

# CSDS 452 Causality and Machine Learning

## **Lecture 16: Causal Explanation**

Instructor: Jing Ma

Fall 2024, CDS@CWRU

# Outline

- Causal explanation: general introduction
- Decision based interpretability
  - Counterfactual explanation
  - Recourse
- Data-based interpretability
- Model-based interpretability (attribution)

# Case study

*Edward*



**Age:** 28

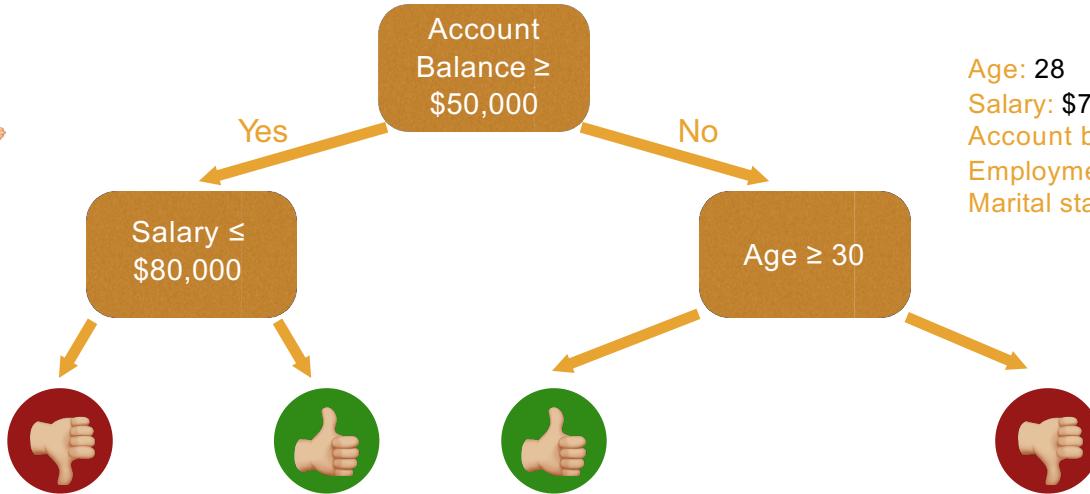
**Salary:** \$75,000

**Account balance:** \$25,000

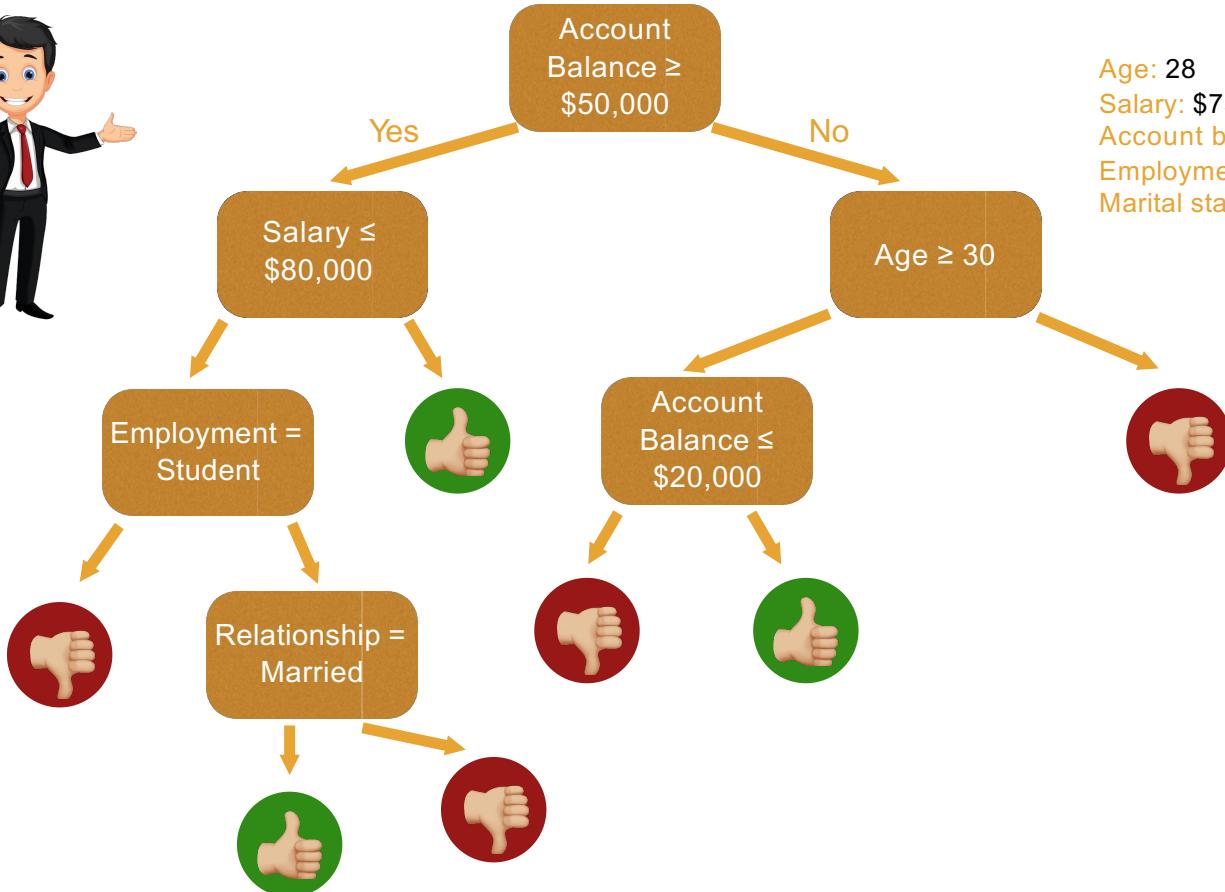
**Employment:** software eng.

**Marital status:** single

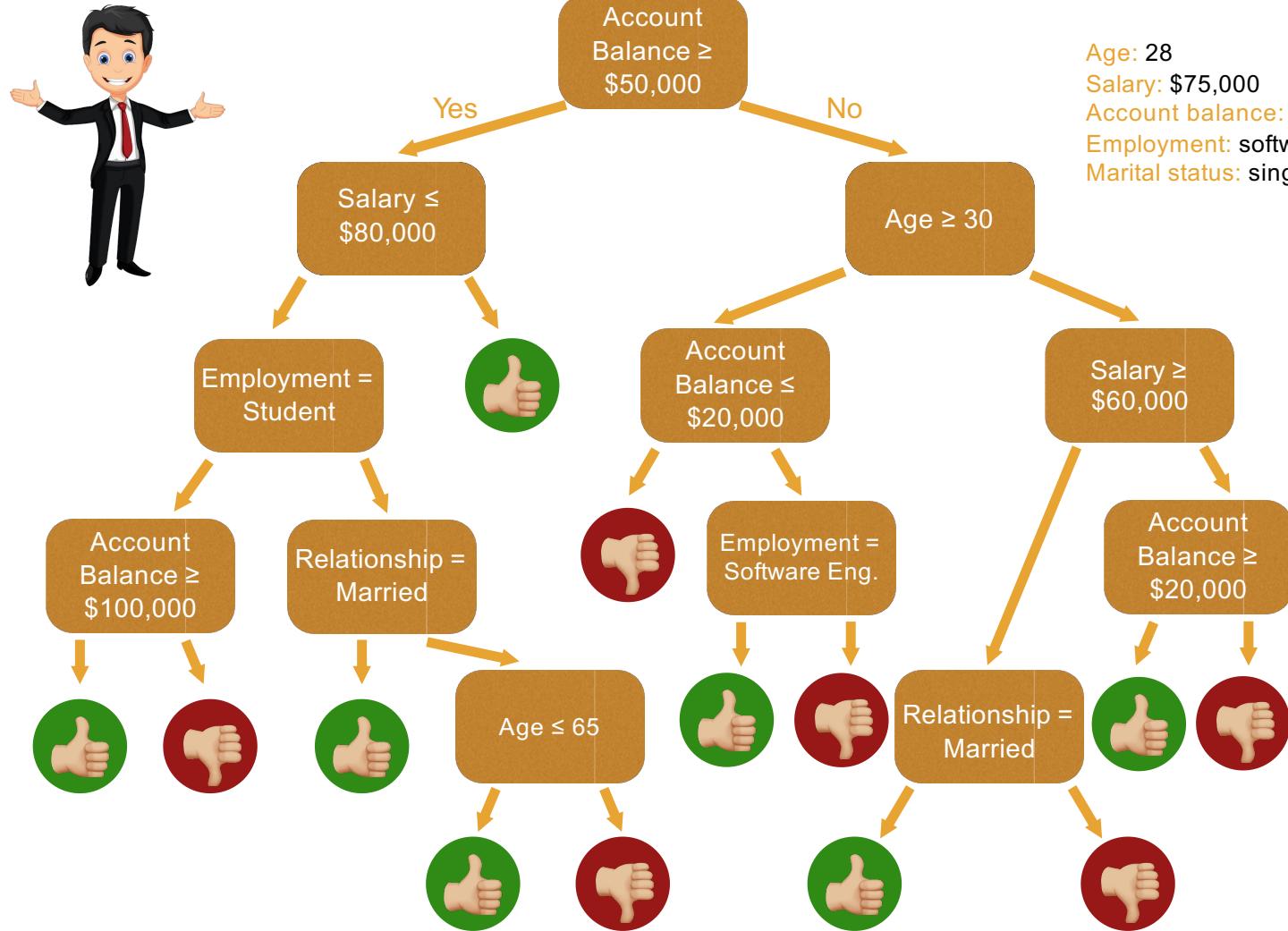




Age: 28  
Salary: \\$75,000  
Account balance: \\$25,000  
Employment: software eng.  
Marital status: single



Age: 28  
Salary: \$75,000  
Account balance: \$25,000  
Employment: software eng.  
Marital status: single



# Stakeholders



What are the **most influential features** towards the decision?

Is the system “fair” by relying on sensitive attributes such as *age* and *marital status*?

I didn’t get the loan; **why not?** And **what should I do** to get it next time I apply?

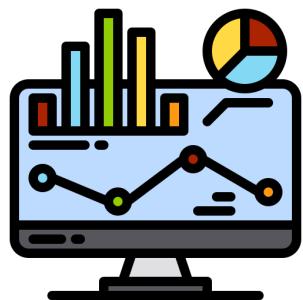
# eXplainable AI (XAI)

Interpretable ML    The ability to explain or to present in understandable terms to a human [Doshi-Velez & Kim, 2017]

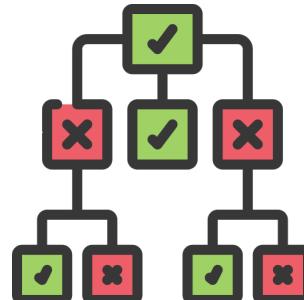


Goals of XAI    Draw insights from data, models, and/or decisions

Causal Insights    Answer what, how, and why (not) questions



Data-based Interpretability



Model-based Interpretability

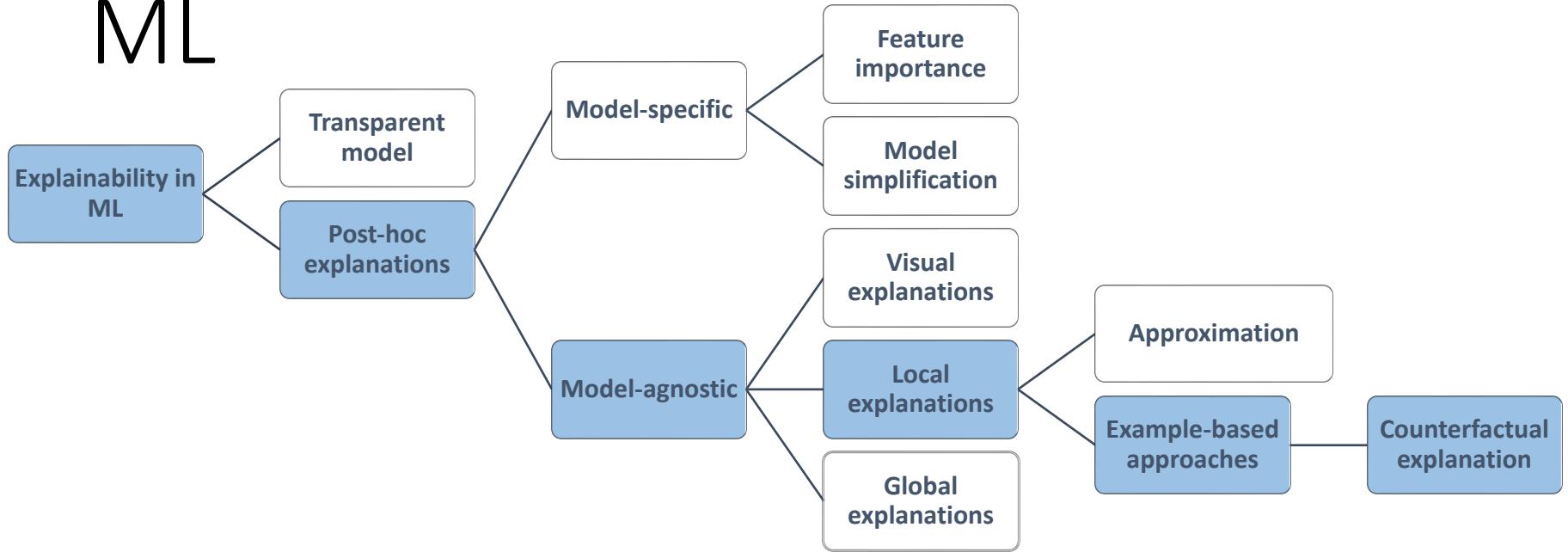


Decision-based Interpretability

# Outline

- Causal explanation: general introduction
- Decision based interpretability
  - Counterfactual explanation
  - Recourse
- Data-based interpretability
- Model-based interpretability (attribution)

# A Broad Picture of Explainability in ML



- This figure describes the role of the “major trend” of counterfactual explanation (CFE) in ML explainability
- Some CFE methods are in different categories (e.g., model specific)

[1] Verma S, Dickerson J, Hines K. Counterfactual explanations for machine learning: A review[J]. arXiv preprint arXiv:2010.10596, 2020.

