

IMPERIAL

Causal Discovery for Trustworthy AI

Fabrizio Russo

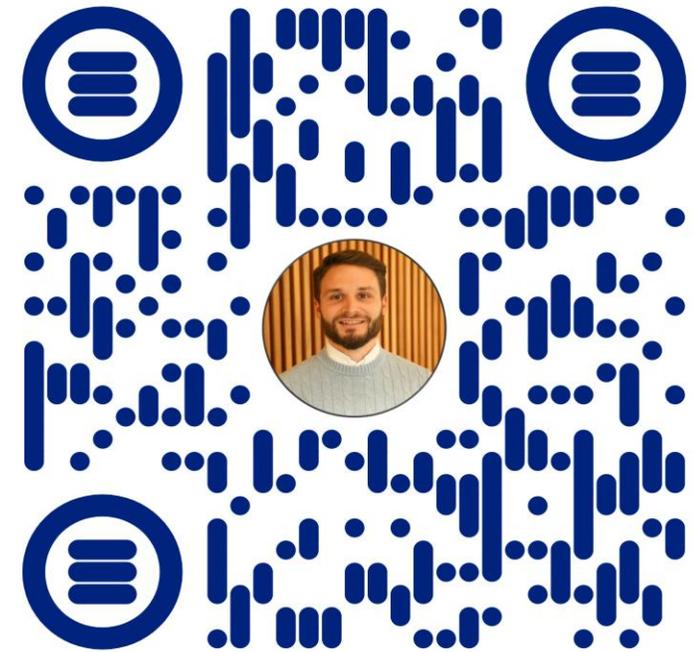
Ethics, Fairness and Explanation in AI

19/02/2026 Guest Lecture

My Journey

Across Finance And Artificial Intelligence

- 2009-2012 Rome – Tor Vergata: BSc Economics & Finance
Madrid – Autónoma
- 2013-2014 London – Learning & Exploring
- 2014-2015 London – General Electric Capital: Credit Risk Analyst
- 2015-2020 London – 4most Europe: Consultant – Head of Data Science
- 2016-2018 London – Birkbeck College: MSc Applied Statistics
- 2020-2024 London – Imperial College: PhD Safe and Trusted AI
- 2025-2026 London – Imperial College: Research Associate
- 2026 June London – Imperial College: Research Fellow



Agenda

- 1. Motivation & Background**
 - i. Causal Models
 - ii. Causal Discovery
- 2. Causal Graphs for Contestable Neural Networks**
- 3. Argumentative Causal Discovery**
 - i. LLM as Imperfect Experts
- 4. Recap & Conclusion**

Structural Causal Model (Pearl, 2009)

Causal Graphs + Structural Equation

Learned edge function between A and B

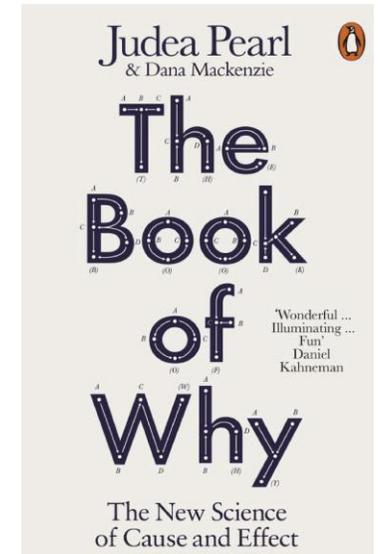
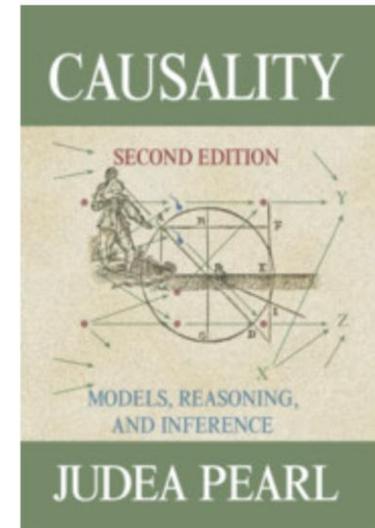
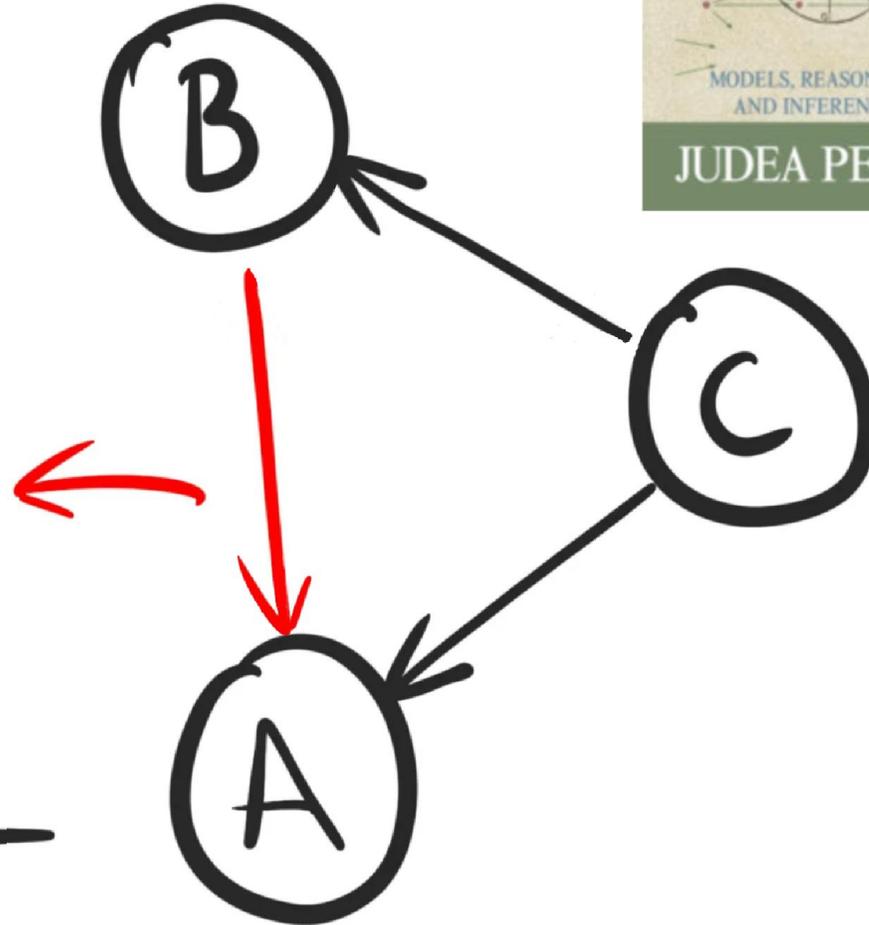
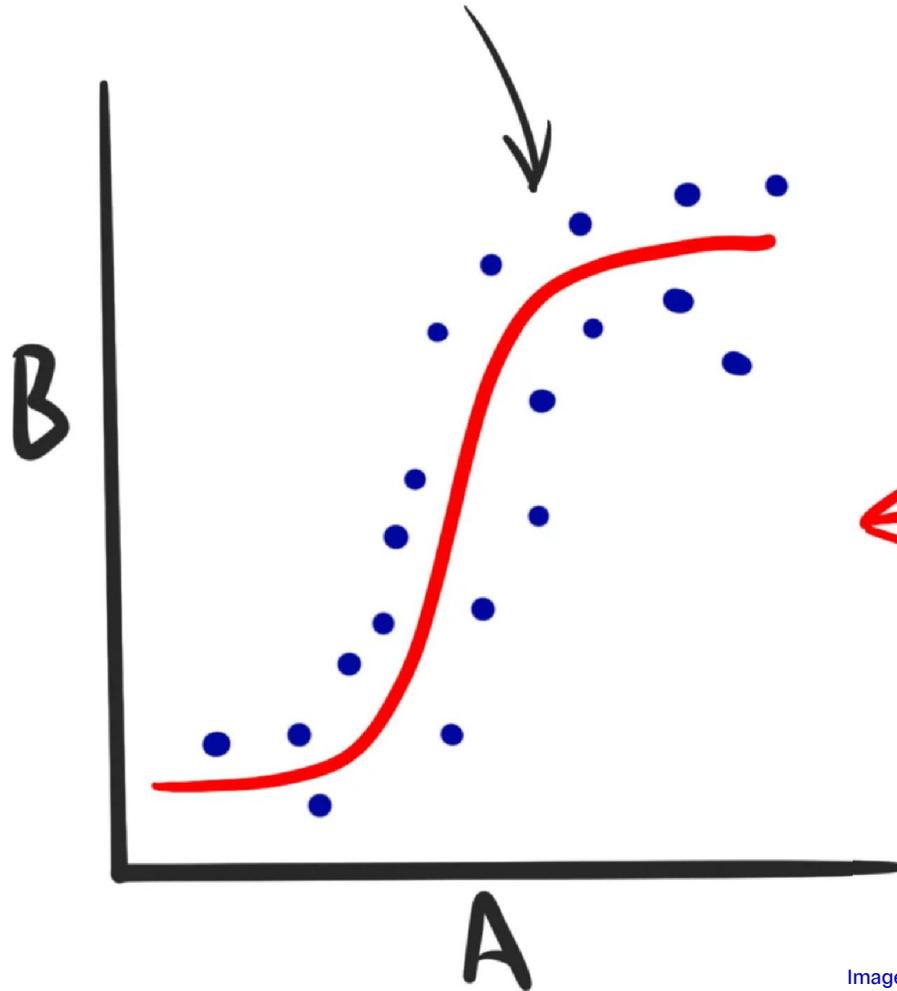


Image from <https://towardsdatascience.com/how-to-understand-the-world-of-causality-c698cdc9f27c>

Prediction vs Intervention

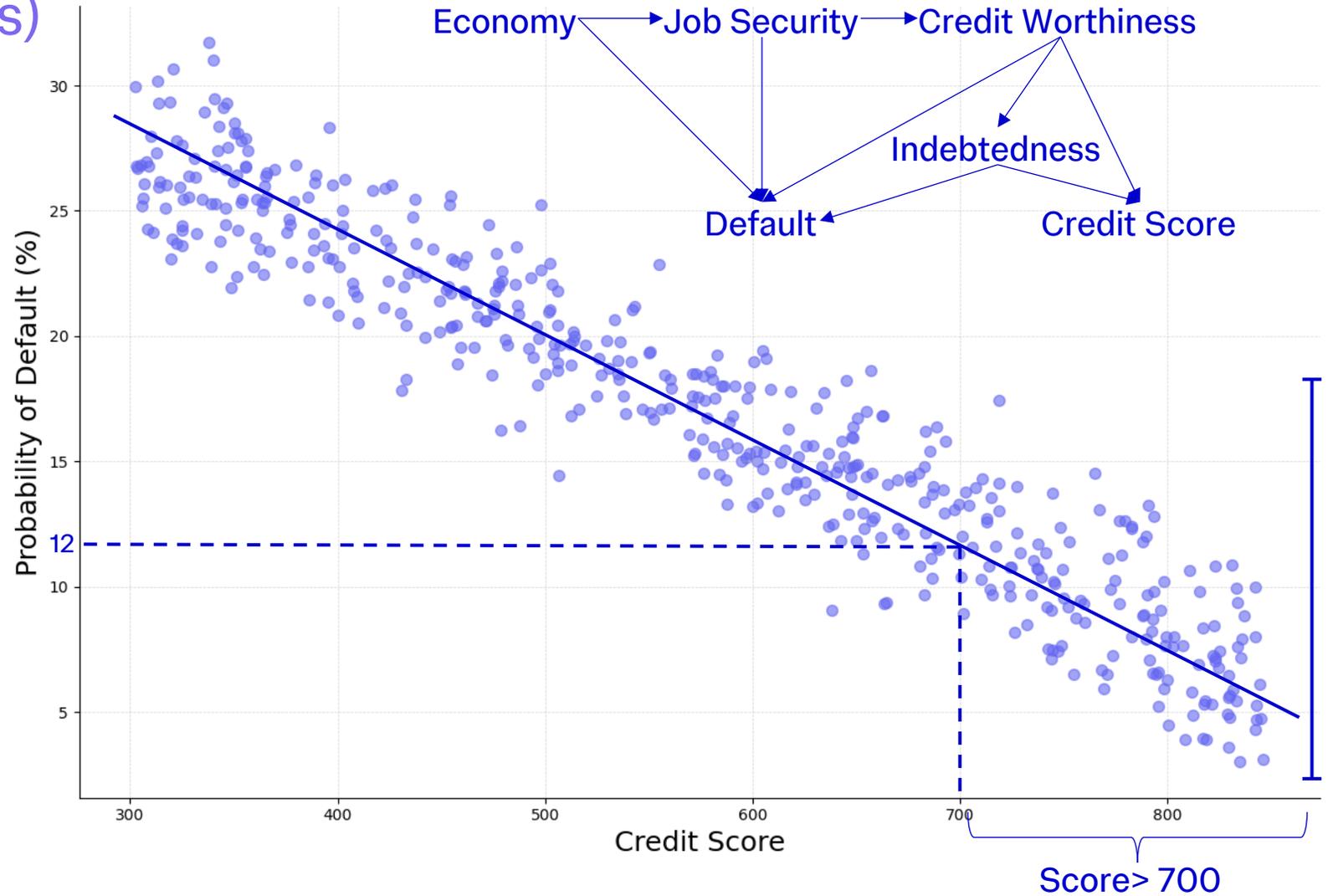
Causal Assumptions (graphs)

Predictive Models

- What if we *observe* a credit score of 700?
- What if we *observe* a PD of 12%?

