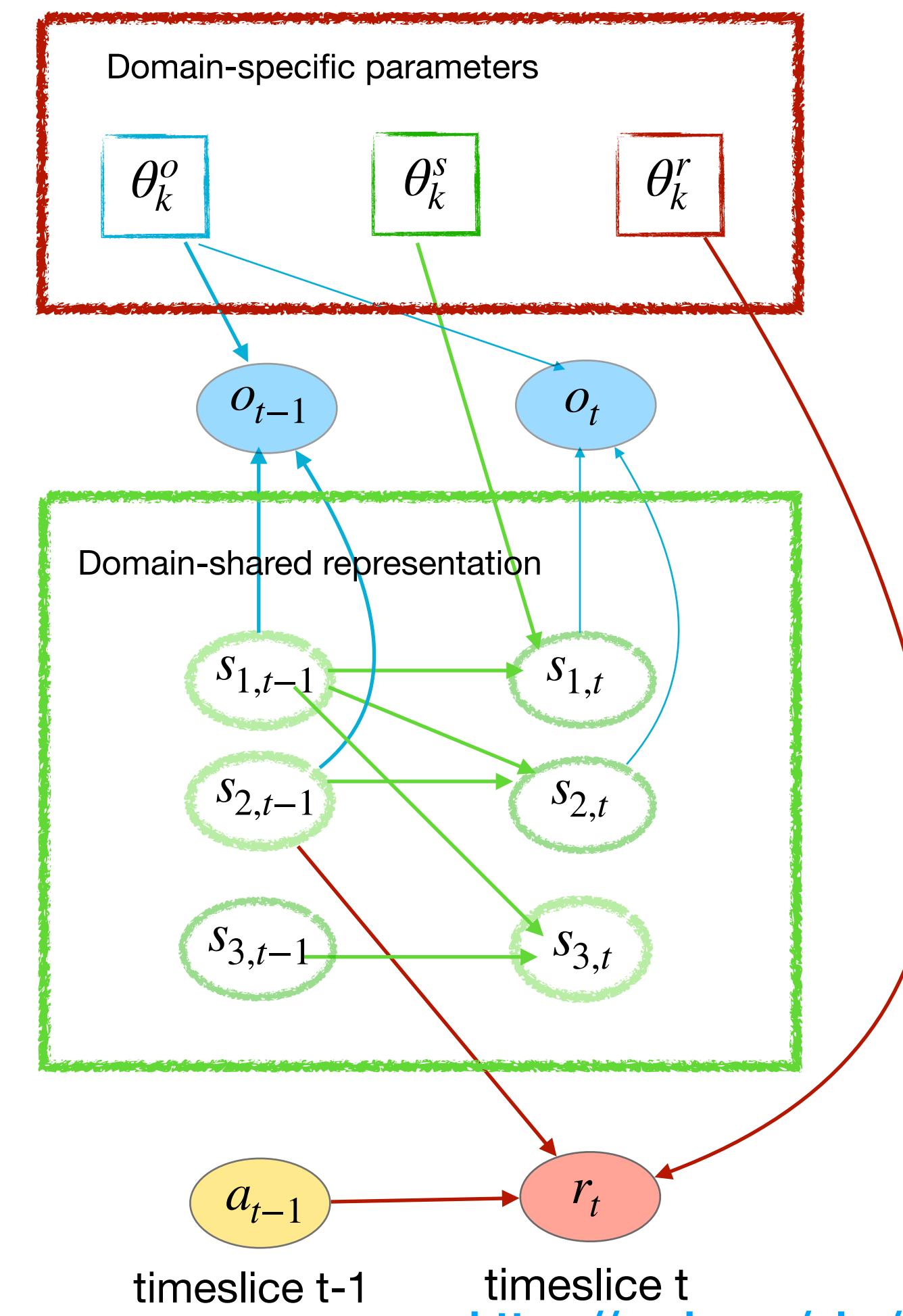
Domain 1

 $\{o_t, a_t, r_t\}_{t=0, \dots, T}$ 

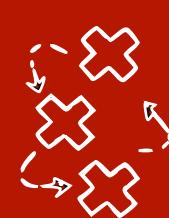
Domain 2

 $\{o_t, a_t, r_t\}_{t=0, \dots, T}$   
...