# Outline

- Causal explanation: general introduction
- Decision based interpretability
  - Counterfactual explanation
  - Recourse
- Data-based interpretability
- Model-based interpretability (attribution)

# Counterfactual Explanation

- **Key question:** what **small changes** could be made to the input features of an instance to change its prediction?
- Motivation example:
  - In loan application, for an applicant who was rejected, what small change (e.g., improve education level) can the applicant make to achieve a desired outcome?

Rejected for a Personal Loan?



# Basic Desiderata

- **Validity:** generated counterfactuals have desired labels

$$\arg \min_{x'} d(x, x') \text{ subject to } f(x') = y' \quad \begin{array}{l} \text{Desired label} \\ \text{Original goal} \end{array}$$
$$\arg \min_{x'} \max_{\lambda} \lambda(f(x') - y')^2 + d(x, x') \quad \begin{array}{l} \text{Differentiable, unconstrained objective} \\ \uparrow \\ \text{Distance (e.g., L1/L2)} \end{array}$$

- **Sparsity:** a counterfactual ideally should change smaller number of features in order to be most effective

# Other Desiderata (e.g.)

- **Data Manifold closeness:** a generated counterfactual is realistic in the sense that it is near the training data  $l(x'; \mathcal{X}) \leftarrow$  training set
  - include a penalty for adhering to the data manifold defined by the training set

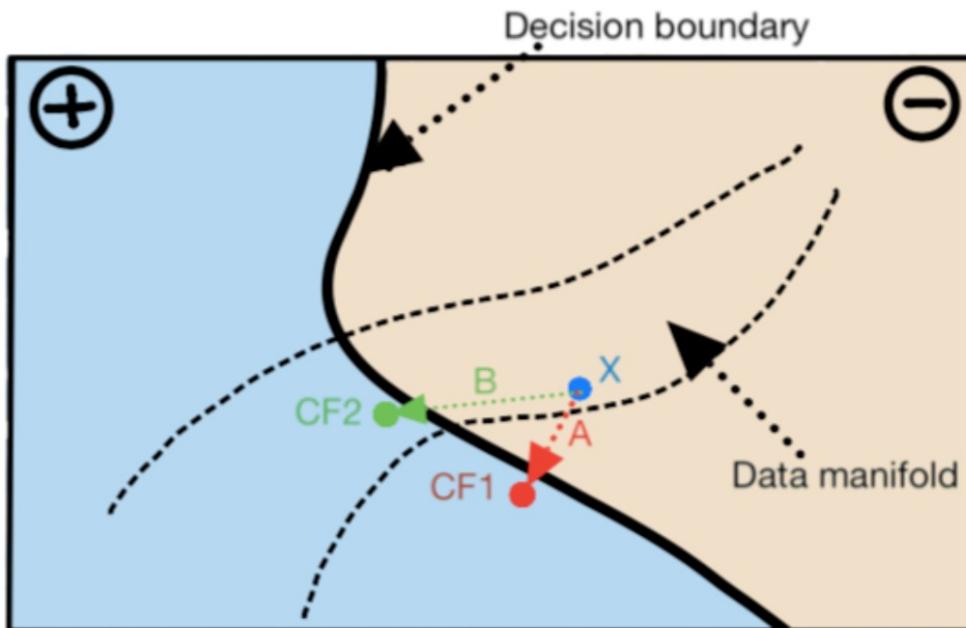


Figure: Two possible paths for a datapoint (shown in blue), originally classified in the negative class, to cross the decision boundary. The end points of both the paths (shown in red and green) are valid counterfactuals for the original point. Note that the red path is the shortest, whereas the green path adheres closely to the manifold of the training data, but is longer.

# Properties of the Approaches

- **Model access**
  - access to complete model internals
  - access to gradients
  - access to only the prediction function (black-box)
- **Model agnostic**
  - some algorithms are model-specific, e.g., only work for those models like tree ensembles or with gradients
- **Optimization amortization**
  - **Amortized Inference:** whether the algorithm can generate counterfactuals for **multiple input datapoints** without optimizing separately
  - **Multiple counterfactual:** whether the algorithm can generate **multiple counterfactual** for a single input datapoint
- **Counterfactual (CF) attributes**
  - sparsity, data manifold adherence, and causality

# Key Properties of Existing CFE Works

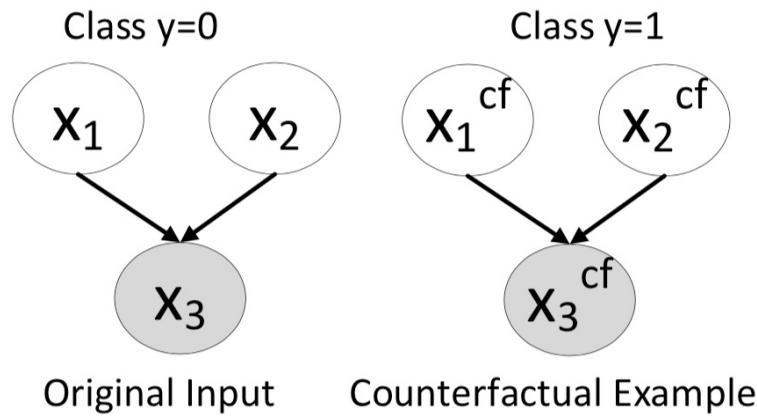
Paper	Assumptions		Optimization amortization			CF attributes		CF opt. problem attributes		
	Model access	Model domain	Amortized Inference	Multiple CF	Sparsity	Data manifold	Causal relation	Feature preference	Categorical func	dist.
[72]	Black-box	Agnostic	No	No	Changes iteratively	No	No	Yes	-	
[111]	Gradients	Differentiable	No	No	L1	No	No	No	-	
[104]	Complete	Tree ensemble	No	No	No	No	No	No	-	
[74]	Black-box	Agnostic	No	No	L0 and post-hoc	No	No	No	-	
[57]	Black-box	Agnostic	No	Yes	Flips min. split nodes	No	No	No	Indicator	
[29]	Gradients	Differentiable	No	No	L1	Yes	No	No	-	
[56]	Black-box	Agnostic	No	No	No	No	No	No <sup>2</sup>	-	
[95]	Complete	Linear	No	Yes	L1	No	No	No	N.A. <sup>3</sup>	
[107]	Complete	Linear	No	No	Hard constraint	No	No	Yes	-	
[98]	Black-box	Agnostic	No	Yes	No	No	No	Yes	Indicator	
[30]	Black-box or gradient	Differentiable	No	No	L1	Yes	No	No	-	
[91]	Black-box	Agnostic	No	No	No	No	No	No	-	
[61]	Gradients	Differentiable	No	No	No	Yes	No	No	-	
[90]	Gradients	Differentiable	No	No	No	No	No	No	-	
[113]	Black-box	Agnostic	No	No	Changes one feature	No	No	No	-	
[85]	Gradients	Differentiable	No	Yes	L1 and post-hoc	No	No	No	Indicator	
[89]	Black-box	Agnostic	No	No	No	Yes <sup>4</sup>	No	No	-	

# Until now, no causality is involved

- In causal inference, **counterfactuals** refer to a different version that does not actually happen, e.g.,
  - What if I hadn't taken the medicine? (I took the medicine)
  - If I were a boy, I would ... (I am a girl)
- The word “**counterfactual**” is often overused!
  - In many papers, counterfactual just refer to anything different than the current one without causality involved.
- Obviously, a counterfactual should take **causal relations** into consideration

# Desiderata: Causality

- **Causality:** changing one feature in the real world affects other features
  - E.g., getting a new educational degree necessitates increases the age



**Standard Proximity Loss:**  $\text{dist}(x_3, x_3^{cf})$   
**Causal Proximity Loss:**  $\text{dist}(\text{f}(x_1^{cf}, x_2^{cf}), x_3^{cf})$

Structural equation in causal model

# Evaluation

- Commonly used datasets
  - Image - MNIST
  - Tabular - Adult income, German credit, Compas recidivism, etc.
- Metrics
  - **Validity:** the **ratio of the counterfactuals that have the desired class label** to the total number of counterfactuals.
  - **Proximity:** the **distance** of a counterfactual from the input datapoint.
  - **Sparsity:** the number of modified features
  - **Diversity:** diversity is encouraged by maximizing the distance between the multiple counterfactuals
  - **Causal constraint:** whether the counterfactuals **satisfy the causal relation** between features.

# Counterfactual Explanation on Graph

- Counterfactual explanations for graphs: the minimal perturbation to the input (graph) data such that the prediction changes.
- Motivation & applications:
  - Drug discovery: if the model predicts it does not have this desirable property, CFE can help identify the minimal change one should make to this molecule
- Challenges:
  - CFE contains different modalities
  - Perturbation on graph structure is non-differentiable
  - Some patterns in graph structure (e.g., circles) are hard to be learned and perturbed for CFE

# CF-GNNExplainer: Method

- CF-GNNExplainer<sup>[1]</sup>
  - based on GNN
  - focus on node classification task
  - focus on perturbing graph structure
- Main idea of the method
  - iteratively remove edges (learn a perturbation matrix) from the original adjacency matrix based on matrix sparsification techniques
  - track of the perturbations that lead to a change in prediction
  - return the perturbation with the smallest change w.r.t. the number of edges

$$\mathcal{L} = \mathcal{L}_{pred}(v, \bar{v} | f, g) + \beta \mathcal{L}_{dist}(v, \bar{v}),$$

Prediction model

Original data    Counterfactual data for the node    CFE generator    Element-wise difference

# CF-GNNExplainer: Evaluation

- Dataset
  - tree-cycles, tree-grids, ba-shapes
  - consists of (i) a base graph, (ii) motifs that are attached to random nodes of the base graph, and (iii) additional edges that are randomly added to the overall graph.
- Task: predict whether the nodes are part of the motif
- Metric:
  - **Fidelity**: the proportion of nodes where the original predictions match the prediction for the explanations
  - **Explanation Size**: the number of removed edges
  - **Sparsity**: the proportion of edges that are removed
  - **Accuracy**: the proportion of explanations that are “correct”.
    - only compute accuracy for nodes that are predicted as in the motifs
    - consider an explanation to be correct if it exclusively involves edges that are inside the motifs

# CF-GNNExplainer: Evaluation

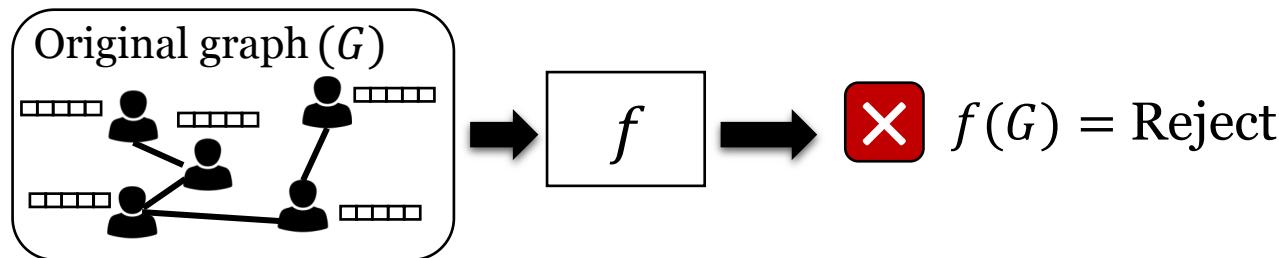
- **Baselines:**
  - **Random**: randomly remove edges in each subgraph, repeat  $K$  times and keep track of the most minimal perturbation resulting in a CF
  - **Only-1hop**: only keeps all edges in the 1-hop ego graph
  - **Rm-1hop**: removes all edges in the 1-hop ego graph
  - **GNNExplainer**: remove the subgraph generated by GNNExplainer (the most relevant subgraph for the prediction)