Causal Models

- What if we accept only customers with 700+ score?
- What if we *do* (perform an action/intervention in an environment)?



Prediction vs Intervention

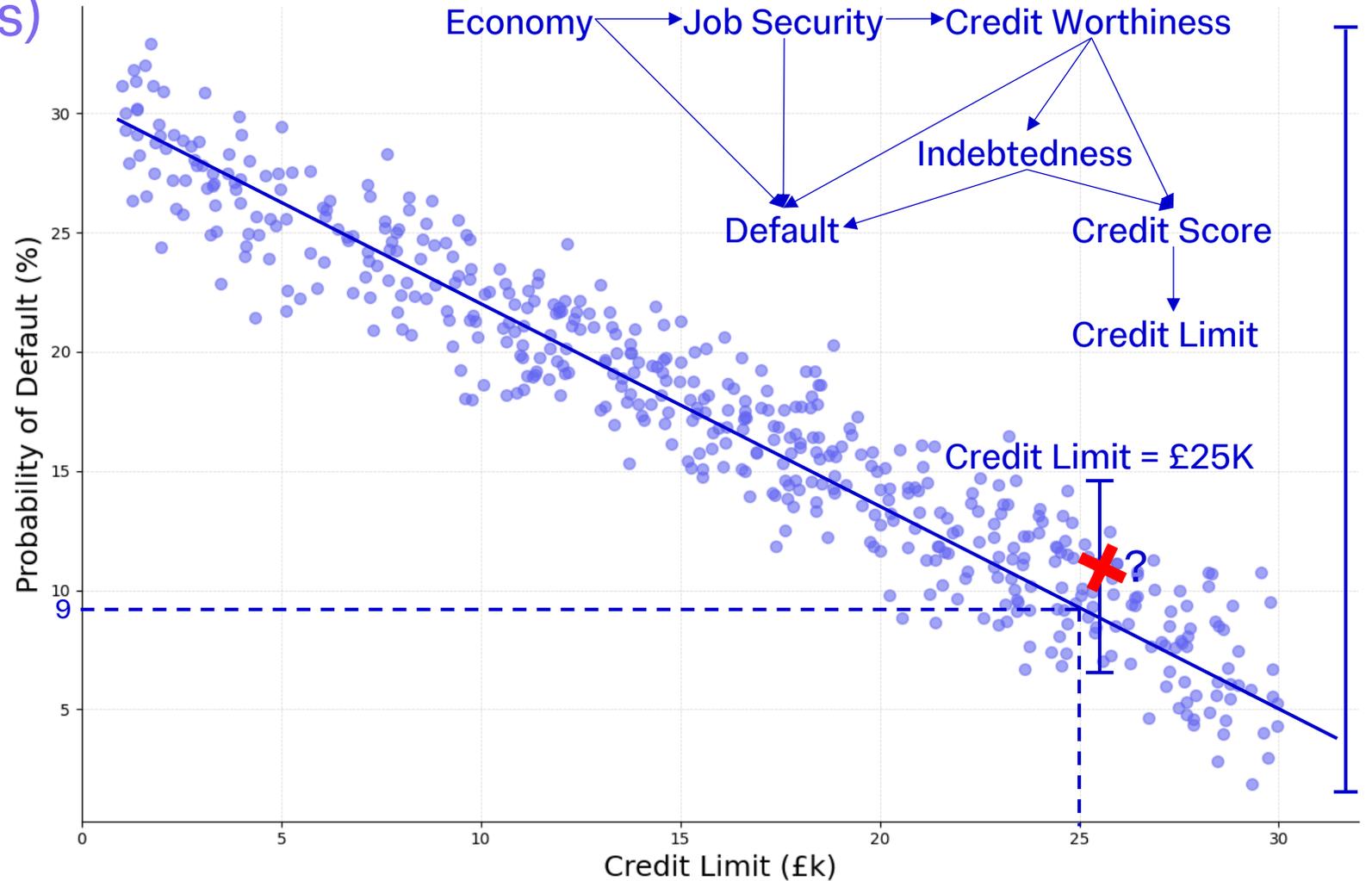
Causal Assumptions (graphs)

Predictive Models

- What if we observe a credit limit of £25k?
- What if we observe a PD of 9%?

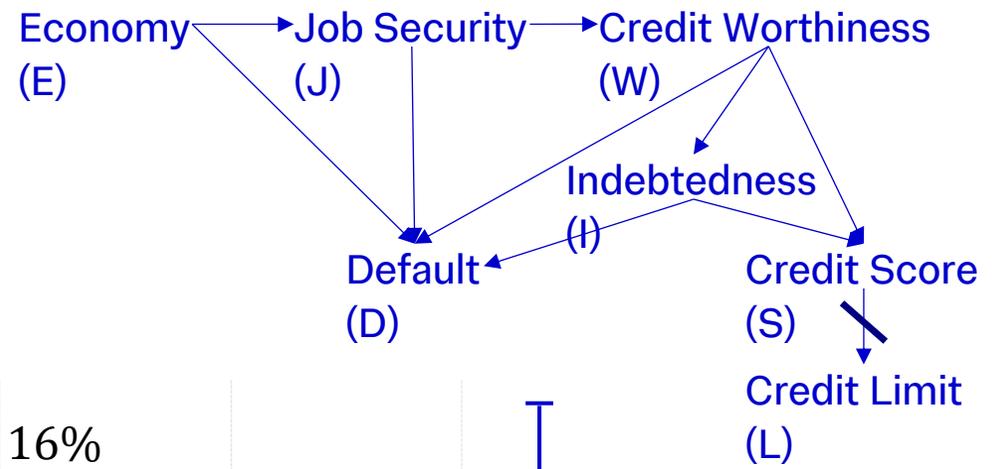
Causal Models

- What if we *intervene* on Credit Limit?
- What's the *causal* effect on the PD?



Causal Models

Structural Equations



Exogenous
 $E := f_E(U_E)$

Endogenous:

$J := f_J(E, U_J)$

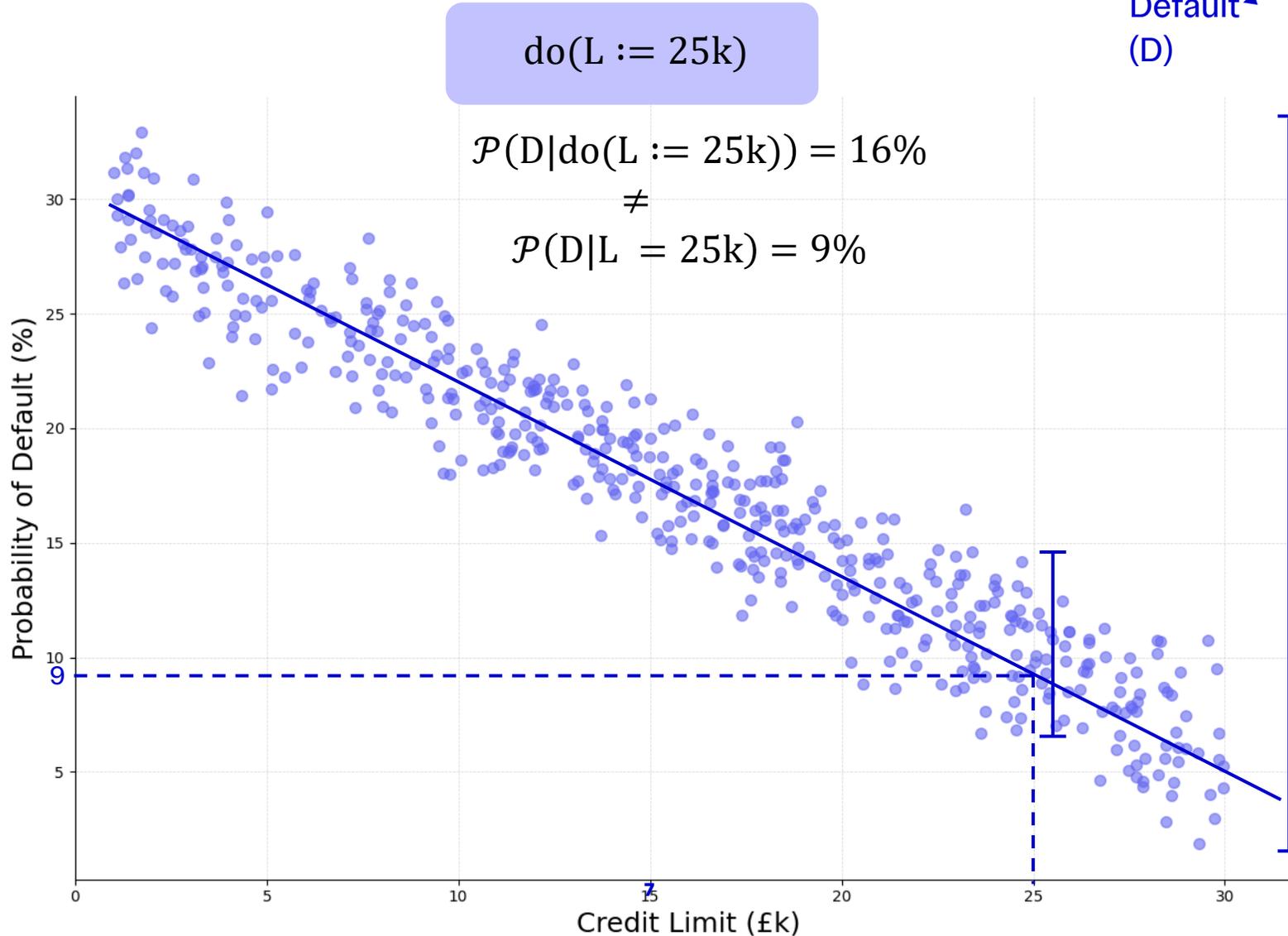
$W := f_W(J, U_W)$

$I := f_I(W, U_I)$

$D := f_D(E, J, W, U_D)$

$S := f_S(W, I, U_S)$

$L := 25k$



Machine Learning vs Causal Models

High-Level Comparison

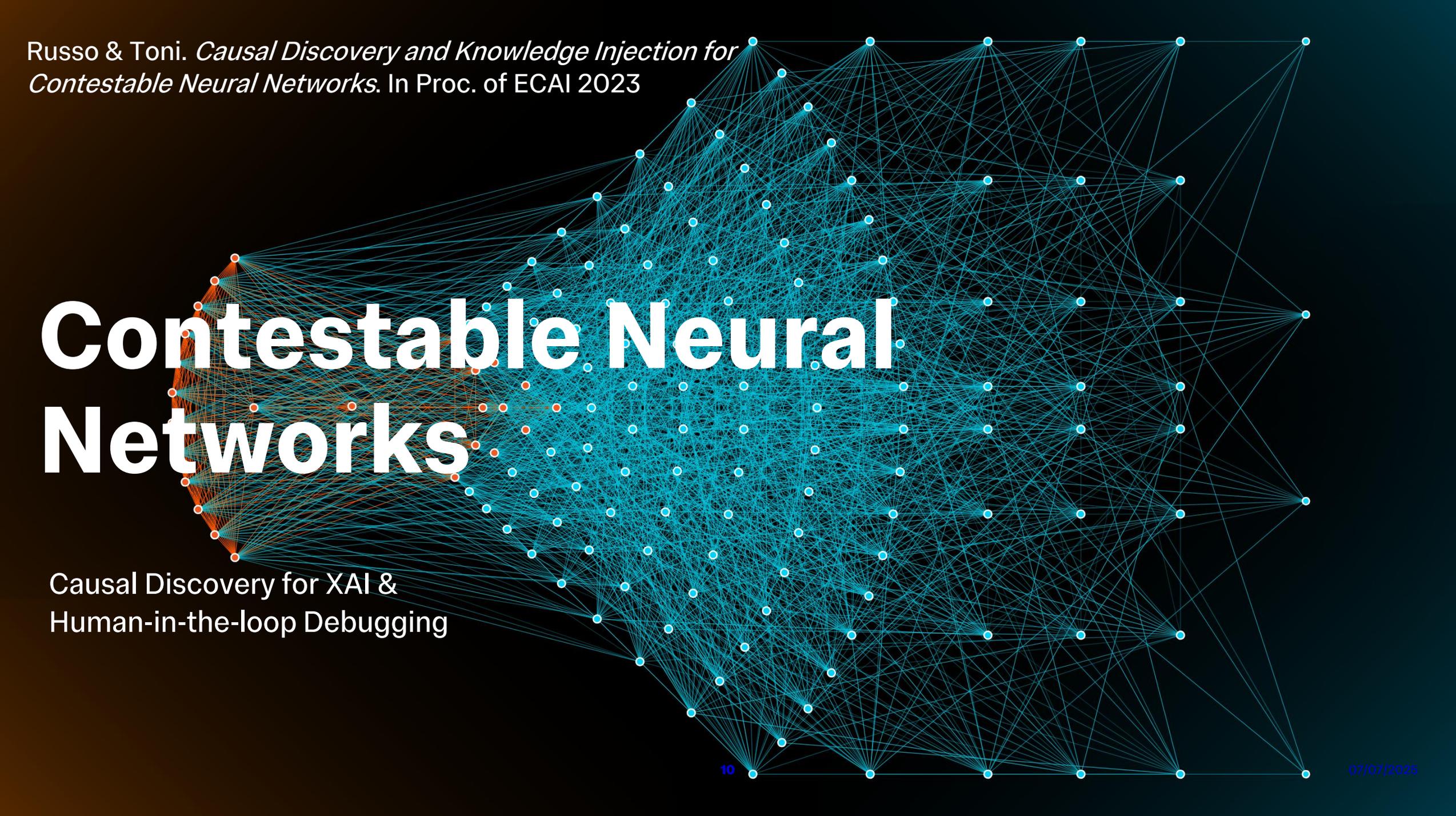
Attribute	$\mathcal{P}(D L = 25k)$ Machine Learning	$\mathcal{P}(D \text{do}(L := 25k))$ Causal Models
Interpretability	Limited	High
Predictive Accuracy	High	Moderate
Generalisation	Moderate	High
Actionable Insights	Limited	High
Data Requirements	High	Moderate
Scalability	High	Low
Expert Knowledge	Moderate	High



Image from <https://towardsdatascience.com/how-to-understand-the-world-of-causality-c09edc0f27c>

Russo & Toni. *Causal Discovery and Knowledge Injection for Contestable Neural Networks*. In Proc. of ECAI 2023

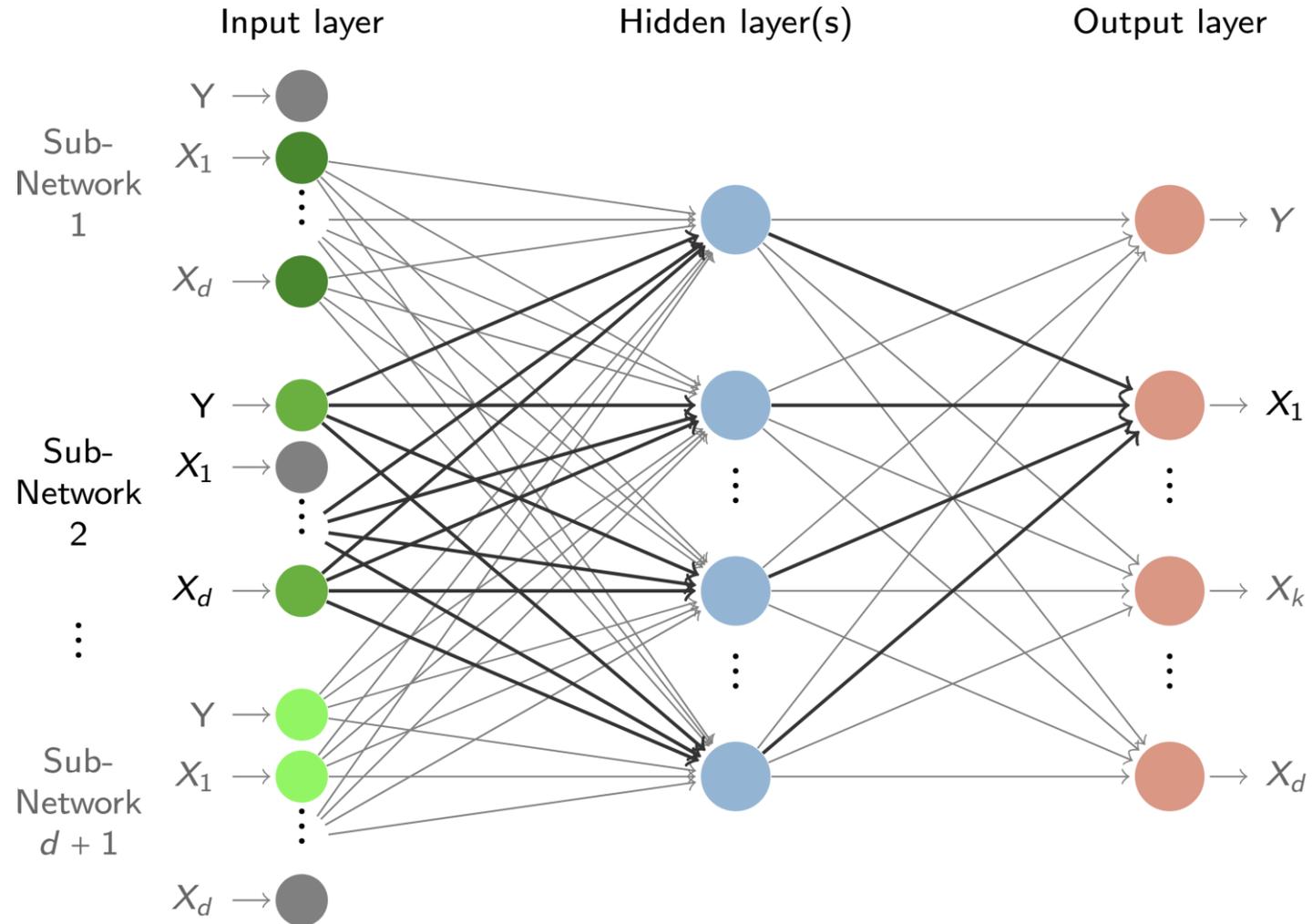
Contestable Neural Networks



Causal Discovery for XAI &
Human-in-the-loop Debugging

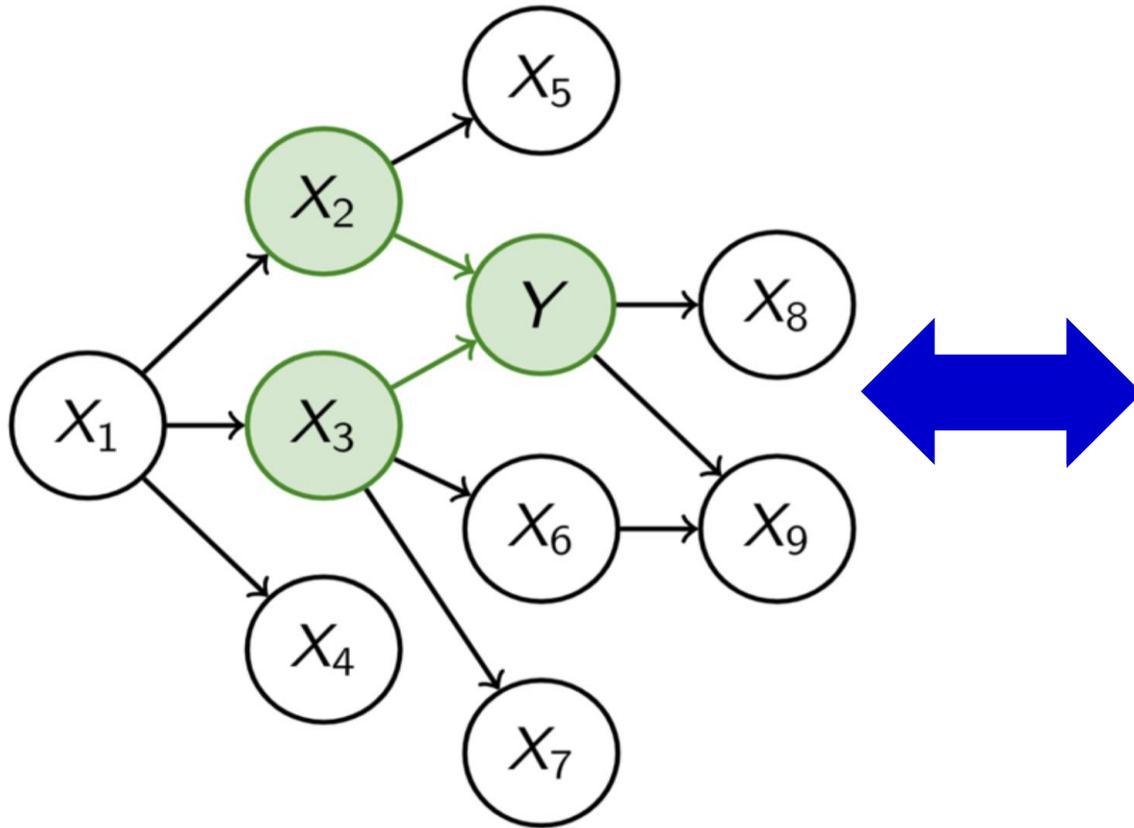
Joint Neural Network Structure

(Kyono, Zhang and van der Schaar 2020)