Domain n

 $\{o_t, a_t, r_t\}_{t=0, \dots, T}$ 

**When we learn from images, we cannot identify the causal variables, so what we learn is not necessarily causal... but it is still useful**



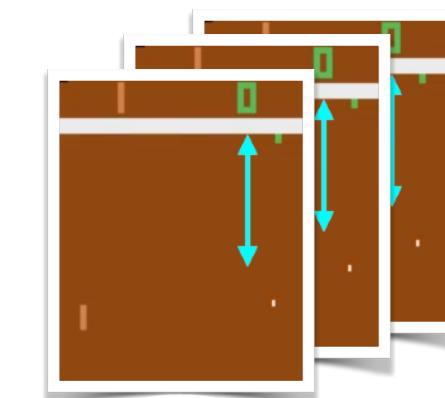
# AdaRL: What, Where, and How to Adapt in Transfer RL

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

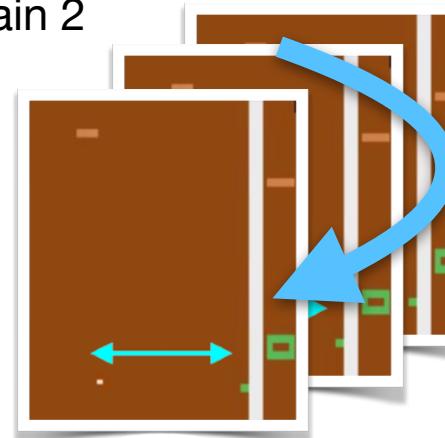
ICLR 2022

Source domains

Domain 1



Domain 2

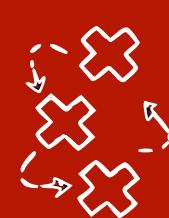


...

Estimate graph over  
estimated  $s_k, \theta_k$

Identify  $s_t^{min}, \theta_t^{min}$   
from the estimated  
graph

Learn optimal  
policy  $\pi^*(s_k^{min}, \theta_k^{min})$   
on source domains



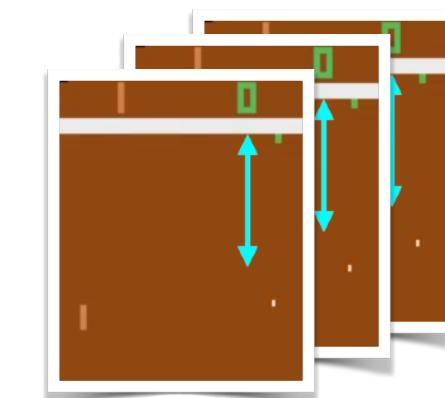
# AdaRL: What, Where, and How to Adapt in Transfer RL

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

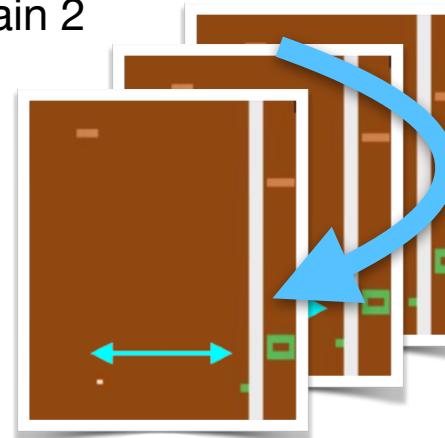
ICLR 2022

Source domains

Domain 1



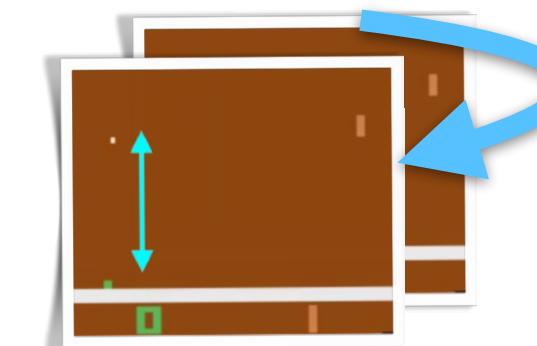
Domain 2

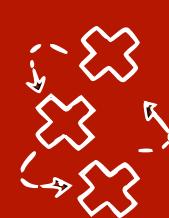


...