**Table 1: Results comparing our method to the baselines. Below each metric, ▼ indicates a low value is desirable, while ▲ indicates a high value is desirable.**

Metric	TREE-CYCLES				TREE-GRID				BA-SHAPES			
	Fid. ▼	Size ▼	Spars. ▲	Acc. ▲	Fid. ▼	Size ▼	Spars. ▲	Acc. ▲	Fid. ▼	Size ▼	Spars. ▲	Acc. ▲
RANDOM	<b>0.00</b>	4.70	0.79	0.63	<b>0.00</b>	9.06	0.75	0.77	<b>0.00</b>	503.31	0.58	0.17
ONLY-1HOP	0.32	15.64	0.13	0.45	0.32	29.30	0.09	0.72	0.60	504.18	0.05	0.18
RM-1HOP	0.46	2.11	0.89	—	0.61	2.27	0.92	—	0.21	10.56	0.97	<b>0.99</b>
GNNEPLAINER	0.55	6.00	0.57	0.46	0.34	8.00	0.68	0.74	0.81	6.00	0.81	0.27
CF-GNNEXPLAINER	0.21	<b>2.09</b>	<b>0.90</b>	<b>0.94</b>	0.07	<b>1.47</b>	<b>0.94</b>	<b>0.96</b>	0.39	<b>2.39</b>	<b>0.99</b>	0.96

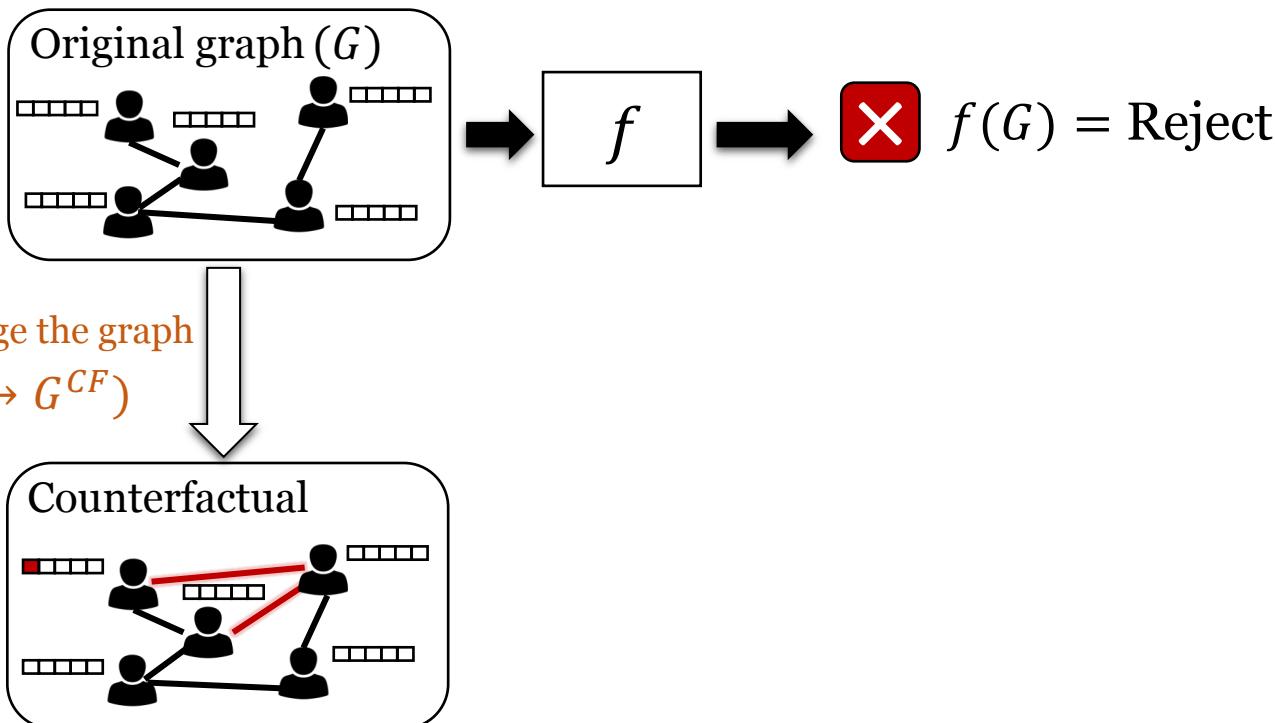
# Counterfactual Explanation for Graph Classification

Example: a graph ML model  $f$  trained for grant application decision-making.



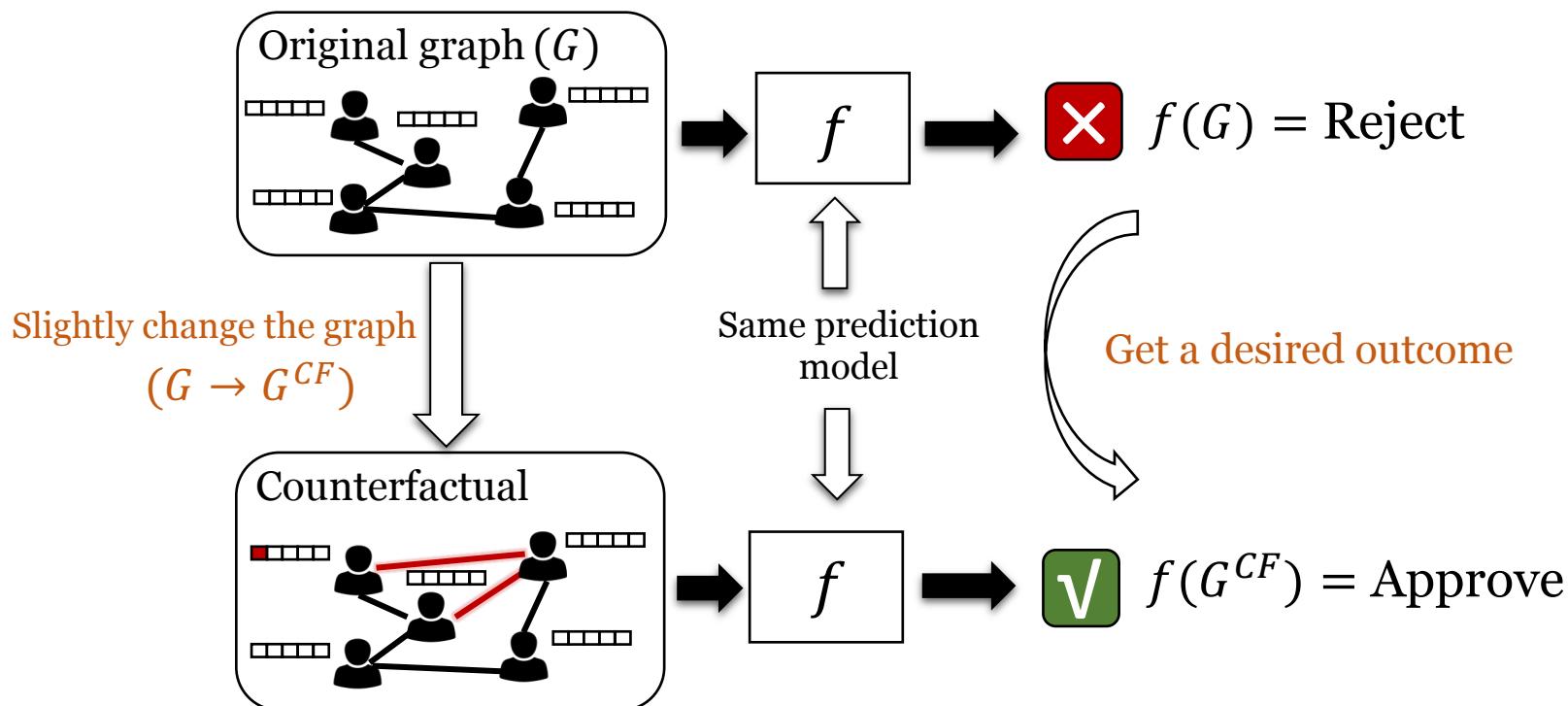
# Counterfactual Explanation for Graph Classification

Example: a graph ML model  $f$  trained for grant application decision-making.



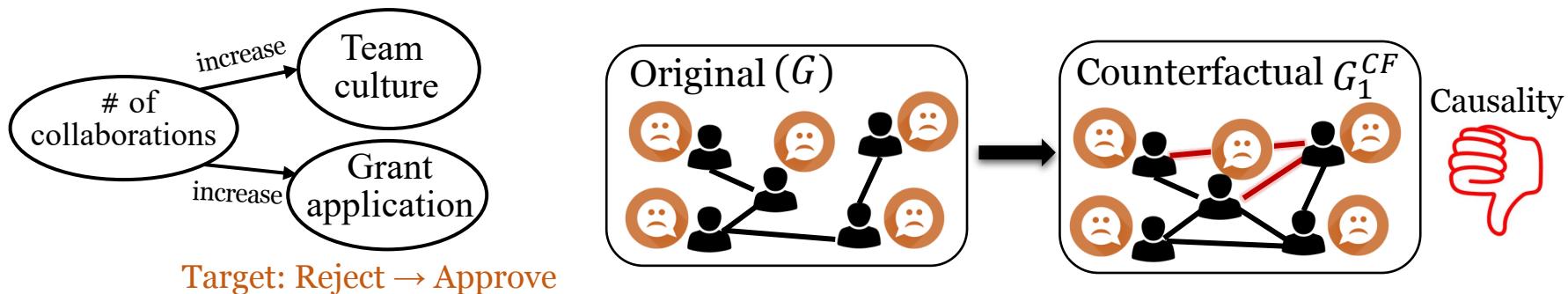
# Counterfactual Explanation for Graph Classification

Example: a graph ML model  $f$  trained for grant application decision-making.



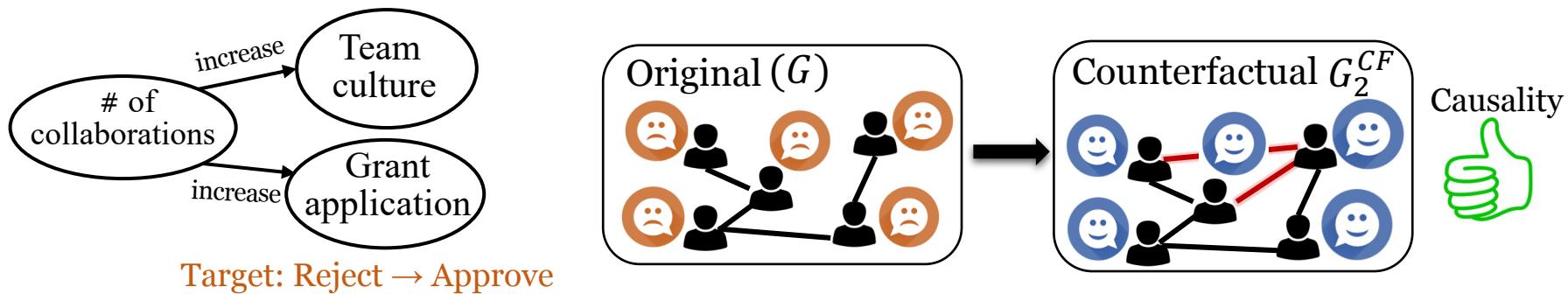
# Challenges of CFE on Graphs

- Optimization
  - The space of perturbation operations on graphs (e.g., add/remove nodes/edges) is discrete, disorganized, and vast
- Generalization
  - Most existing CFE methods on graphs generate counterfactuals for each graph separately, and cannot generalize to unseen graphs
- Causality
  - It is challenging to generate counterfactuals that are consistent with the underlying causality