Objective

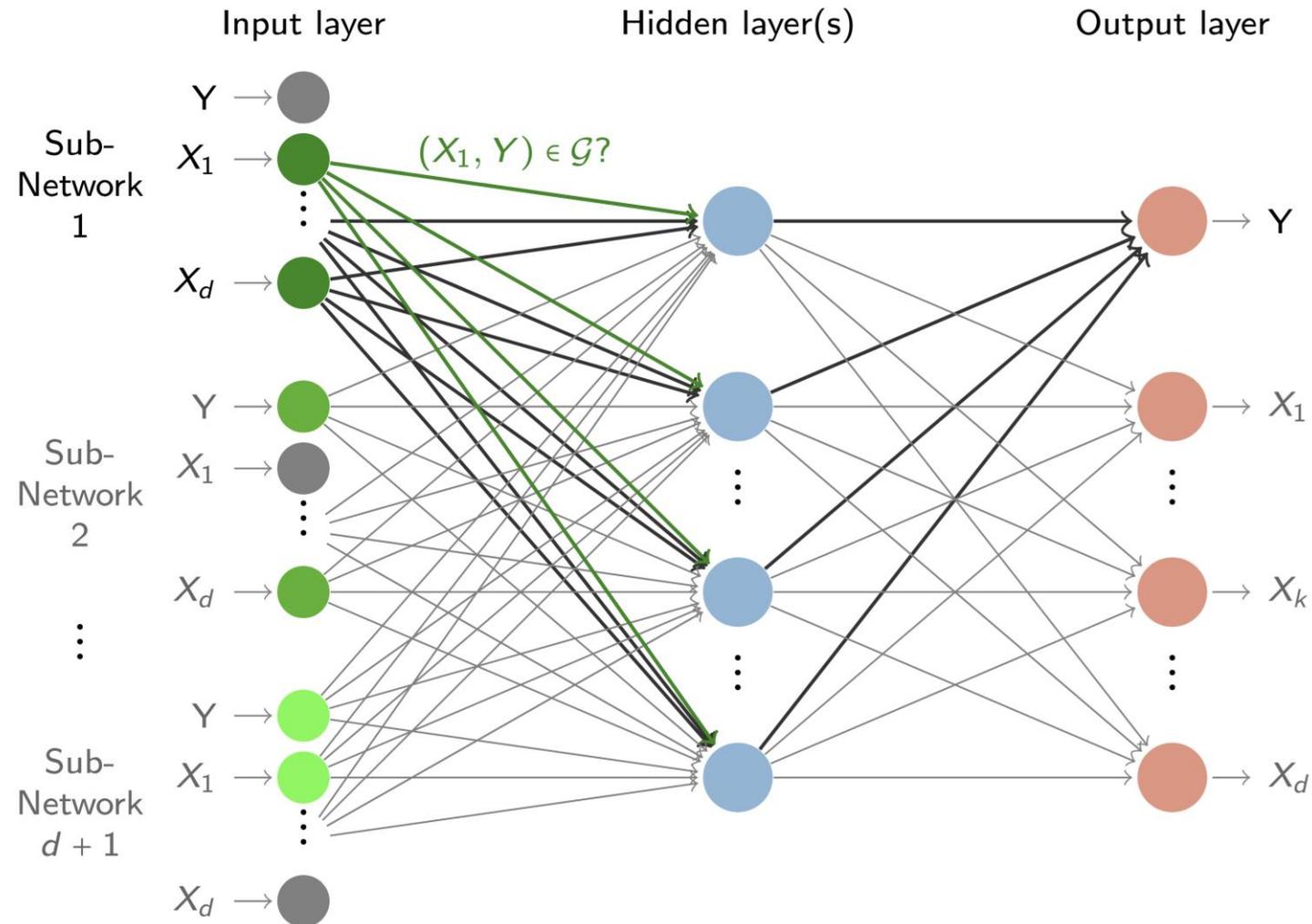
Capture Causal Relations



	Y	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9
Y	0.0	0.005	0.017	0.008	0.002	0.042	0.02	0.005	0.059	0.05
X_1	0.006	0.0	0.063	0.054	0.068	0.009	0.006	0.013	0.006	0.008
X_2	0.088	0.036	0.0	0.022	0.019	0.124	0.008	0.011	0.006	0.008
X_3	0.087	0.034	0.021	0.0	0.024	0.005	0.107	0.104	0.006	0.009
X_4	0.009	0.032	0.02	0.023	0.0	0.01	0.013	0.01	0.005	0.005
X_5	0.026	0.006	0.017	0.004	0.004	0.0	0.012	0.002	0.005	0.018
X_6	0.025	0.006	0.008	0.011	0.005	0.017	0.0	0.014	0.002	0.114
X_7	0.029	0.003	0.007	0.011	0.002	0.024	0.029	0.0	0.011	0.01
X_8	0.036	0.002	0.004	0.003	0.004	0.006	0.009	0.006	0.0	0.006
X_9	0.024	0.003	0.003	0.004	0.003	0.005	0.079	0.01	0.004	0.0

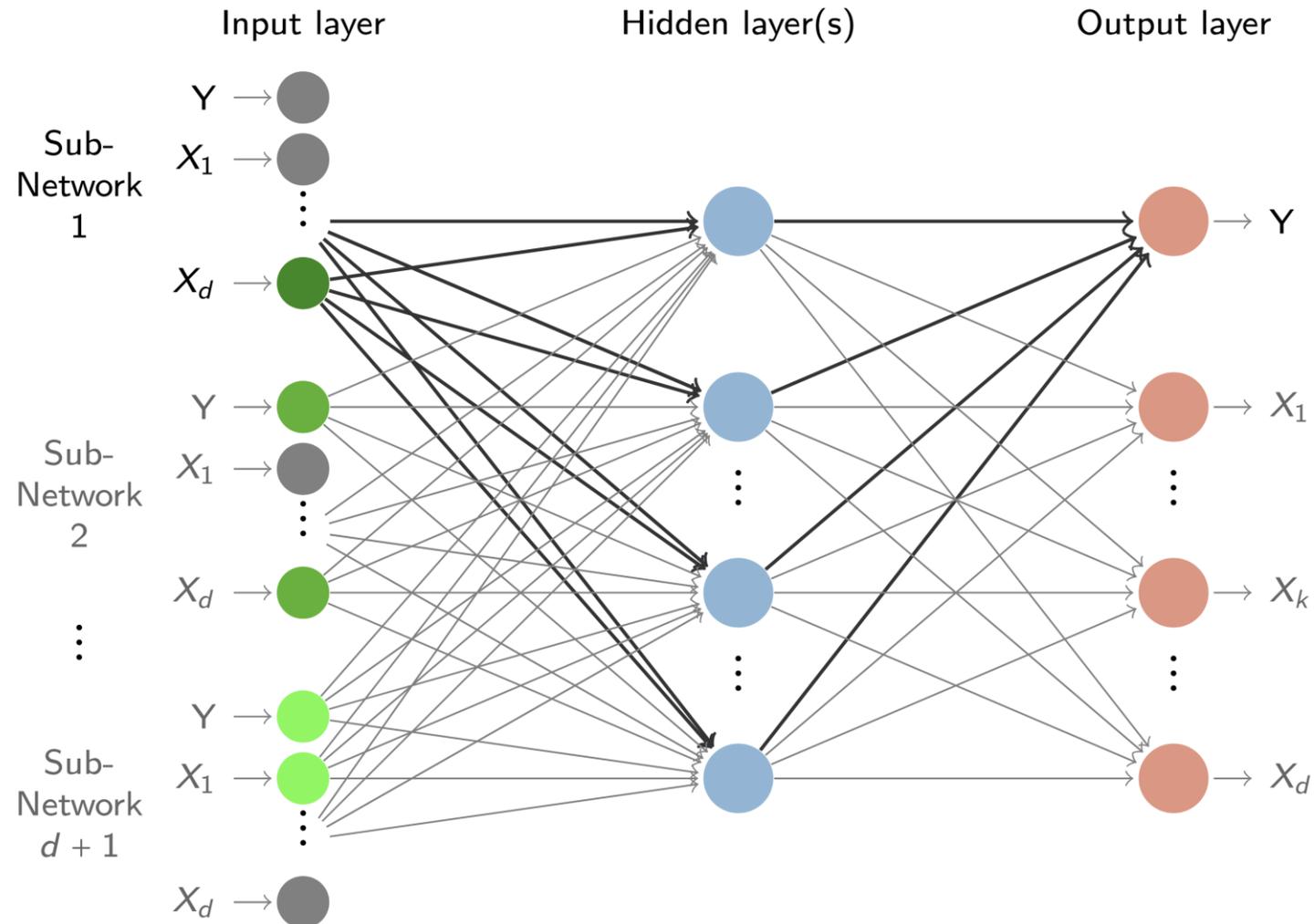
Encode Causality in the Network Structure

(Kyono, Zhang and van der Schaar 2020)



Encode Causality in the Network Structure

(Kyono, Zhang and van der Schaar 2020)



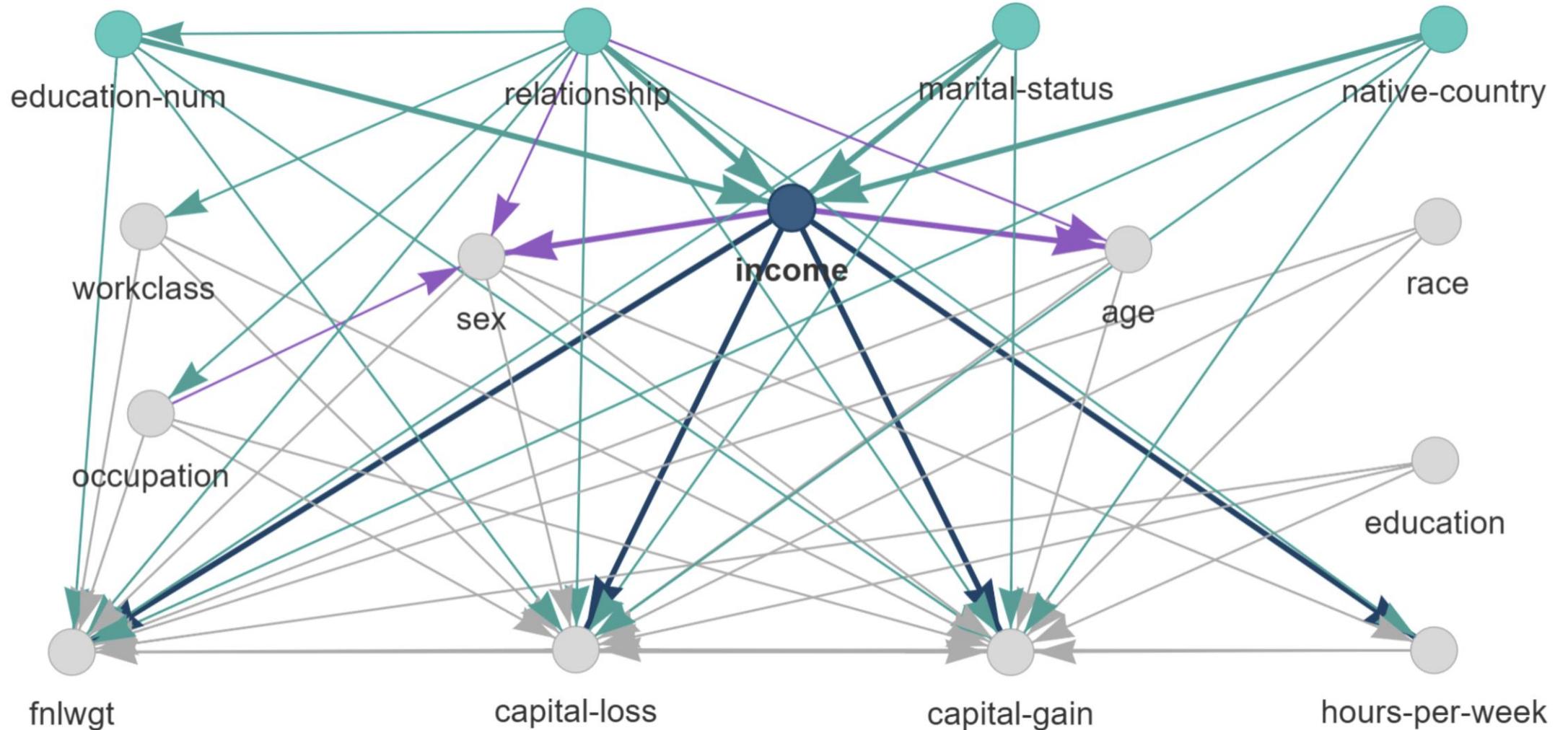
Income Prediction Case Study

Adult Dataset (Becker and Kohavi, 1996)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)	(15)	
(1)			0.0	0.2				0.1					36.1	3.7	21.3	income
(2)													1.6	0.3	11.5	race
(3)								0.1					1.9	0.4	16.3	sex
(4)													0.8	0.1	2.0	age
(5)	0.0												1.1	0.2	4.9	native-country
(6)			0.0										0.6	0.2	4.0	occupation
(7)													1.2	0.3	6.4	workclass
(8)													0.8	0.1	3.6	hours-per-week
(9)													0.5	0.2	5.5	education
(10)	0.0												2.7	0.4	6.9	education-num
(11)	0.1												2.0	0.3	18.1	marital-status
(12)	0.0		0.0	0.1		0.0	0.0	0.1		0.0			2.6	0.5	15.0	relationship
(13)															0.1	capital-gain
(14)													0.2		0.5	capital-loss
(15)																fnlwgt

Income Prediction Case Study

Adult Dataset (Becker and Kohavi, 1996)