Estimate graph over  
estimated  $s_k, \theta_k$ Identify  $s_t^{min}, \theta_t^{min}$   
from the estimated  
graphLearn optimal  
policy  $\pi^*(s_k^{min}, \theta_k^{min})$   
on source domains

Target domain

 $\{o_t, a_t, r_t\}_{t=0, \dots, T}$ Use model to  
estimate  $s_{target}^{min}, \theta_{target}^{min}$   
with few samplesApply policy  
 $\pi^*(s_{target}^{min}, \theta_{target}^{min})$ **Simplifying  
assumption: no  
new edges in  
target domain**

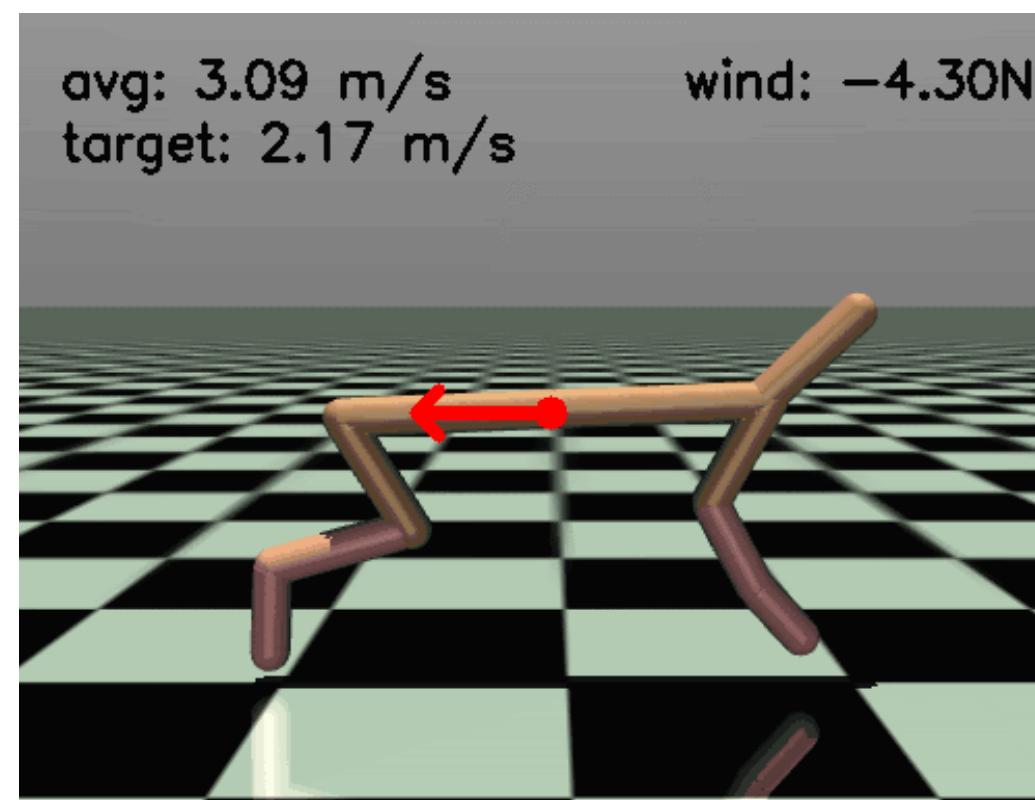


# FansRL: Factored Adaptation for Non-Stationary Reinforcement Learning

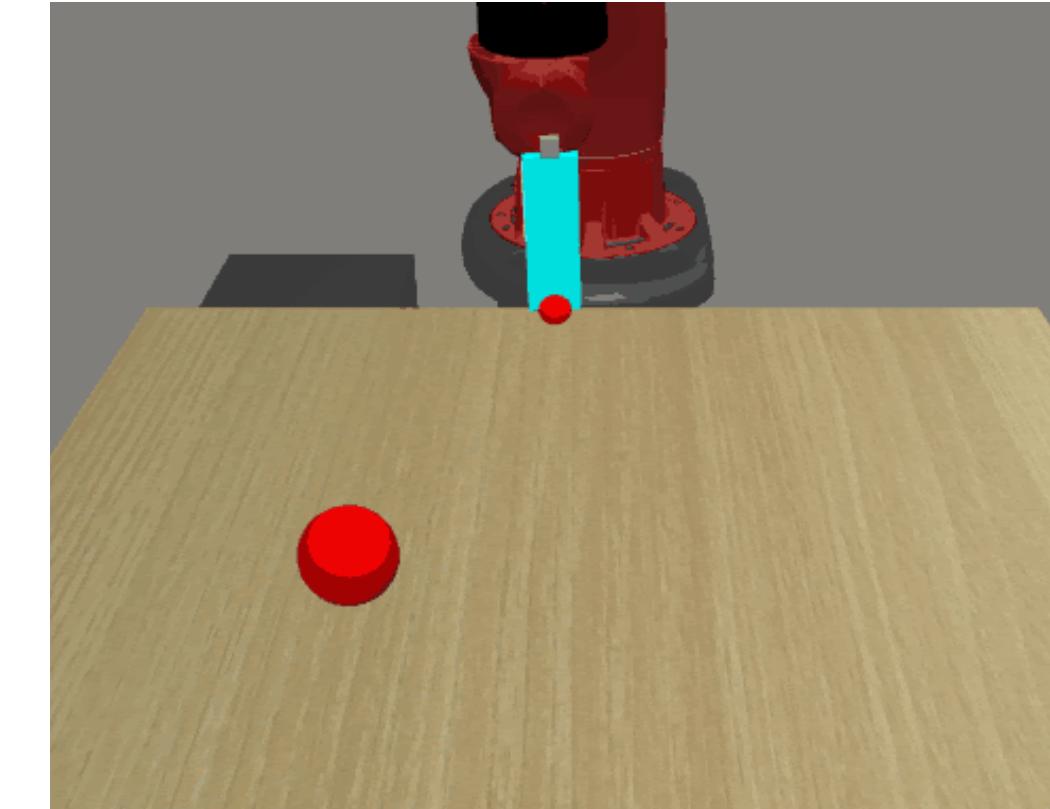
Fan Feng, Biwei Huang, Kun Zhang, Sara Magliacane