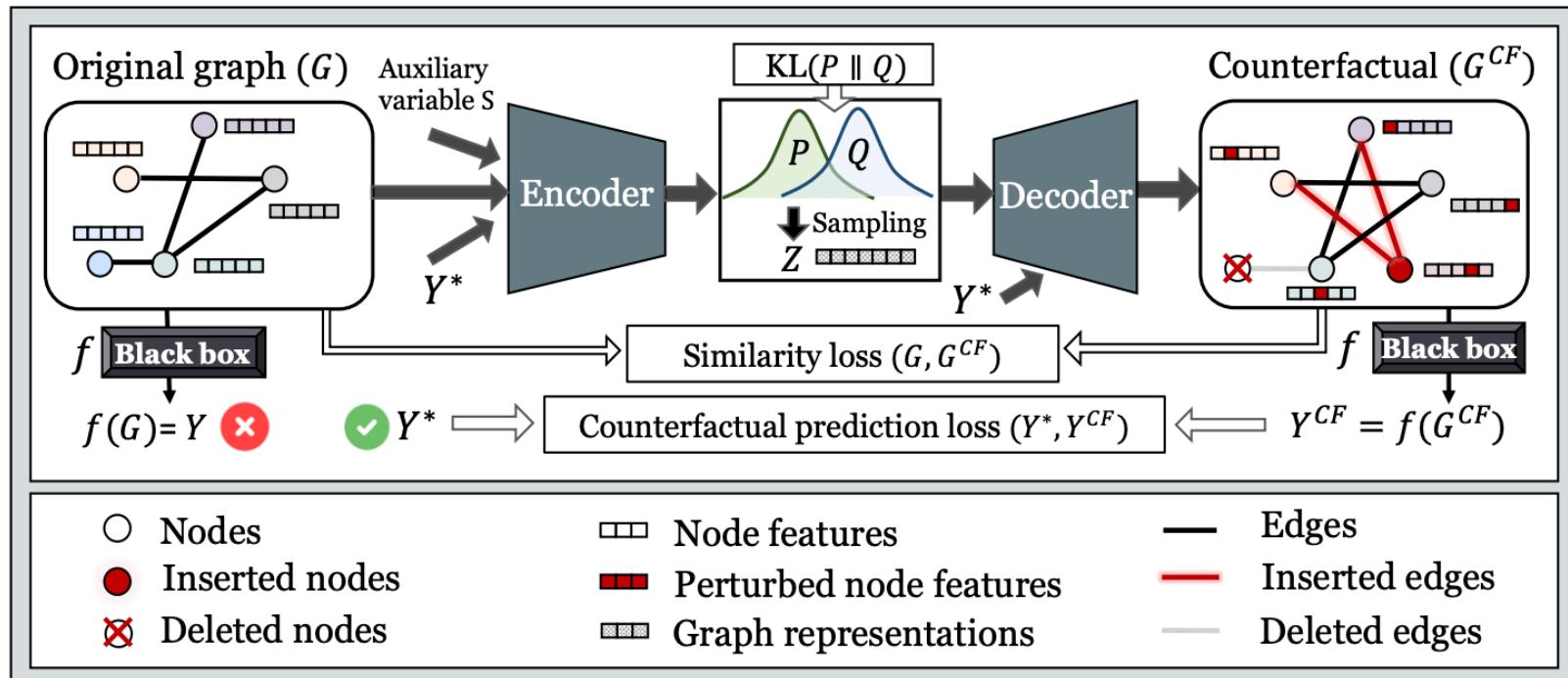
# Challenges of CFE on Graphs

- Optimization
  - The space of perturbation operations on graphs (e.g., add/remove nodes/edges) is discrete, disorganized, and vast
- Generalization
  - Most existing CFE methods on graphs generate counterfactuals for each graph separately, and cannot generalize to unseen graphs
- Causality
  - It is challenging to generate counterfactuals that are consistent with the underlying causality



# CLEAR – CFE Generator for Graphs

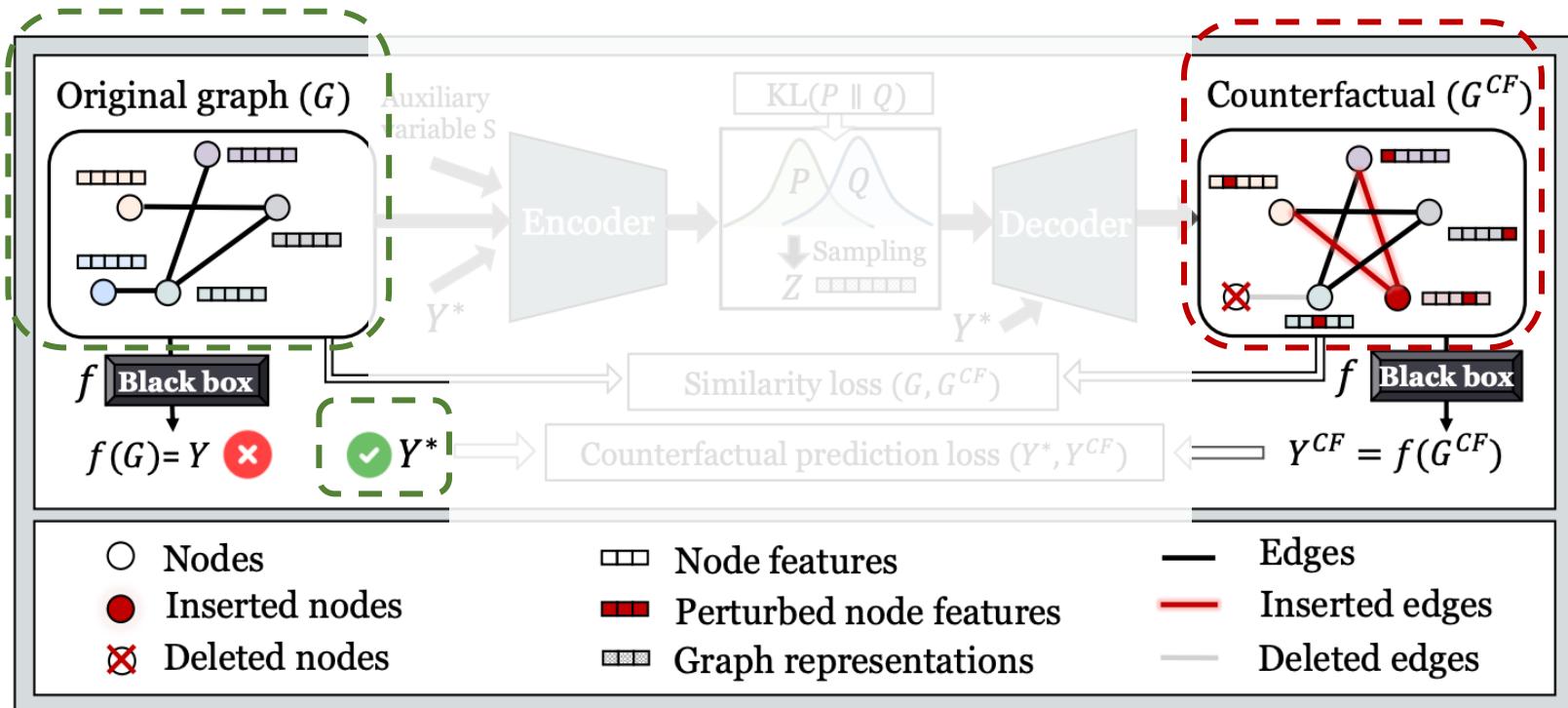
- Backbone CLEAR-VAE enables **optimization** and **generalization** on graph data
- On top of CLEAR-VAE, CLEAR promotes the **causality** of CFEs with an auxiliary variable  $S$



# CLEAR – CFE Generator for Graphs

- CLEAR-VAE:

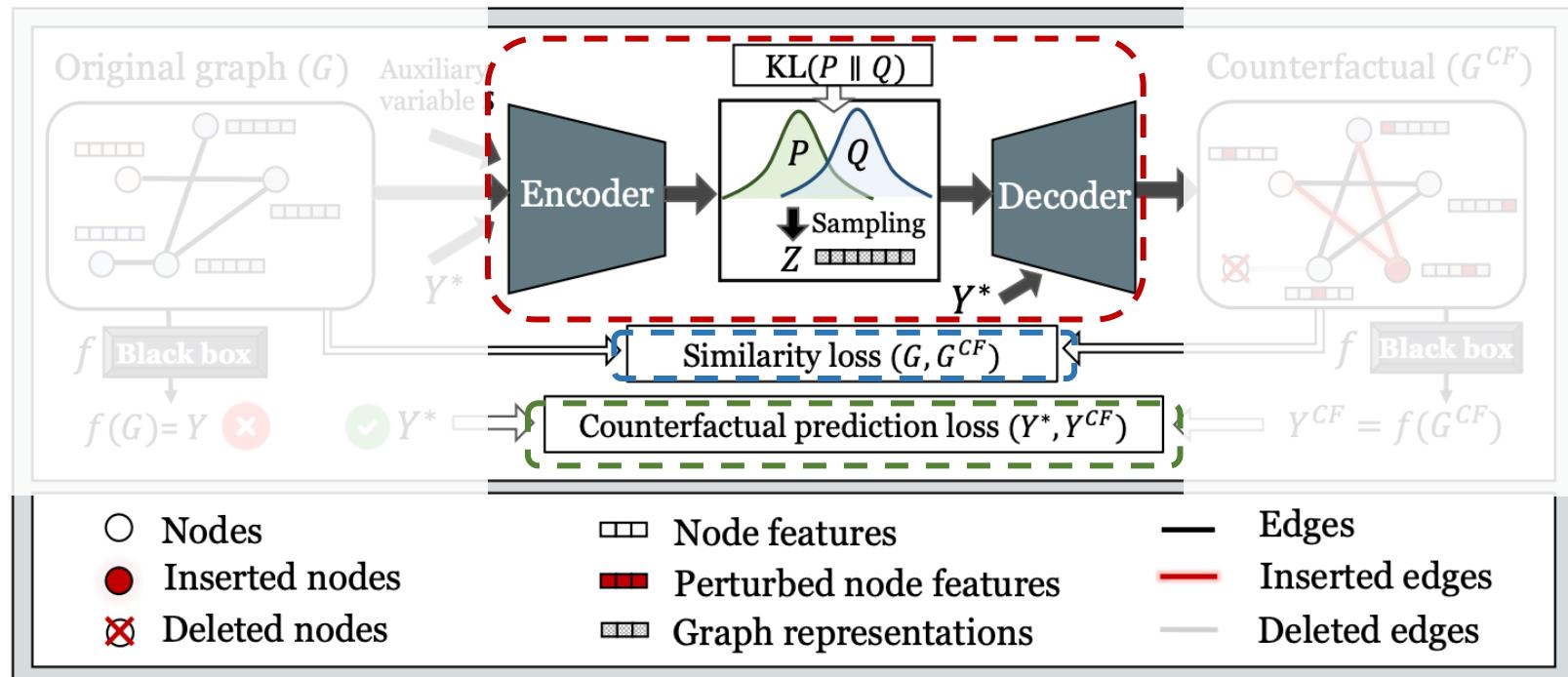
- Input: original graph ( $G$ ), the desired label ( $Y^*$ )
- Output: counterfactual ( $G^{CF}$ )



# CLEAR – CFE Generator for Graphs

- CLEAR-VAE: graph variational autoencoder (VAE) based, with representation  $Z$  learned in the bottleneck layer.
- Loss:

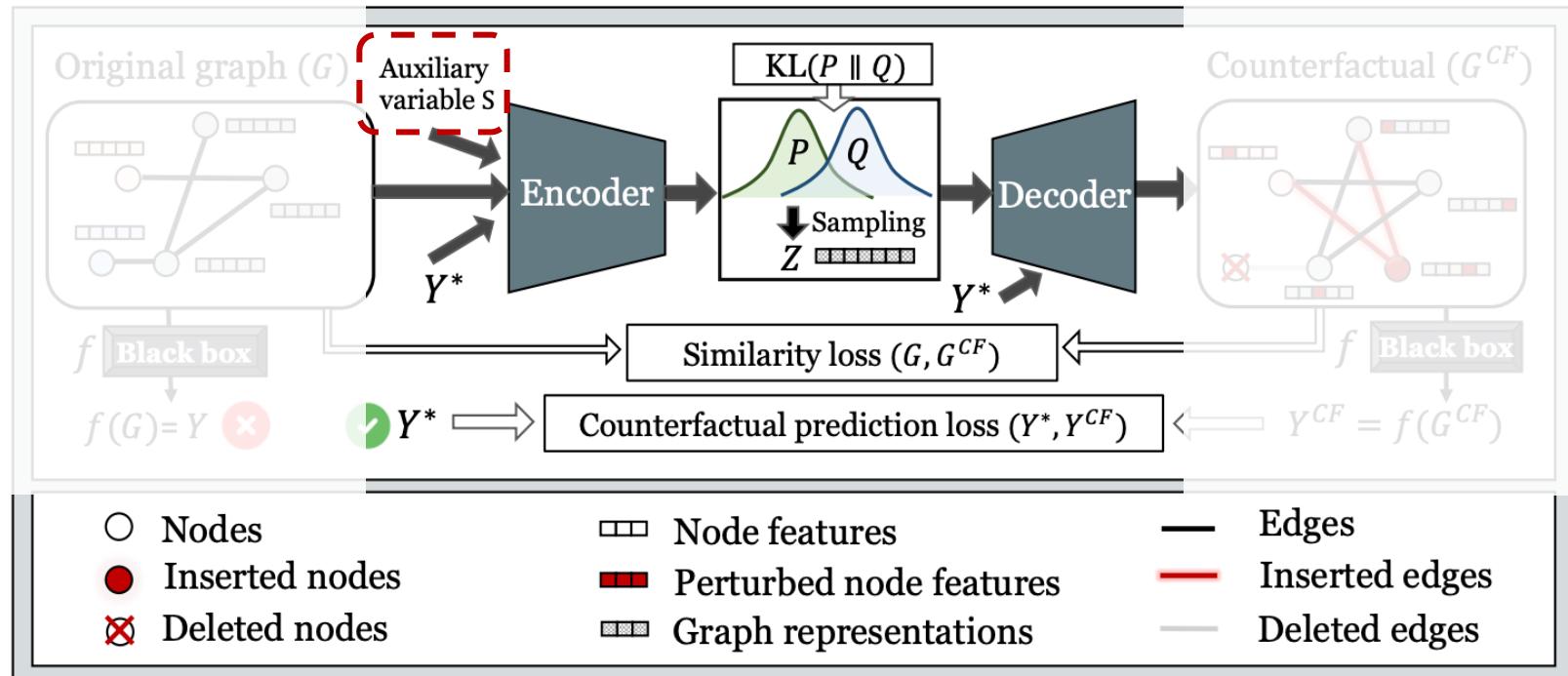
$$\mathcal{L} = \mathbb{E}_Q[d(G, G^{CF})] + \alpha \cdot l(f(G^{CF}), Y^*) + \text{KL}(Q(Z|G, Y^*)\|P(Z|G, Y^*))$$