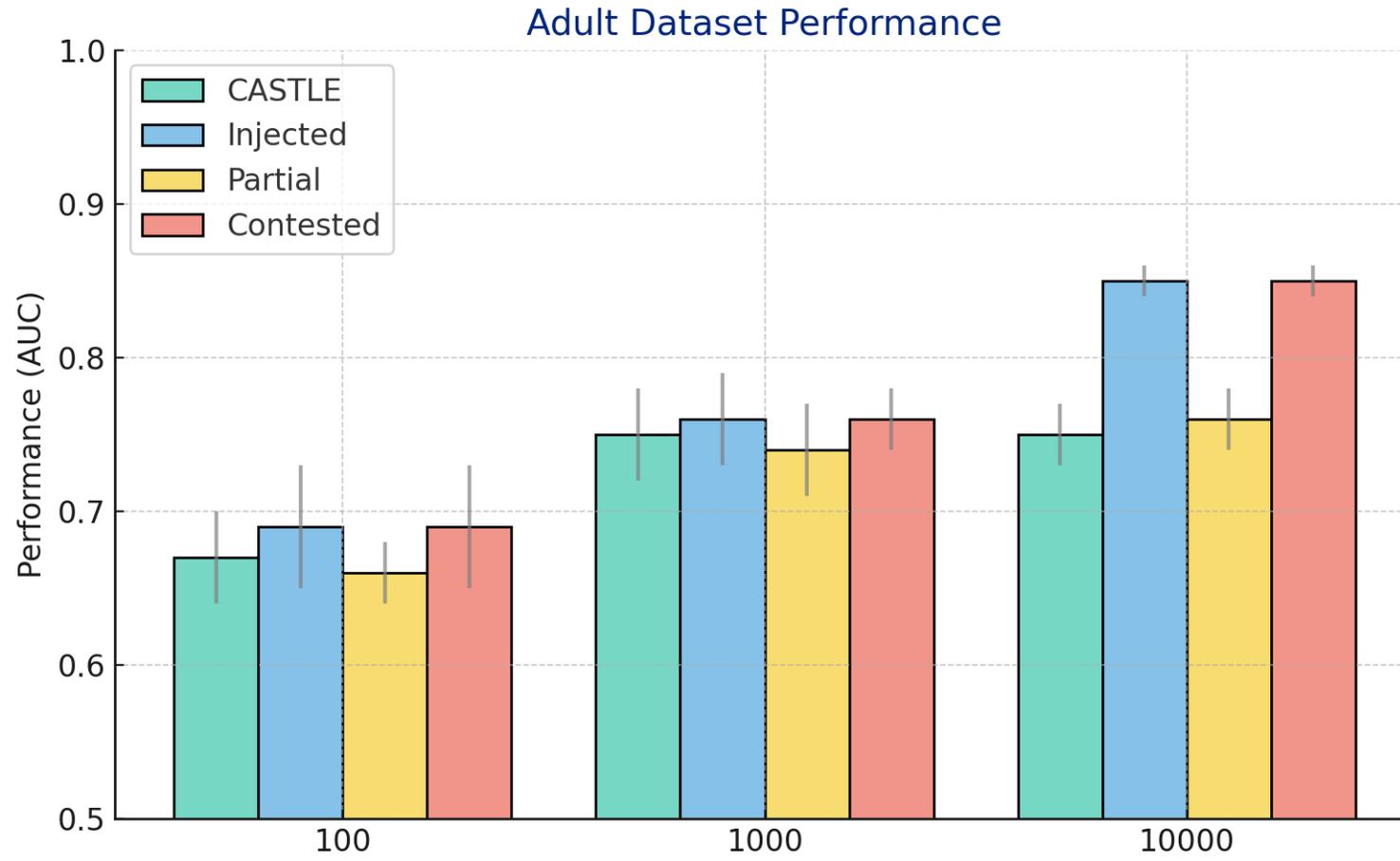
Income Prediction Case Study

Adult Dataset (Becker and Kohavi, 1996)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)	(15)	
(1)																income
(2)	■					■	■	■	■	■	■	■	■	■	■	race
(3)	■					■	■	■	■	■	■	■	■	■	■	sex
(4)	■					■	■	■	■	■	■	■	■	■	■	age
(5)	■					■	■	■	■	■	■	■	■	■	■	native-country
(6)	■					□	■	■	□	□	□	□	■	■	□	occupation
(7)	■					■	□	■	■	■	■	■	■	■	■	workclass
(8)	■					■	■	□	□	□	□	□	■	■	□	hours-per-week
(9)	■					■	■	■	□	■	■	■	■	■	■	education
(10)	■					■	■	■	■	□	■	■	■	■	■	education-num
(11)	■					■	■	■	■	■	□	■	■	■	■	marital-status
(12)	■					■	■	■	■	■	■	□	■	■	■	relationship
(13)	■					□	□	□	□	□	□	□	□	□	□	capital-gain
(14)	■					□	□	□	□	□	□	□	□	□	□	capital-loss
(15)	■					■	■	■	■	■	■	■	■	■	□	fnlwgt

Income Prediction Case Study

Adult Dataset (Becker and Kohavi, 1996)



Takeaways

Causal Graphs to...



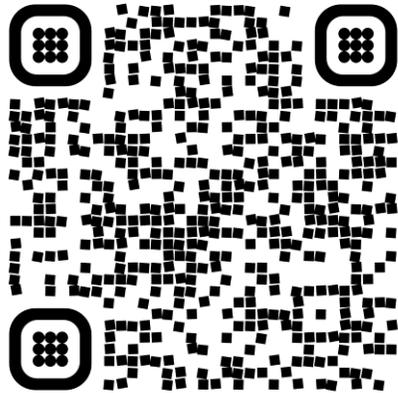
Explain



Contest



Improve

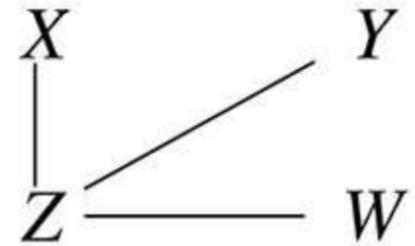
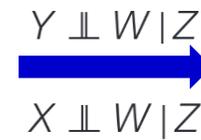
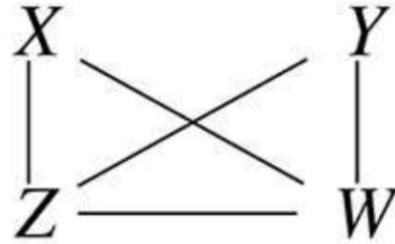
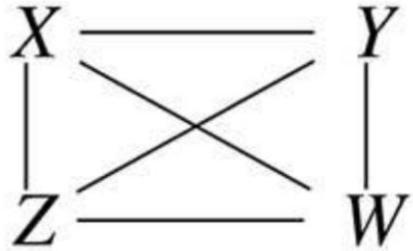
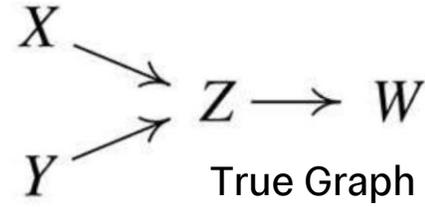


Russo & Toni. *Causal Discovery and Knowledge Injection for Contestable Neural Networks*. In Proc. of ECAI 2023

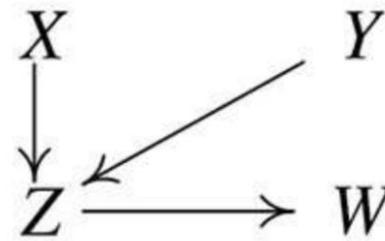
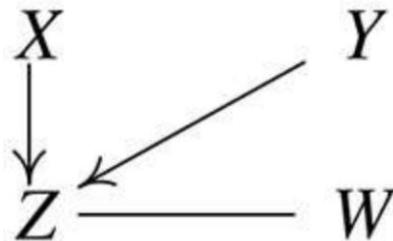
How do we **reliably** build causal graphs from data?

The Peter-Clark Algorithm

(Spirtes et al, 1993) Example from Glymour et al, 2019



$X \perp Y$
 $X \not\perp Y | Z$



Income Prediction Example



$E \perp\!\!\!\perp I$	$p = 0.00$	$\mathcal{S} = 1.00$
$E \perp\!\!\!\perp O$	$p = 0.00$	$\mathcal{S} = 1.00$
$R \perp\!\!\!\perp O$	$p = 0.00$	$\mathcal{S} = 1.00$
$O \perp\!\!\!\perp I$	$p = 0.00$	$\mathcal{S} = 1.00$
$R \perp\!\!\!\perp E$	$p = 0.46$	$\mathcal{S} = 0.71$
$R \perp\!\!\!\perp I$	$p = 0.05$	$\mathcal{S} = 0.52$
$E \perp\!\!\!\perp I \mid \{R\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$E \perp\!\!\!\perp I \mid \{O\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$E \perp\!\!\!\perp O \mid \{R\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$R \perp\!\!\!\perp O \mid \{I\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$E \perp\!\!\!\perp O \mid \{I\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$R \perp\!\!\!\perp O \mid \{E\}$	$p = 0.00$	$\mathcal{S} = 0.50$

$O \perp\!\!\!\perp I \mid \{E\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$O \perp\!\!\!\perp I \mid \{R\}$	$p = 0.00$	$\mathcal{S} = 0.50$
$R \perp\!\!\!\perp E \mid \{O\}$	$p = 0.53$	$\mathcal{S} = 0.38$
$R \perp\!\!\!\perp I \mid \{O\}$	$p = 0.03$	$\mathcal{S} = 0.35$
$R \perp\!\!\!\perp E \mid \{O\}$	$p = 0.33$	$\mathcal{S} = 0.32$
$R \perp\!\!\!\perp E \mid \{I\}$	$p = 0.05$	$\mathcal{S} = 0.25$
$R \perp\!\!\!\perp E \mid \{O, I\}$	$p = 0.39$	$\mathcal{S} = 0.00$
$R \perp\!\!\!\perp I \mid \{E, O\}$	$p = 0.00$	$\mathcal{S} = 0.00$
$E \perp\!\!\!\perp O \mid \{R, I\}$	$p = 0.00$	$\mathcal{S} = 0.00$
$R \perp\!\!\!\perp I \mid \{E, O\}$	$p = 0.00$	$\mathcal{S} = 0.00$
$R \perp\!\!\!\perp O \mid \{E, I\}$	$p = 0.03$	$\mathcal{S} = 0.00$

Income Prediction Example



True Causal Graph



Majority-PC (Colombo and Maathias, 2012)

No Directed Acyclic Graph (DAG) is compatible with *all* the tests at the same time

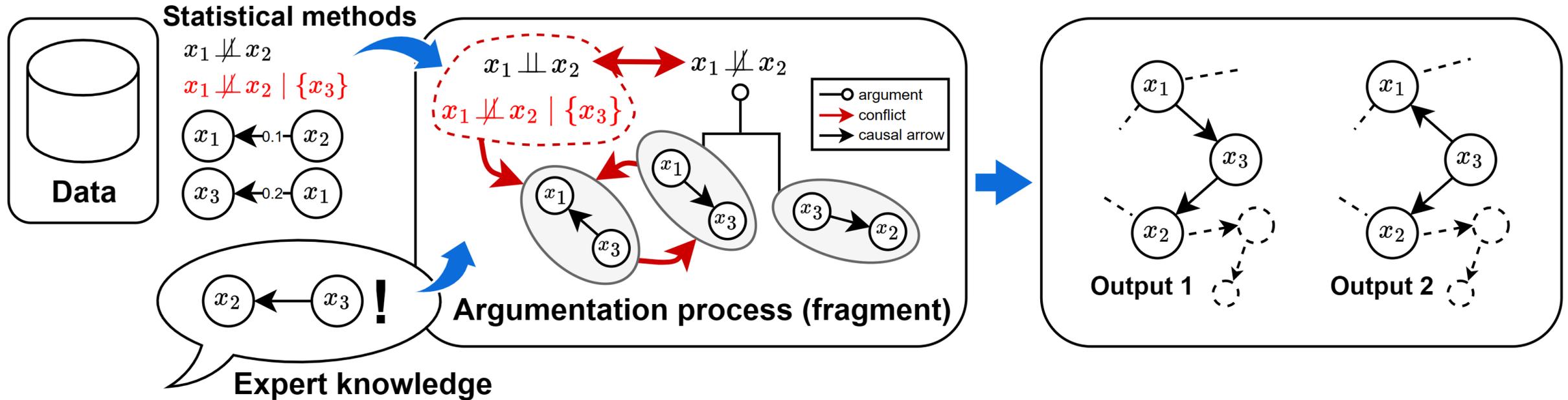
Argumentative Causal Discovery

A Debate about Causality

⌘ Fabrizio Russo (Imperial College London)
 ⌘ Anna Rapberger (Imperial College London)
 ⌘ Francesca Toni (Imperial College London)



In Proceedings of the 21st International
 Conference on Principles of Knowledge
 Representation and Reasoning – Reasoning,
 Learning & Decision-Making Track. Pages 938–949.