NeurIPS 2022

- **Task:** RL agent has to learn a policy that is robust to different types of non-stationarity, including **multiple simultaneous changes of different types**



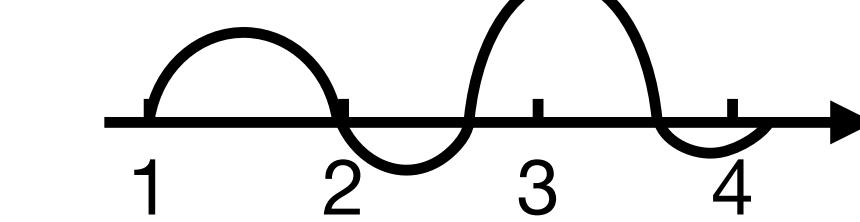
Non-stationary environments  
(wind changes)



Non-stationary rewards  
(target changes)

Continuous

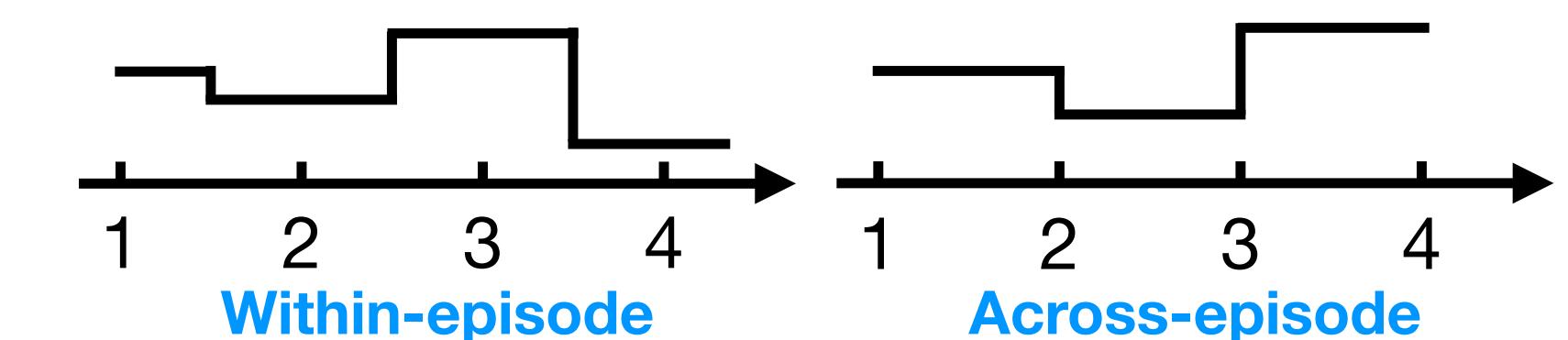
episode



Different functions, e.g. sine, linear, damping

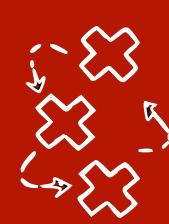
Discrete

episode



Within-episode

Across-episode



# Conclusions

- Graphical models and d-separation [Pearl 1988] are a principled way to reason about **invariances and distribution shift**
- Not a new observation, known since [Schoelkopf et al 2012]
- Even with **unknown causal graphs, Missing data/zero-shot settings**
- Often we **do not need to reconstruct the causal graph**, we only need to infer missing conditional independences
- These ideas seem empirically useful even if we **cannot guarantee** that we are **learning the true causal variables or the true causal graph**