# CLEAR – CFE Generator for Graphs

- CLEAR: To further promote causality, we leverage an auxiliary  $S$  to better identify the underlying causal model
- Final loss:

$$\mathcal{L} = \mathbb{E}_Q[d(G, G^{CF}) + \alpha \cdot l(f(G^{CF}), Y^*)] + \text{KL}(Q(Z|G, S, Y^*)\|P(Z|G, S, Y^*))$$



# Experiments

- Observations:

- CLEAR achieves good performance in **validity** and **proximity**. (Effective CFE)
  - CLEAR can **time-efficiently** generate counterfactuals for a new graph.
  - CLEAR significantly outperforms all the baselines in **causality**.
- Ratio of counterfactuals which satisfy the causal constraints of interest

Ratio of counterfactuals which achieve the desired outcome

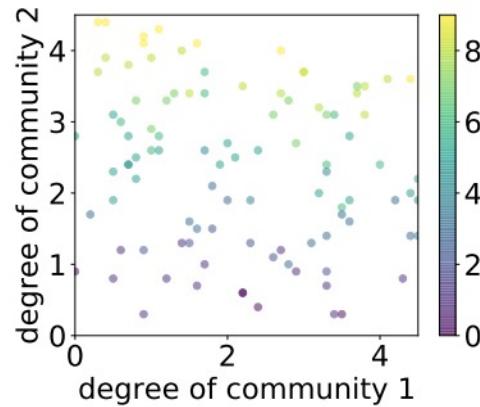
Similarity between original and counterfactual graphs

**Table: The performance of different methods of CFEs on graphs.**

Datasets	Methods	Validity ( $\uparrow$ )	Proximity $X$ ( $\uparrow$ )	Proximity $A$ ( $\uparrow$ )	Causality ( $\uparrow$ )	Time ( $\downarrow$ )
Community	Random	0.53 $\pm$ 0.05	N/A	0.77 $\pm$ 0.02	0.52 $\pm$ 0.06	0.20 $\pm$ 0.01
	EG-IST	0.53 $\pm$ 0.05	N/A	0.66 $\pm$ 0.03	0.13 $\pm$ 0.06	0.27 $\pm$ 0.03
	EG-RMV	0.55 $\pm$ 0.04	N/A	<b>0.85 <math>\pm</math> 0.01</b>	0.03 $\pm$ 0.02	0.15 $\pm$ 0.01
	GNNEExplainer	0.52 $\pm$ 0.06	N/A	0.71 $\pm$ 0.01	0.05 $\pm$ 0.00	2.87 $\pm$ 0.08
	CF-GNNEExplainer	0.90 $\pm$ 0.04	N/A	0.72 $\pm$ 0.00	0.14 $\pm$ 0.02	25.14 $\pm$ 1.22
	MEG	0.88 $\pm$ 0.04	N/A	0.71 $\pm$ 0.01	0.10 $\pm$ 0.03	27.29 $\pm$ 1.32
	<b>CLEAR(ours)</b>	<b>0.94 <math>\pm</math> 0.02</b>	<b>0.91 <math>\pm</math> 0.01</b>	<b>0.77 <math>\pm</math> 0.00</b>	<b>0.65 <math>\pm</math> 0.03</b>	<b>0.01 <math>\pm</math> 0.01</b>
Ogbg-molhiv	Random	0.48 $\pm$ 0.09	N/A	0.87 $\pm$ 0.02	0.46 $\pm$ 0.1	0.17 $\pm$ 0.02
	EG-IST	0.48 $\pm$ 0.09	N/A	0.83 $\pm$ 0.03	0.46 $\pm$ 0.09	0.19 $\pm$ 0.04
	EG-RM	0.483 $\pm$ 0.09	N/A	<b>0.96 <math>\pm</math> 0.01</b>	0.47 $\pm$ 0.09	0.17 $\pm$ 0.04
	GNNEExplainer	0.50 $\pm$ 0.01	N/A	0.92 $\pm$ 0.00	0.48 $\pm$ 0.10	2.78 $\pm$ 0.10
	CF-GNNEExplainer	0.54 $\pm$ 0.02	N/A	0.92 $\pm$ 0.01	0.49 $\pm$ 0.02	27.93 $\pm$ 1.20
	MEG	0.49 $\pm$ 0.03	N/A	0.93 $\pm$ 0.01	0.50 $\pm$ 0.10	22.39 $\pm$ 2.20
	<b>CLEAR(ours)</b>	<b>0.98 <math>\pm</math> 0.01</b>	<b>0.92 <math>\pm</math> 0.02</b>	<b>0.95 <math>\pm</math> 0.01</b>	<b>0.64 <math>\pm</math> 0.02</b>	<b>0.01 <math>\pm</math> 0.00</b>
IMDB-M	Random	0.50 $\pm$ 0.04	N/A	0.67 $\pm$ 0.01	0.43 $\pm$ 0.08	0.19 $\pm$ 0.01
	EG-IST	0.56 $\pm$ 0.12	N/A	0.67 $\pm$ 0.06	0.45 $\pm$ 0.07	0.16 $\pm$ 0.03
	EG-RM	0.45 $\pm$ 0.11	N/A	<b>0.75 <math>\pm</math> 0.03</b>	0.53 $\pm$ 0.08	0.18 $\pm$ 0.02
	GNNEExplainer	0.43 $\pm$ 0.10	N/A	0.62 $\pm$ 0.02	0.50 $\pm$ 0.02	2.46 $\pm$ 0.50
	CF-GNNEExplainer	0.95 $\pm$ 0.02	N/A	0.74 $\pm$ 0.02	0.51 $\pm$ 0.02	22.21 $\pm$ 1.42
	MEG	0.90 $\pm$ 0.02	N/A	0.72 $\pm$ 0.02	0.51 $\pm$ 0.02	24.12 $\pm$ 1.08
	<b>CLEAR(ours)</b>	<b>0.96 <math>\pm</math> 0.01</b>	<b>0.99 <math>\pm</math> 0.00</b>	<b>0.75 <math>\pm</math> 0.01</b>	<b>0.73 <math>\pm</math> 0.01</b>	<b>0.01 <math>\pm</math> 0.00</b>

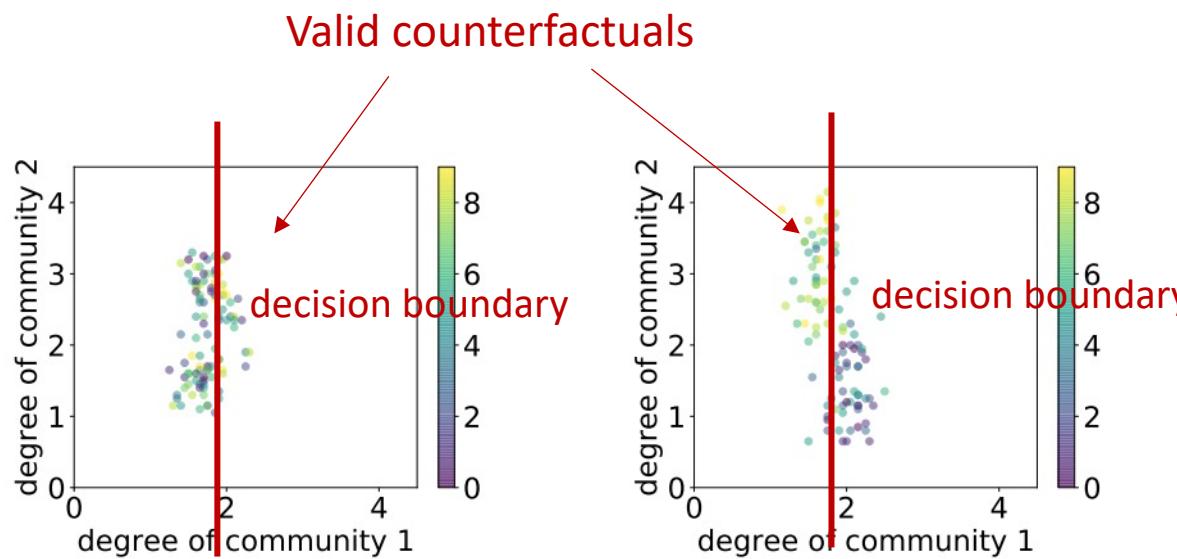
# Experiments

- Observations:
  - Both CLEAR-VAE and CLEAR can generate valid counterfactual



(a) Original data

Colors denote  
different values of  $S$



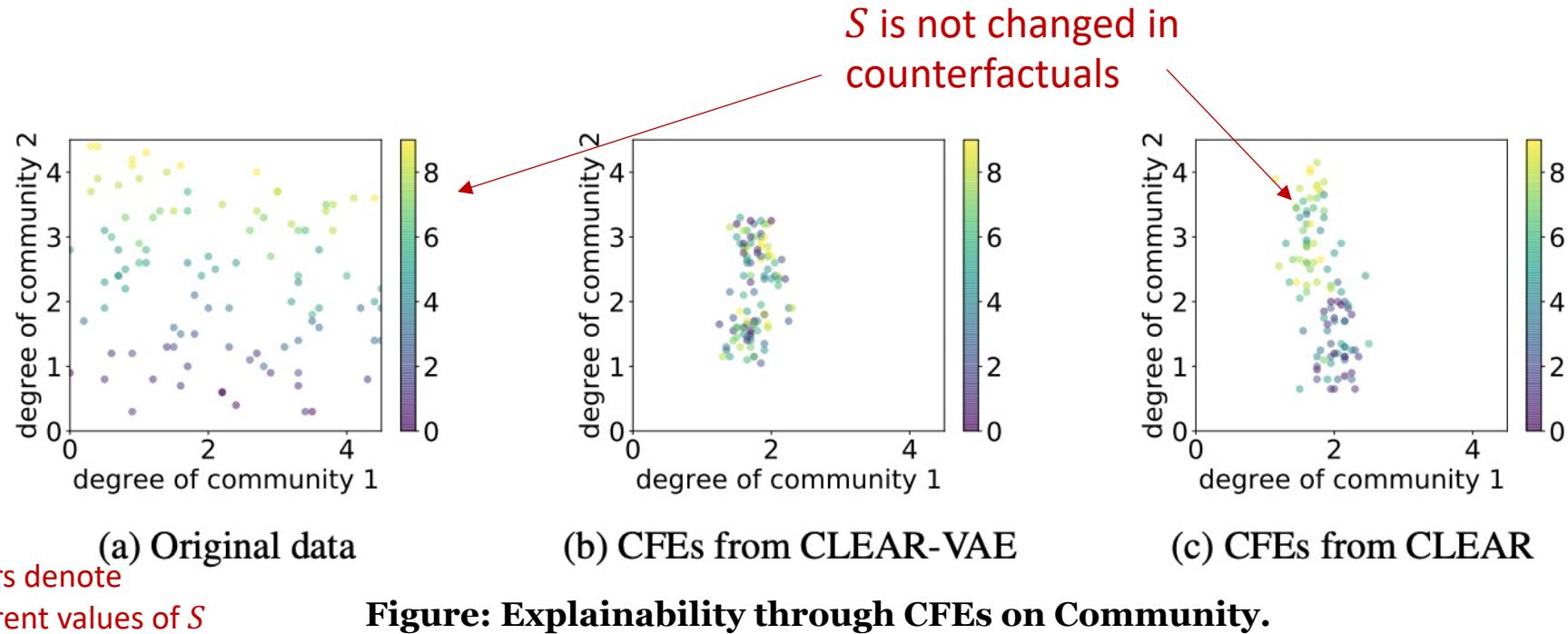
(b) CFEs from CLEAR-VAE

(c) CFEs from CLEAR

**Figure: Explainability through CFEs on Community.**

# Experiments

- Observations:
  - Both CLEAR-VAE and CLEAR can generate valid counterfactual
  - Compared with CLEAR-VAE, the counterfactuals generated by CLEAR are more consistent with the underlying causal model