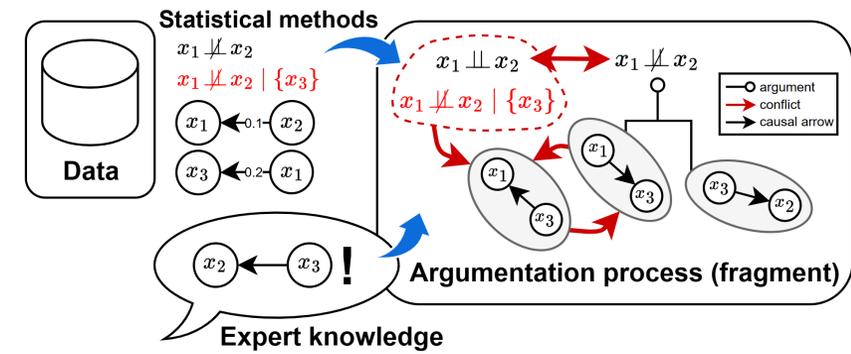
Causal ABA (Assumption-based Argumentation)

Assumptions

$$\{arr_{xy} \mid x, y \in \mathbf{V}, x \neq y\} \cup \{noe_{xy} \mid x, y \in \mathbf{V}, x \neq y\}$$

$$\{(x \perp\!\!\!\perp y \mid \mathbf{Z}) \mid \mathbf{Z} \subseteq \mathbf{V}, x, y \in \mathbf{V} \setminus \mathbf{Z}, x \neq y\}.$$

$$\{bp_{p|z} \mid p \text{ is a } x\text{-}y\text{-path}, \mathbf{Z} \subseteq \mathbf{V} \setminus \{x, y\}\}$$



Rules

- $\bar{a} \leftarrow b, a \neq b, a, b \in \{arr_{xy}, arr_{yx}, noe_{xy}\}, x, y \in \mathbf{V} \implies$ **Can only choose one assumption**

- $\overline{arr_{x_i x_{i+1}}} \leftarrow arr_{x_1 x_2}, \dots, arr_{x_{k-1} x_k}$ for each sequence $x_1 \dots x_k$ with $x_1 = x_k$, for each $1 \leq i < k \implies$ **Cycles are not allowed**

$$dpath_{xy} \leftarrow arr_{xy}$$

$$dpath_{xz} \leftarrow dpath_{xy}, arr_{yz}$$

$$e_{xy} \leftarrow arr_{xy}$$

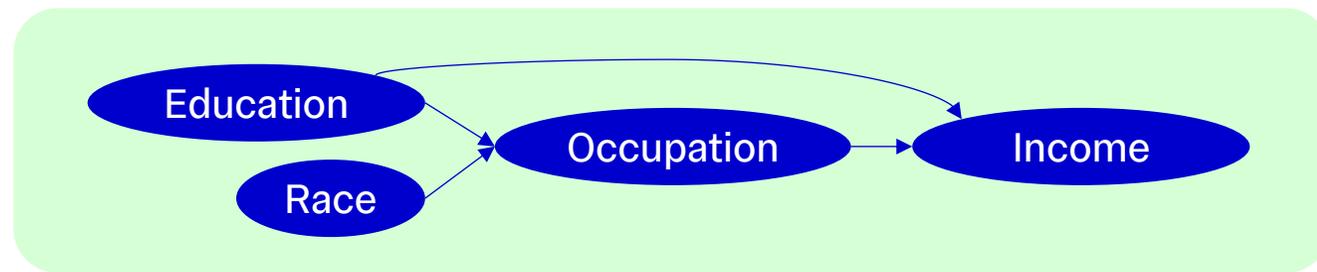
$$e_{xy} \leftarrow arr_{yx} \quad \overline{noe_{xy}} \leftarrow e_{xy}$$

$$\overline{x \perp\!\!\!\perp y \mid \mathbf{Z}} \leftarrow pz$$

- $x \perp\!\!\!\perp y \mid \mathbf{Z} \leftarrow bp_{p_1|z}, \dots, bp_{p_k|z}$ where p_1, \dots, p_k denote all paths between x and y ;

- $ap_{p|z} \leftarrow p$ for each \mathbf{Z} -active p

Income Prediction Example



$E \not\perp I$	$p = 0.00$	$S = 1.00$
$E \not\perp O$	$p = 0.00$	$S = 1.00$
$R \not\perp O$	$p = 0.00$	$S = 1.00$
$O \not\perp I$	$p = 0.00$	$S = 1.00$
$R \perp E$	$p = 0.46$	$S = 0.71$
$R \not\perp I$	$p = 0.05$	$S = 0.52$
$E \not\perp I \mid \{R\}$	$p = 0.00$	$S = 0.50$
$E \not\perp I \mid \{O\}$	$p = 0.00$	$S = 0.50$
$E \not\perp O \mid \{R\}$	$p = 0.00$	$S = 0.50$
$R \not\perp O \mid \{I\}$	$p = 0.00$	$S = 0.50$
$E \not\perp O \mid \{I\}$	$p = 0.00$	$S = 0.50$
$R \not\perp O \mid \{E\}$	$p = 0.00$	$S = 0.50$

$O \not\perp I \mid \{E\}$	$p = 0.00$	$S = 0.50$
$O \not\perp I \mid \{R\}$	$p = 0.00$	$S = 0.50$
$R \perp E \mid \{O\}$	$p = 0.53$	$S = 0.38$
$R \not\perp I \mid \{O\}$	$p = 0.03$	$S = 0.35$
$R \perp E \mid \{O\}$	$p = 0.33$	$S = 0.32$
$R \perp E \mid \{I\}$	$p = 0.05$	$S = 0.25$
$R \perp E \mid \{O, I\}$	$p = 0.39$	$S = 0.00$
$R \not\perp I \mid \{E, O\}$	$p = 0.00$	$S = 0.00$
$E \not\perp O \mid \{R, I\}$	$p = 0.00$	$S = 0.00$
$R \not\perp I \mid \{E, O\}$	$p = 0.00$	$S = 0.00$
$R \not\perp O \mid \{E, I\}$	$p = 0.03$	$S = 0.00$

Argumentative Causal Discovery

Robust & Interactive



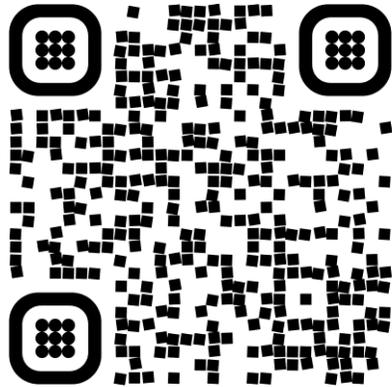
Sound & Complete



Robust to Errors but
Computationally Demanding



Data Check &
Stakeholder Engagement

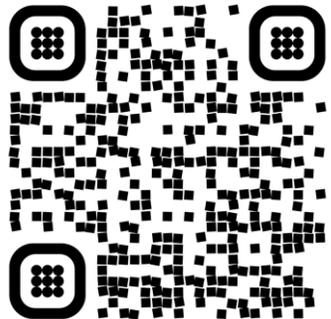
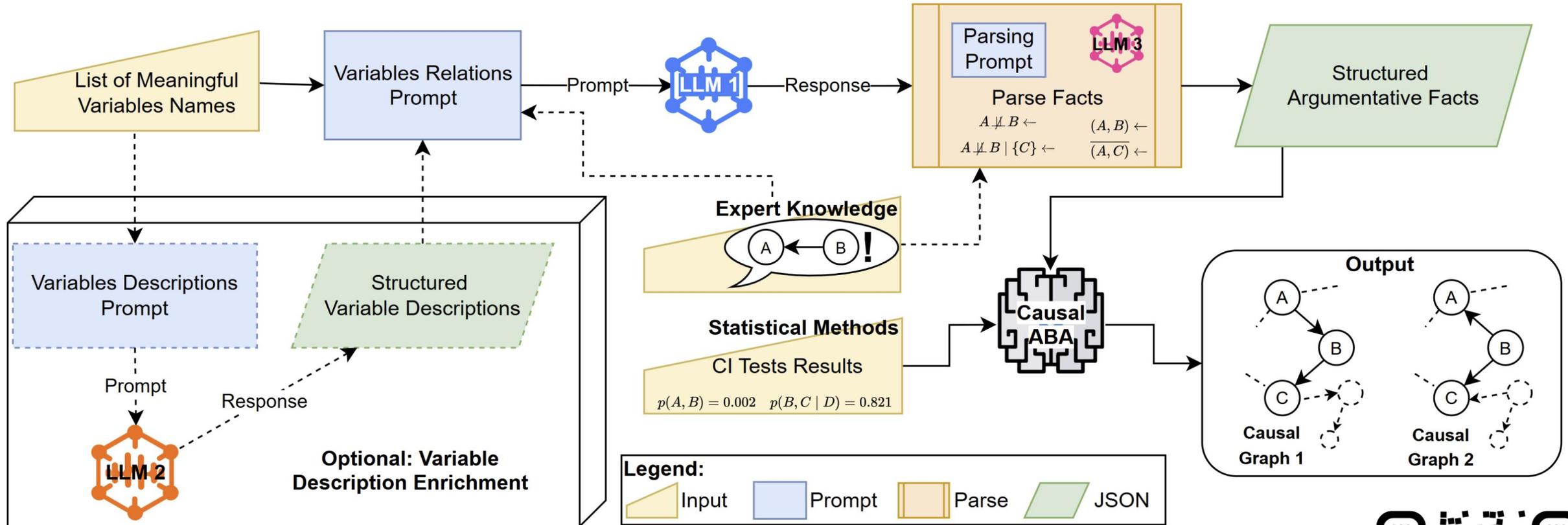


Russo, Rapberger & Toni. *Argumentative Causal Discovery*.
In Proc. of KR 2024.

Causal ABA Ensemble

LLMs as Imperfect Experts

Li, Z. and Russo, F. *Leveraging Large Language Models for Causal Discovery: a Constraint-based, Argumentation-driven Approach.* Arxiv 2602.16481



Causal Discovery Literature

From the 90s



Review of Causal Discovery Methods Based on Graphical Models

Clark Glymour, Kun Zhang* and Peter Spirtes

Department of Philosophy, Carnegie Mellon University, Pittsburgh, PA, United States

A fundamental task in various disciplines of science, including biology, is to find underlying causal relations and make use of them. Causal relations can be seen if interventions are properly applied; however, in many cases they are difficult or even impossible to conduct. It is then necessary to discover causal relations by analyzing statistical properties of purely observational data, which is known as causal discovery or causal structure search. This paper aims to give an introduction to and a brief review of the computational methods for causal discovery that were developed in the past three decades, including constraint-based and score-based methods and those based on functional causal models, supplemented by some illustrations and applications.

Keywords: directed graphical causal models, causal discovery, conditional independence, statistical independence, structural equation models, non-Gaussian distribution, non-linear models

D'ya Like DAGs? A Survey on Structure Learning and Causal Discovery

MATTHEW J. VOWELS, NECATI CIHAN CAMGOZ, and RICHARD BOWDEN,
CVSSP, University of Surrey, U.K.

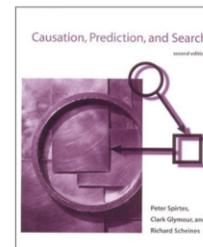
Causal reasoning is a crucial part of science and human intelligence. In order to discover causal relationships from data, we need structure discovery methods. We provide a review of background theory and a survey of methods for structure discovery. We primarily focus on modern, continuous optimization methods, and provide reference to further resources such as benchmark datasets and software packages. Finally, we discuss the assumptive leap required to take us from structure to causality.