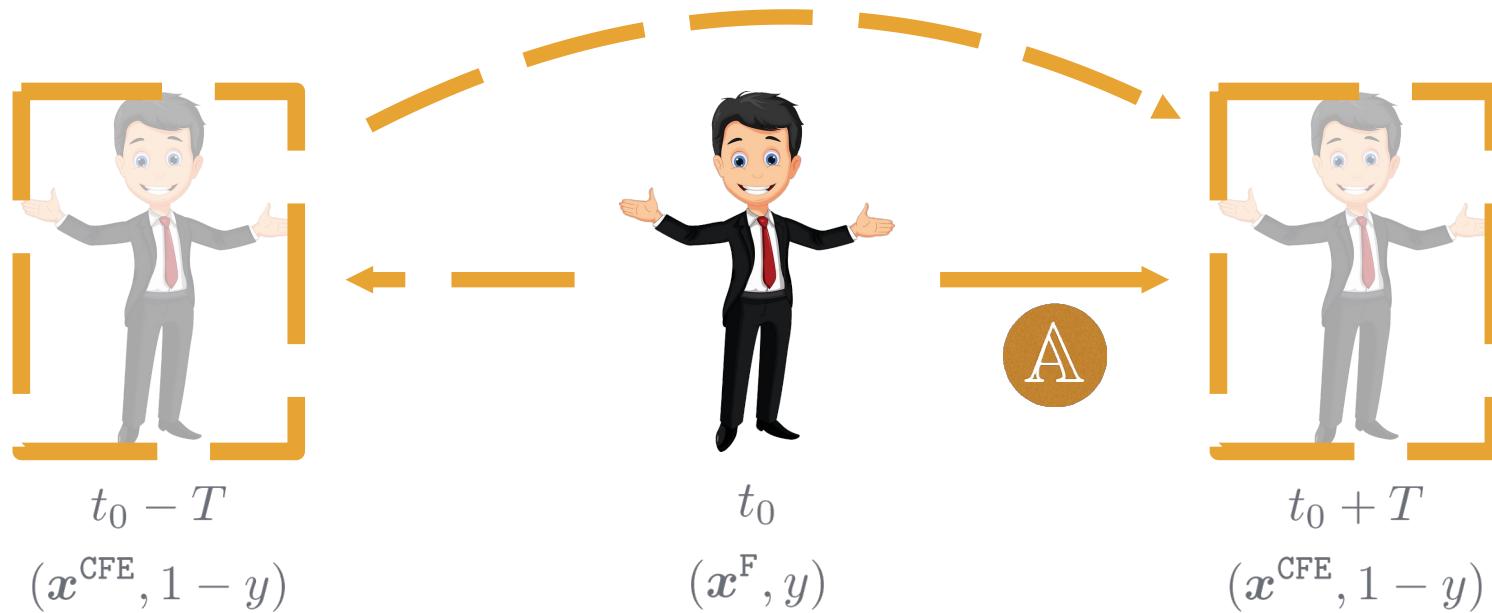


# Outline

- Causal explanation: general introduction
- Decision based interpretability
  - Counterfactual explanation
  - Recourse
- Data-based interpretability
- Model-based interpretability (attribution)

# Algorithmic Recourse

One of the primary objectives of “**explanations** as a means to help a data-subject **act** rather than merely **understand.**” [Wachter et al., 2017]



Algorithmic  
Recourse

the **actions** required for, or “the systematic process of reversing unfavorable decisions by algorithms & bureaucracies across a range of counterfactual scenarios” [Venkata-subramanian & Alfano, 2020]

# CFE-based Recourse

Counterfactual Explanation

$$\begin{aligned} \boldsymbol{x}_*^{\text{CFE}} &\in \operatorname{argmin}_{\boldsymbol{x}} \quad \text{dist}(\boldsymbol{x}, \boldsymbol{x}^F) \\ \text{s.t. } h_{\theta}(\boldsymbol{x}) &\neq h_{\theta}(\boldsymbol{x}^F) \\ \boldsymbol{x} &\in \mathcal{P}\text{lausible} \end{aligned}$$

[Wachter et al., 2017]

CFE-based Recourse

$$\begin{aligned} \boldsymbol{\delta}^* &\in \operatorname{argmin}_{\boldsymbol{\delta}} \quad \text{cost}(\boldsymbol{\delta}; \boldsymbol{x}^F) \\ \text{s.t. } h_{\theta}(\boldsymbol{x}^{\text{CFE}}) &\neq h_{\theta}(\boldsymbol{x}^F) \\ \boldsymbol{x}^{\text{CFE}} &= \boldsymbol{x}^F + \boldsymbol{\delta} \\ \boldsymbol{x}^{\text{CFE}} &\in \mathcal{P}\text{lausible} \\ \boldsymbol{\delta} &\in \mathcal{F}\text{easible} \end{aligned}$$

[Ustun et al., 2019]

Question: Do CFE-based recourse actions translate to **optimal** and **feasible** real-world actions for recourse?

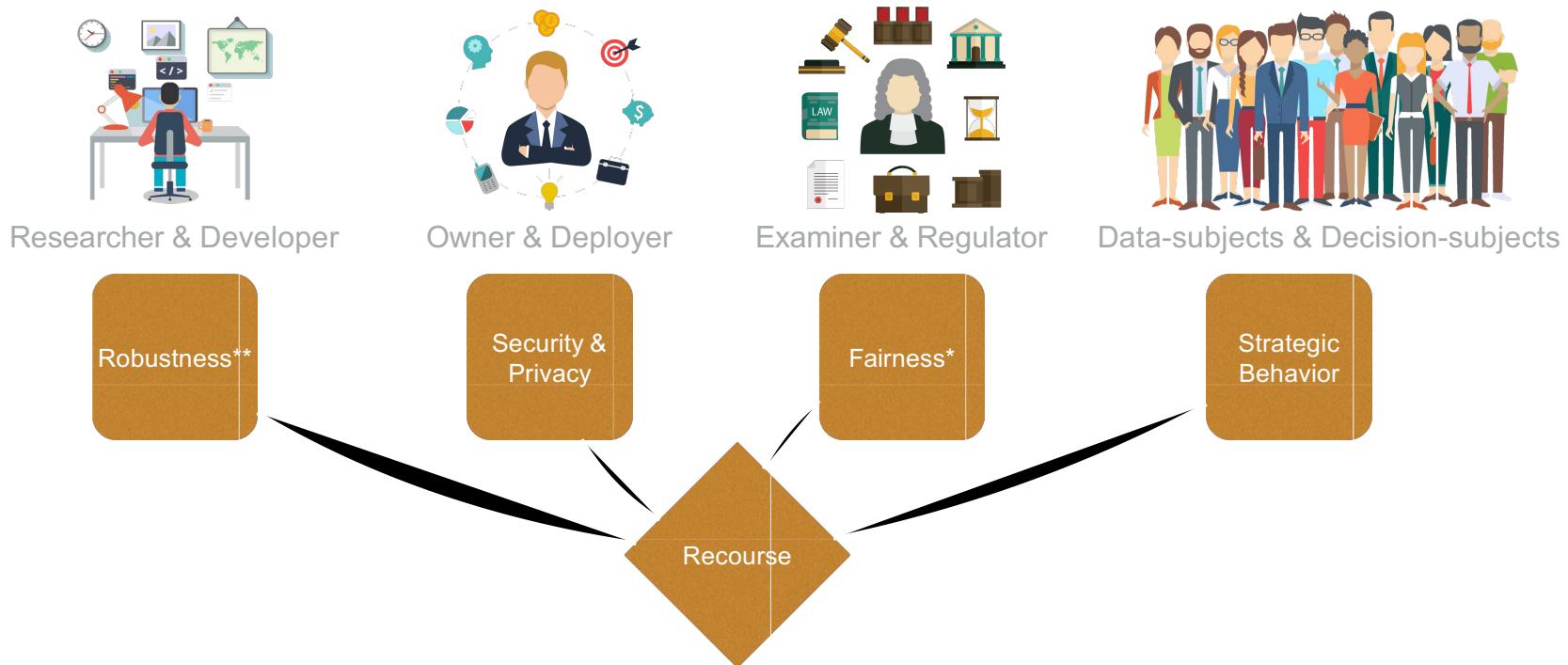
# Sub-optimal Recourse

- Issue: Optimal changes may not be feasible!
  - e.g., “decrease your age”



Question: Do CFE-based recourse actions translate to **optimal** and **feasible** real-world actions for recourse?

# Recourse & Other Ethical Desiderata



\*[JvK, Karimi, et al., AAAI (oral), 2022]

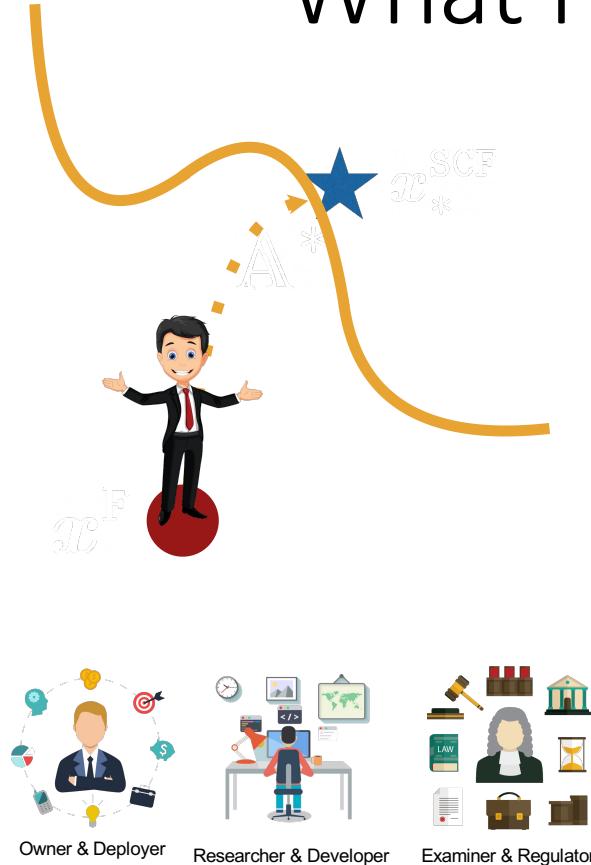
\*\*[RDO, Karimi, et al., ICML (spotlight), 2022]

# Recourse Summary

- In general, counterfactual explanations **do not, imply feasible actions.**
- In general, **recourse can only be guaranteed under perfect causal knowledge.**
- **Diversity of stakeholder needs** illustrates a tension between many desirable system properties.
- Trade-offs between different needs may require new techniques, perhaps on the **cross-disciplinary expertise.**

Algorithm	Formulation								Solution								
	Goal	Model	Actionability		Plausibility		Extra	Data types	Tools	Access	Properties						
			TB	KB	DF	OT	uncond.	cond.	dom.	dens.	proto.	diver.	spar.	grid	camera	code	
(2014.03) SEDC [119]	E	●	●	●	●	●	○	○	○	○	○	●	heuristic	query	●	●	●
(2015.08) OAE [48]	E	●	●	●	●	●	●	●	●	●	●	●	ILP	white-box	●	●	●
(2016.05) HCLS [101, 103]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt/heuristic	gradient/query	●	●	●
(2017.06) Feature Tweaking [169]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	white-box	○	○	○
(2017.11) CF Expl. [177]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	○	○	○
(2017.12) Growing Spheres [105]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	○	○	○
(2018.02) CEM [52]	E	●	●	●	●	●	●	●	●	●	●	●	FISTA	class prob.	●	●	●
(2018.02) POLARIS [188]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	gradient	●	●	●
(2018.05) LORE [74]	E	●	●	●	●	●	●	●	●	●	●	●	gen alg + heuristic	query	●	●	●
(2018.06) Local Foil Trees [172]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2018.09) Actionable Recourse [171]	E	●	●	●	●	●	●	●	●	●	●	●	ILP	white-box	●	●	●
(2018.11) Weighted CFs [71]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2019.01) Efficient Search [158]	E	●	●	●	●	●	●	●	●	●	●	●	MILP	white-box	●	●	●
(2019.04) CF Visual Expl. [70]	E	●	●	●	●	●	●	●	●	●	●	●	greedy search	white-box	●	●	●
(2019.05) MACE [91]	E	●	●	●	●	●	●	●	●	●	●	●	SAT	white-box	●	●	●
(2019.05) DICE [32]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2019.05) CERTIFIAl [162]	E	●	●	●	●	●	●	●	●	●	●	●	gen alg	query	●	●	●
(2019.06) MACEM [53]	E	●	●	●	●	●	●	●	●	●	●	●	FISTA	query	●	●	●
(2019.06) Expl. using SHAP [152]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2019.07) Nearest Observable [181]	E	●	●	●	●	●	●	●	●	●	●	●	brute force	dataset	●	●	●
(2019.07) Guided Prototypes [173]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt/FISTA	gradient/query	●	●	●
(2019.07) REVISE [87]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2019.08) CLEAR [182]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2019.08) MC-BRP [113]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2019.09) FACE [149]	E	●	●	●	●	●	●	●	●	●	●	●	graph + heuristic	query	●	●	●
(2019.10) Action Sequences [150]	E	●	●	●	●	●	●	●	●	●	●	●	program synthesis	class prob.	●	●	●
(2019.10) C-CHVAE [143]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt + heuristic	query + gradient	●	●	●
(2019.11) OCE [114]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt + heuristic	white-box	●	●	●
(2019.12) Model-based CFs [117]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2019.12) LIME-C/SHAP-C [151]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2019.12) EMAP [40]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	dataset/query	●	●	●
(2019.12) PRINCE [65]	E	●	●	●	●	●	●	●	●	●	●	●	graph + heuristic	query	●	●	●
(2019.12) LowProFool [18]	E	○	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.01) ABELE [73]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.02) CEMA [11–13]	E	●	●	●	●	●	●	●	●	●	●	●	gen alg + heuristic	query + data	●	●	●
(2020.02) MINT [92]	R	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.03) VICE [68]	E	●	●	●	●	●	●	●	●	●	●	●	gen alg + heuristic	query + data	●	●	●
(2020.03) Plausible CFs [22]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt/heuristic	gradient/query	●	●	●
(2020.04) SEDC-T [175]	E	●	●	●	●	●	●	●	●	●	●	●	SAT	white-box	●	●	●
(2020.04) MOC [49]	E	○	●	●	●	●	●	●	●	●	●	●	grad opt + gen alg	dataset	●	●	●
(2020.04) SCOUT [179]	E	●	●	●	●	●	●	●	●	●	●	●	gen alg	query	●	●	●
(2020.04) ASP-based CFs [28]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.05) CBR-based CFs [95]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.06) Survival Model CFs [97]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query + data	●	●	●
(2020.06) Probabilistic Recourse [93]	R	●	●	●	●	●	●	●	●	●	●	●	gen alg	query	●	●	●
(2020.06) C-CHVAE [142]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt/brute force	gradient/query	●	●	●
(2020.07) FRACE [189]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.07) DACE [88]	E	●	●	●	●	●	●	●	●	●	●	●	MILP	white-box	●	●	●
(2020.07) CRUDS [56]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient/data	●	●	●
(2020.07) Gradient Boosted-based CFs [5]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	data	●	●	●
(2020.08) Gradual Construction [89]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	class prob.	●	●	●
(2020.08) DECE [42]	E	●	●	●	●	●	●	●	●	●	●	●	grad opt	gradient	●	●	●
(2020.08) Time Series CFs [16]	E	●	●	●	●	●	●	●	●	●	●	●	heuristic	query	●	●	●
(2020.08) PermuteAttack [90]	E	●	●	●	●	●	●	●	●	●	●	●	gen alg	query	●	●	●

# What remains to be done?



## **Underlying model:**

classification / regression —> matching / scheduling?

## **Variables:**

meaningful concepts —> learn / abstract concepts?

## **Stakeholders:**

is recourse enough? —> fair, robust, private recourse?

## **Interaction:**

single-step recommendation —> multi-step interventions?

## **World knowledge:**

implausible or costly actions —> collaborative discussion?

# Outline

- Causal explanation: general introduction
- Decision based interpretability
  - Counterfactual explanation
  - Recourse
- Data-based interpretability (Discovery)
- Model-based interpretability (attribution)

# Explaining a car to a 4-year-old



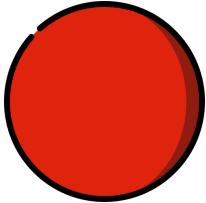
## What's a car?

A car is a machine with **wheels** that can take you from one place to another. It usually has seats for people to sit in and a steering wheel for the driver to control where the car goes. [...] People use cars to go to work, go on vacations, or visit their friends and family.



## What's a “wheel”?

Wheels are **circular**, flat objects that are attached to the bottom of a car. They help the car move by rolling along the ground. Think of a wheel like a **circular** plate that is thick in the middle and thinner on the edges. The edge of the wheel is called the "rim," and it's usually made of [...]



## What's a “circle”?

A circle is a **shape** that is round like a ball or a tire. It has no corners or edges, and all points on its edge are the same distance away from the center. [...] Circles are all around us, from the **sun** and **moon** in the sky to the wheels on a car or **bicycle**.

# Concepts

Explanations are effective when presented in a *language* consistent with that of the explainee

Language is built atop a set of *mutually understandable concepts*

Concepts subsume

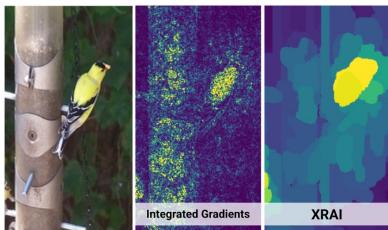
- tone
- colors
- shapes
- Patterns
- processes & demonstrations



good explanations are built atop good concepts

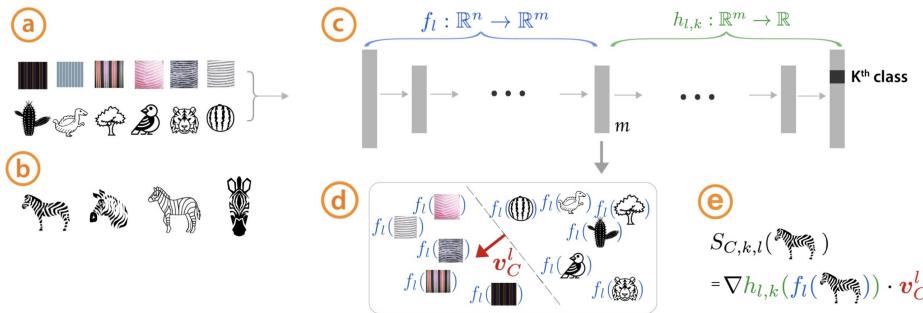
# The *explanation-concept alignment* spectrum

Concepts *unclear*



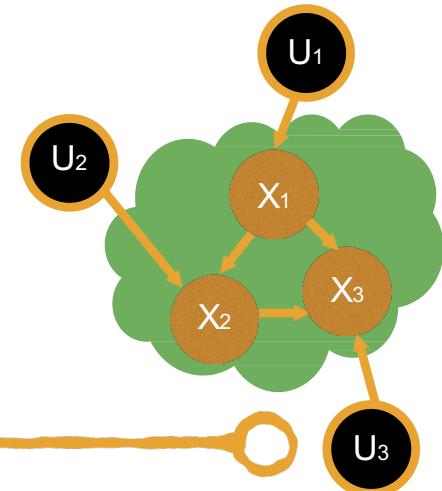
Attribution-based Explanations

Concepts “*extracted*”



Testing with Concept Activation Vectors  
[Kim et al., ICML, 2018]

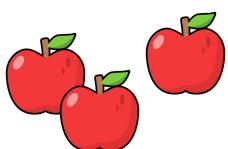
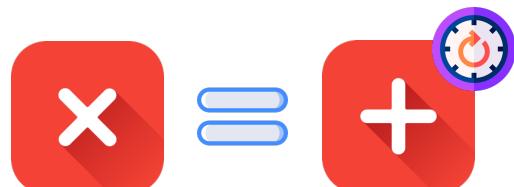
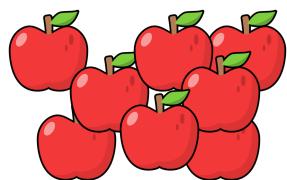
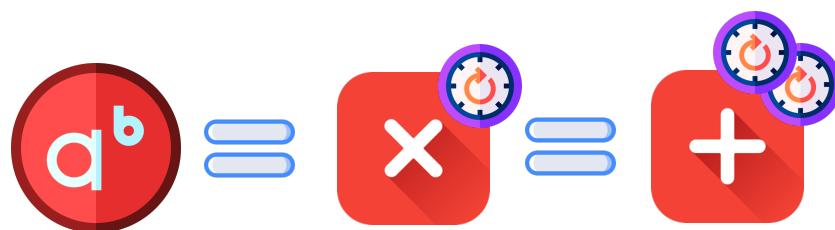
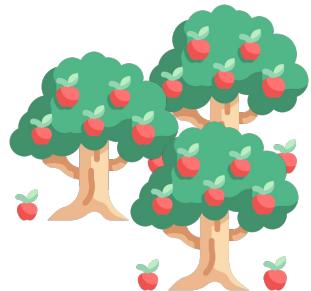
Concepts *clear*



Recourse-based Explanations

Ideally, an AI will expand our concept bank

# Learning concepts from first principles

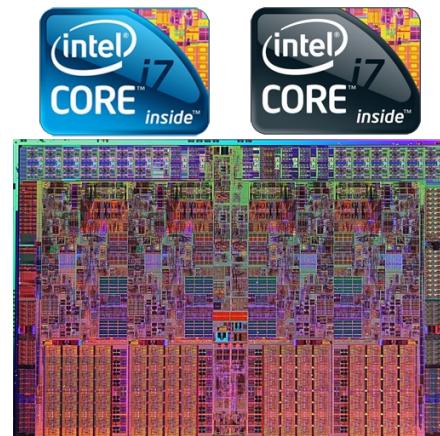
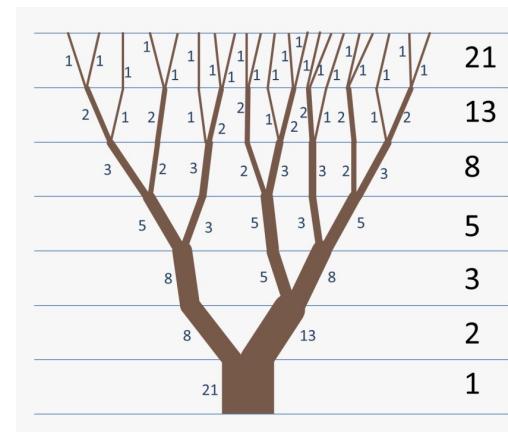
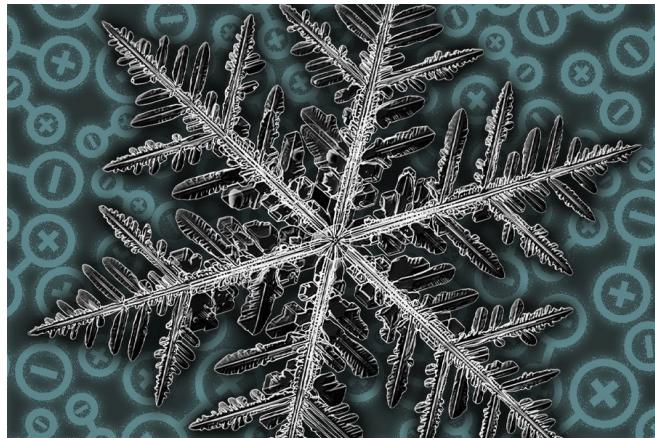


```
def pow(x: int, y: int) -> int:  
    out = 1  
    for _ in range(y)  
        out = mult(out, x)  
    return out
```

```
def mult(x: int, y: int) -> int:  
    out = 0  
    for _ in range(y)  
        out = add(out, x)  
    return out
```

```
def add(x: int, y: int) -> int:  
    return x + y
```

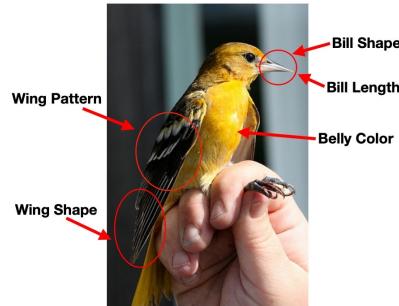
# Abstractions



# Causal Discovery for XAI

High-dimensional set  
of “low-level” features  
(e.g., pixels)

$$X = (X_1, X_2, \dots, X_q)$$



Low-dim set of  
“high-level” features  
(human interpretable)

$$Z = (Z_1, Z_2, \dots, Z_p)$$

Fig. 1: An image of a Baltimore Oriole annotated with interpretable features.