CCS Concepts: • **Mathematics of computing** → **Causal networks**; • **Computing methodologies** → **Machine learning**; **Causal reasoning and diagnostics**;

Additional Key Words and Phrases: Causality, causal discovery, directed acyclic graphs, DAGs, structure learning, survey

ACM Reference format:

Matthew J. Vowels, Necati Cihan Camgoz, and Richard Bowden. 2022. D'ya Like DAGs? A Survey on Structure Learning and Causal Discovery. *ACM Comput. Surv.* 55, 4, Article 82 (November 2022), 36 pages. <https://doi.org/10.1145/3527154>



Adaptive Computation And Machine Learning Series

Causation, Prediction, and Search (Second Edition)

By Peter Spirtes, Clark Glymour, Richard Scheines

The MIT Press

DOI: <https://doi.org/10.7551/mitpress/1754.001.0001>

ISBN electronic: 9780262284158

In Special Collection: CogNet

Publication date: 2001

Recap

& Conclusion

Causal Models

- Causal Discovery
- Structural Equations

Contestable NN

- Method & Process to align NNs to SMEs' causal view

Causal ABA

- Method to increase consistency guarantees and allow integrations and contestability

IMPERIAL

Questions?

Feedback Form

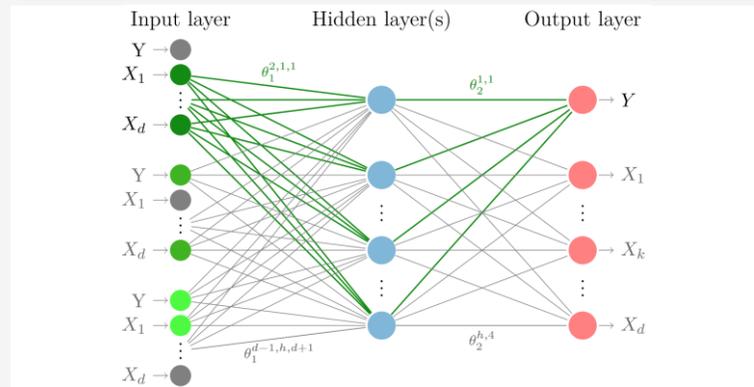


fabrizio@imperial.ac.uk

Appendix

Causal Discovery for Trustworthy AI Methods

Contestable Neural Networks



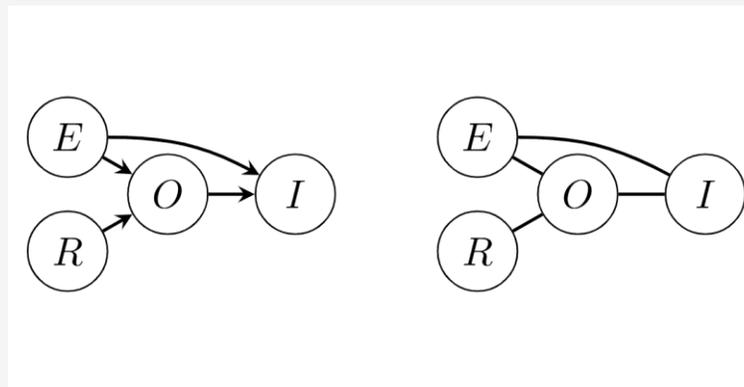
Russo & Toni. *Causal Discovery and Knowledge Injection for Contestable Neural Networks*. In Proc. of ECAI 2023

Learn a causal Graph while fitting a neural network

Expose the learned graph to SMEs and allow them to modify it

Inject the graph back into the NN

Causal Discovery with Shapley Values



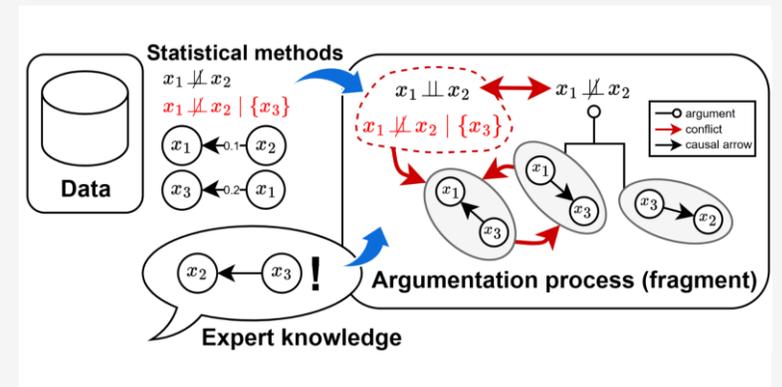
Russo & Toni. *Shapley-PC: Constraint-based Causal Structure Learning with a Shapley Inspired Framework*. In Proc. of CLeaR 2025

Learn a causal graph with asymptotic guarantees

Identify colliders given a skeleton using Shapley Values

Embed the decision rule into the PC algorithm

Argumentative Causal Discovery



Russo, Rapberger & Toni. *Argumentative Causal Discovery*. In Proc. of KR 2024.

Represent DAGs and d-Separation as a debate

Allow external knowledge to be input into the framework

Guarantee consistency of output DAG with a subset of the input

Income Prediction Case Study

Adult dataset (Becker and Kohavi, 1996) UCI repository

Adult Dataset ($ V = 14$)				
Data (N)	CASTLE $ E = 210$	Injected $ E = 46$	Partial $ E = 116$	Contested $ E = 30$
100	0.67 (0.03)	0.69 .	0.66	0.69 .
500	0.72 (0.04)	0.74 *	0.71	0.74 *
1000	0.75 (0.03)	0.76	0.74	0.76
2000	0.74 (0.03)	0.77 ***	0.76 *	0.77 ***
5000	0.75 (0.03)	0.79 ***	0.76	0.79 ***
10000	0.75 (0.02)	0.85 ***	0.76 .	0.85 ***
20000	0.76 (0.02)	0.86 ***	0.77 .	0.86 ***

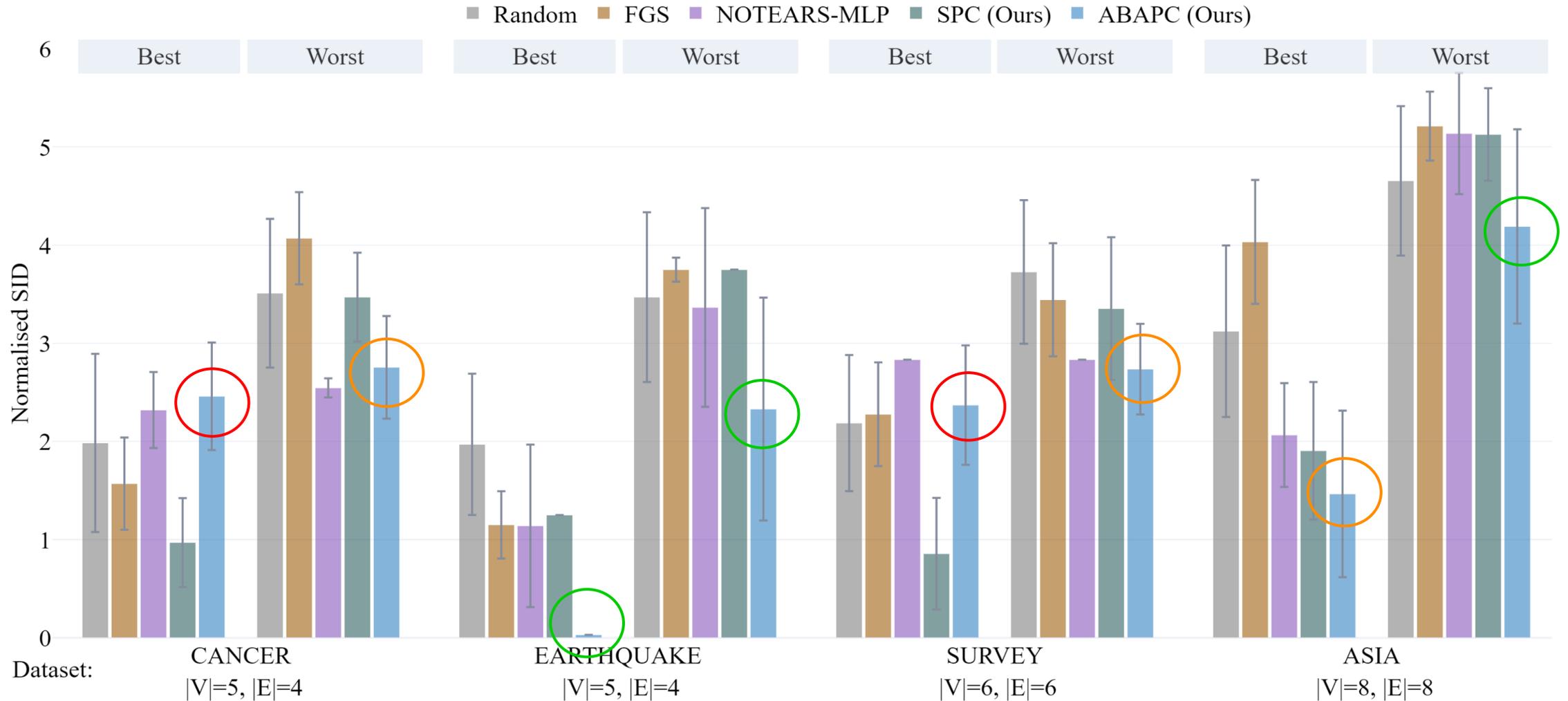
More Datasets

From Credit Risk to House Prices

Metric Data & Method	$N=100$	$N=500$	$N=1K$	$N=2K$	$N=5K$	$N=10K$	$N=20K$	
AUC \uparrow	Adult ($ V = 14$)							
	CASTLE ($ E = 210$)	0.67 ± 0.03	0.72 ± 0.04	0.75 ± 0.03	0.74 ± 0.03	0.75 ± 0.03	0.75 ± 0.02	0.76 ± 0.02
	<i>Injected</i> ($ E = 46$)	0.69.	0.74*	0.76	0.77***	0.79***	0.85***	0.86***
	<i>Partial</i> ($ E = 116$)	0.66	0.71	0.74	0.76*	0.76	0.76.	0.77.
	<i>Refined</i> ($ E = 30$)	0.69.	0.74*	0.76	0.77***	0.79***	0.85***	0.86***
HELOC ($ V = 23$)	CASTLE ($ E = 552$)	0.75 ± 0.02	0.79 ± 0.01	0.78 ± 0.01	0.79 ± 0.01	0.79 ± 0.01	0.80 ± 0.01	
	<i>Injected</i> ($ E = 85$)	0.74	0.78^{***}	0.78	0.78^{***}	0.79	0.79^{***}	
	California ($ V = 8$)							
MSE \downarrow	CASTLE ($ E = 72$)	7.05 ± 12.81	2.33 ± 1.39	2.96 ± 4.12	3.86 ± 3.68	4.91 ± 7.41	1.74 ± 1.70	0.66 ± 0.08
	<i>Injected</i> ($ E = 31$)	2.94^{***}	2.25	1.68^{**}	1.71^{***}	1.51^{***}	1.16^*	1.02^{**}
	Boston ($ V = 14$)							
CASTLE ($ E = 182$)	112.0 ± 91.1	21.95 ± 6.84						
<i>Injected</i> ($ E = 48$)	86.17^{***}	20.45*						

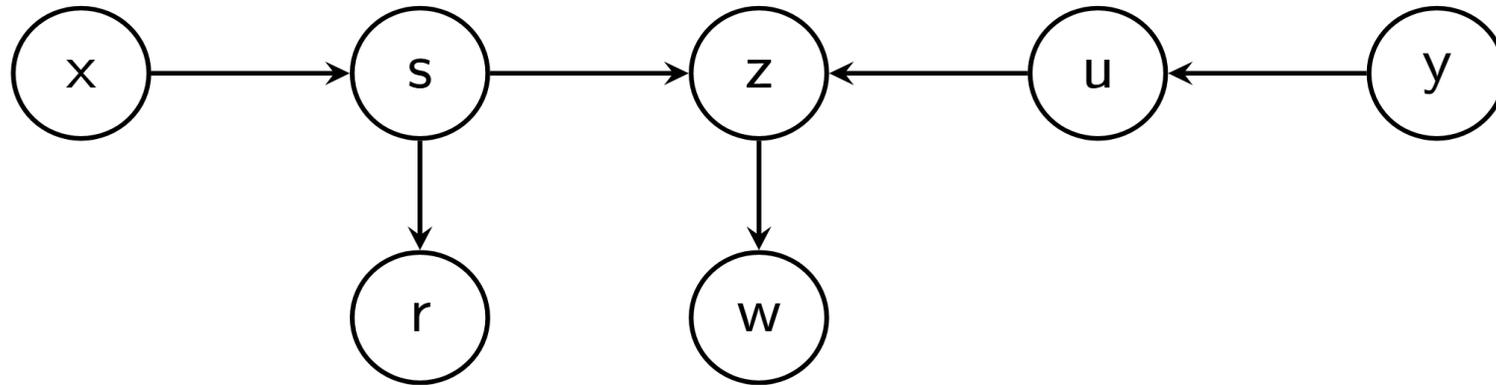
ABAPC Reconstruction Accuracy

on bnlearn datasets (Scutari, 2014)



Connect DAGs and Data: d-Separation

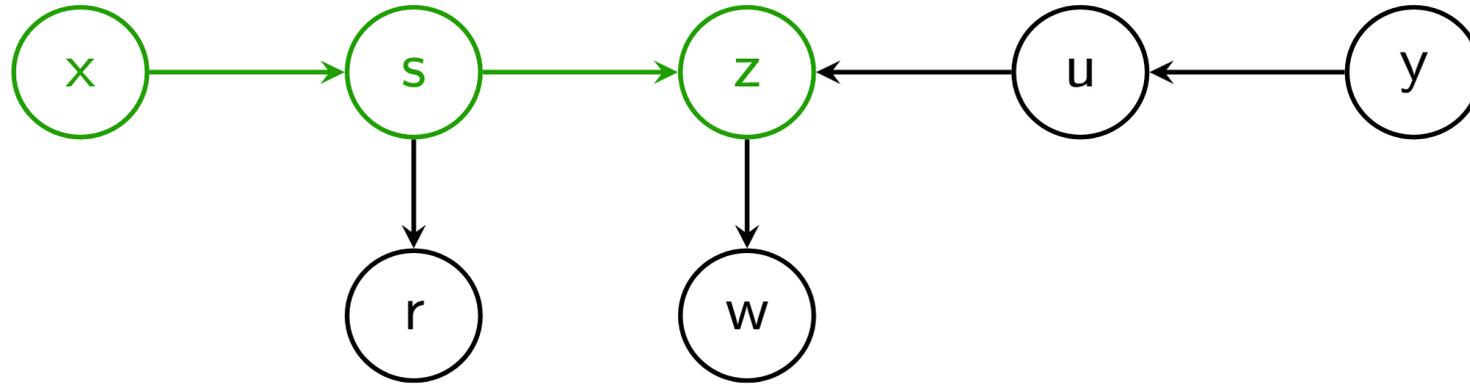
(Pearl 2009)



Unconditional
 $x \not\perp_G z$

Connect DAGs and Data: d-Separation

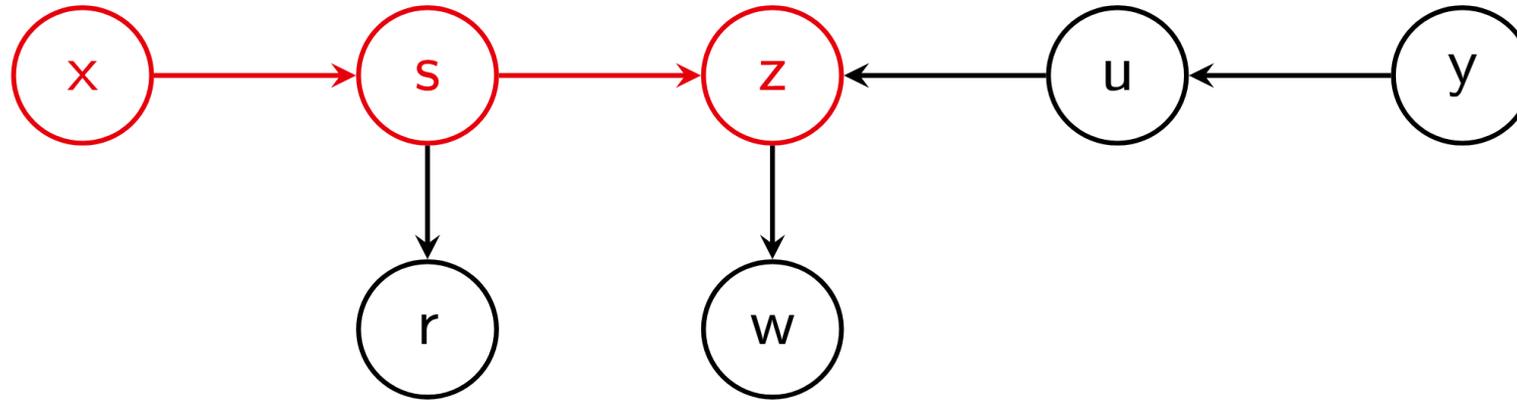
(Pearl 2009)



Unconditional
 $x \not\perp_G z$ **Active**

Connect DAGs and Data: d-Separation

(Pearl 2009)

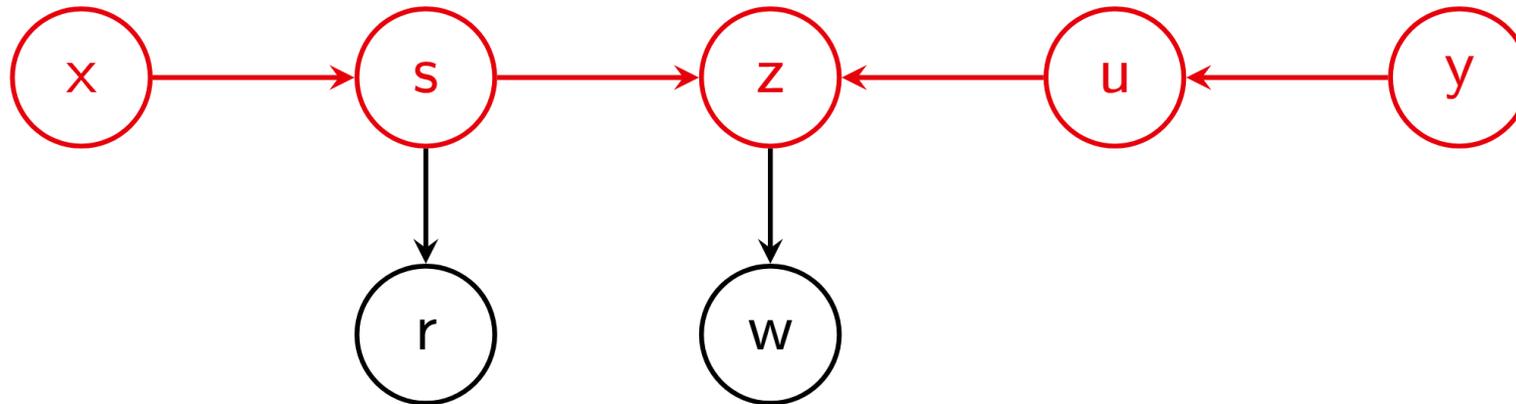


Unconditional
 $x \not\perp_G z$

d-separate by conditioning
 $x \perp_G z \mid s$ **Blocked**

Connect DAGs and Data: d-Separation

(Pearl 2009)



Unconditional

$x \not\perp_G z$

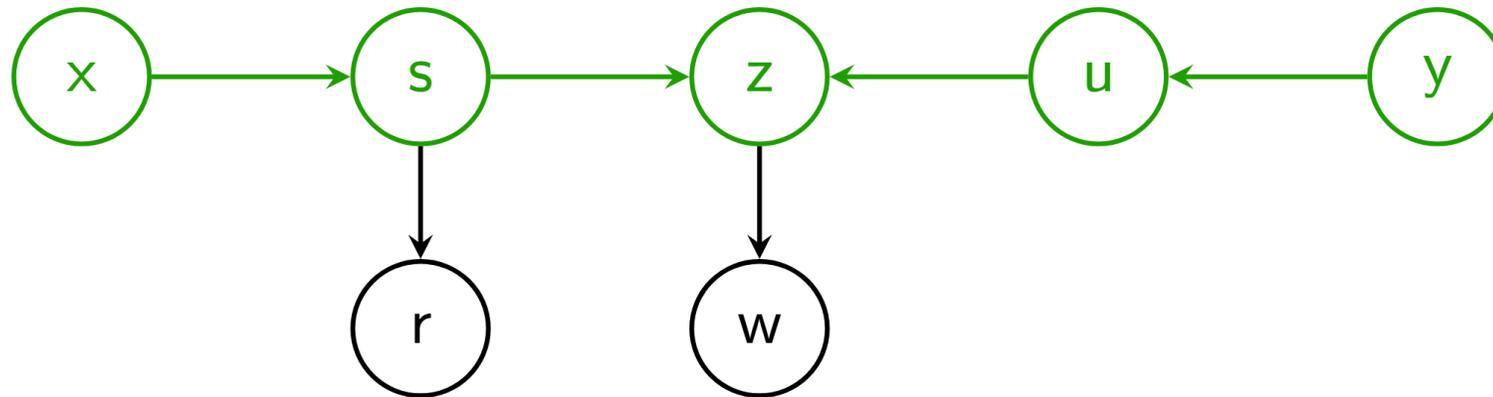
$x \perp_G y$ **Blocked**

d-separate by conditioning

$x \perp_G z \mid s$

Connect DAGs and Data: d-Separation

(Pearl 2009)



Unconditional

$x \not\perp_G z$

$x \perp_G y$

d-separate by conditioning

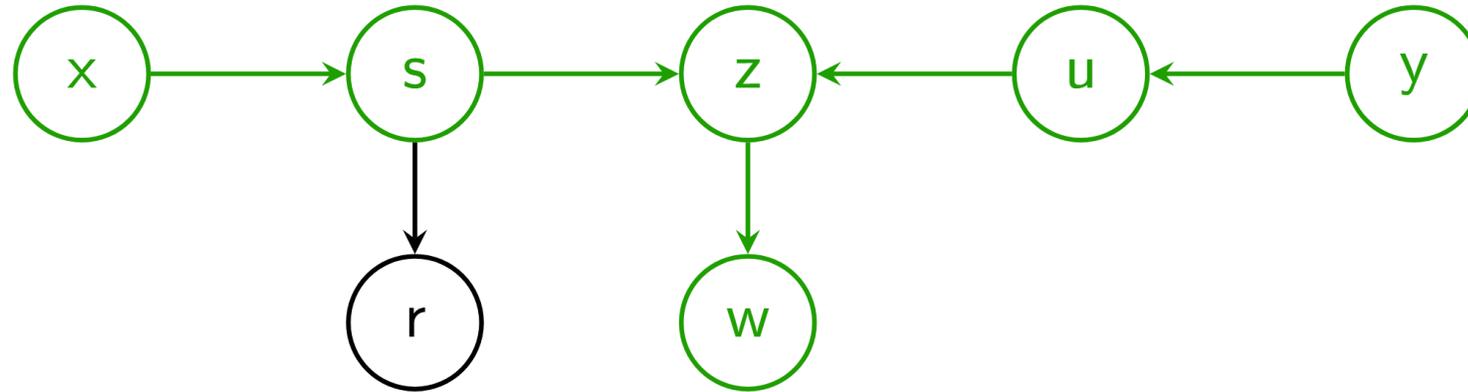
$x \perp_G z \mid s$

d-connect by conditioning

$x \not\perp_G y \mid z$ **Active**

Connect DAGs and Data: d-Separation

(Pearl 2009)



Unconditional

$x \not\perp_G z$

$x \perp_G y$

d-separate by conditioning

$x \perp_G z \mid s$

d-connect by conditioning

$x \not\perp_G y \mid z$

$x \not\perp_G y \mid w$ **Active**