Assume dataset with  $X$  covariates,  $Y$  targets, and  $Z$  annotations

- Train a classifier  $f : X \rightarrow Y$
- Estimate a causal PAG  $G$  over  $V = (Z, Y = f(X))$  ( $q \gg p$ )
- Determine which high-level features  $Z$  are
  - causes
  - potential causes
  - non-causes of  $Y$

PAG Edge	Meaning
$Z_i \rightarrow \hat{Y}$	$Z_i$ is a cause of $\hat{Y}$
$Z_i \leftrightarrow \hat{Y}$	$Z_i$ and $\hat{Y}$ share an unmeasured common cause $Z_i \leftarrow U \rightarrow \hat{Y}$
$Z_i \circ \rightarrow \hat{Y}$	Either $Z_i$ is a cause of $\hat{Y}$ or there is unmeasured confounding, or both

[Sani et al., 2020]

# Sample Experiment

Each  $d \times d$  image may contain

- Horizontal bar (H)
- Vertical bar (V)
- Circle (C)
- Rectangle (R)
- Random pixel noise

$$Y = V \vee C$$

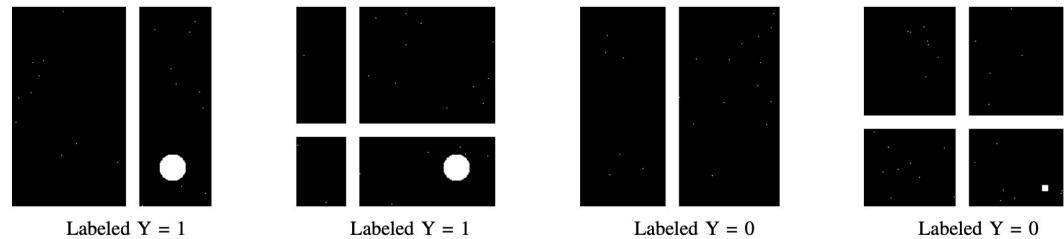


Fig. 3: Simulated image examples with horizontal bars, vertical bars, circles, and rectangles.

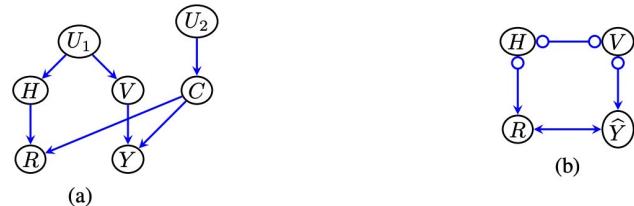
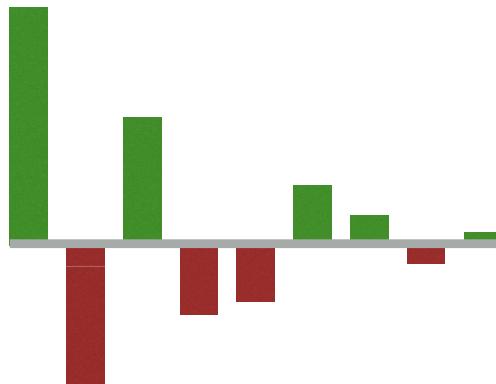


Fig. 4: (a) A causal diagram representing the true data generating process. (b) The PAG learned using FCI with output  $\hat{Y}$  from a convolutional neural network.

# Outline

- Causal explanation: general introduction
- Decision based interpretability
  - Counterfactual explanation
  - Recourse
- Data-based interpretability (Discovery)
- Model-based interpretability (attribution)



# Attribution, causal or not

## Non-causal Attribution

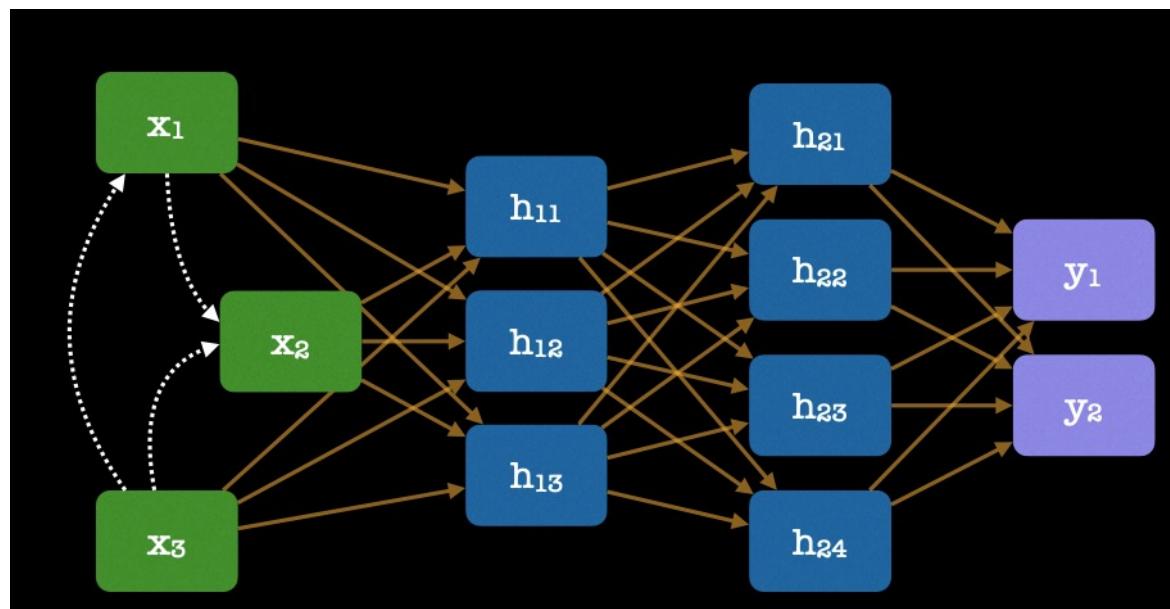
How much would **perturbing** a particular input neuron affect a particular output neuron?

$$E[Y_k | X_i = x_i] - E[Y_k | X_i = x_j]$$

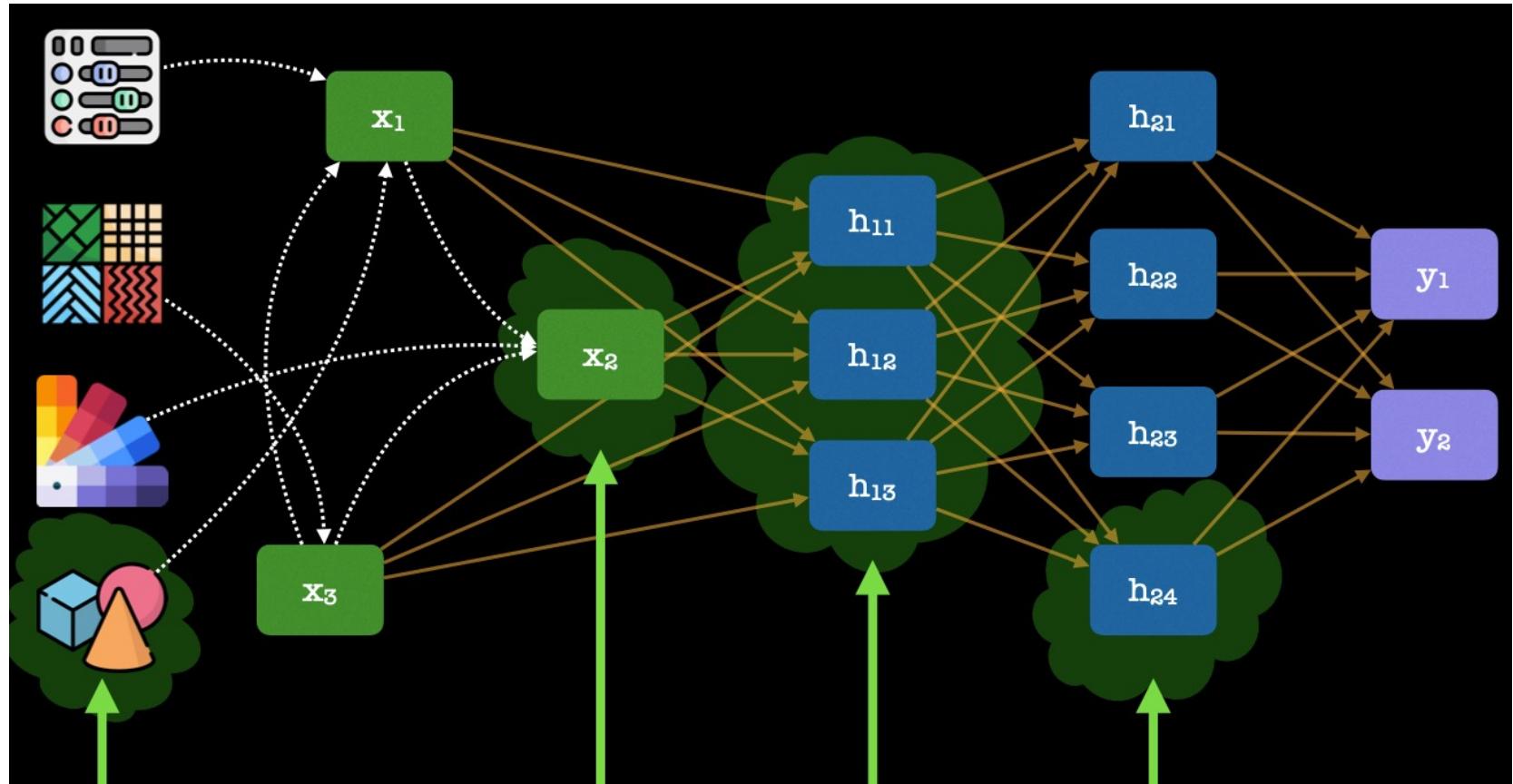
## Causal Attribution

What is the **causal effect** of a particular input neuron on a particular output neuron?

$$E[Y_k | \text{do}(X_i = x_i)] - E[Y_k | \text{do}(X_i = x_j)]$$



$$E[Y_k | \text{do}(X_i = x_i)] - E[Y_k | \text{do}(X_i = x_j)]$$

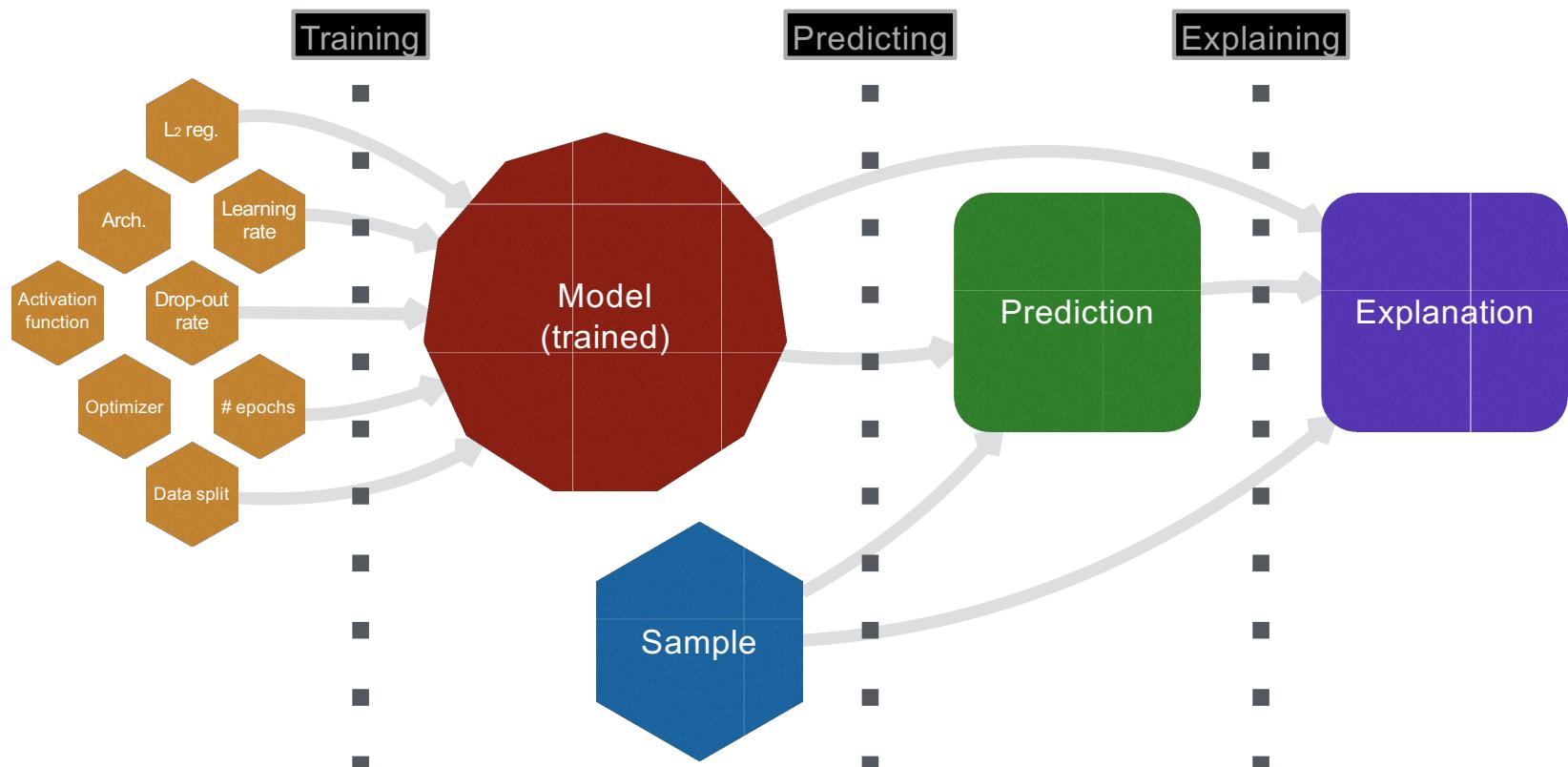




# Evaluating explanations

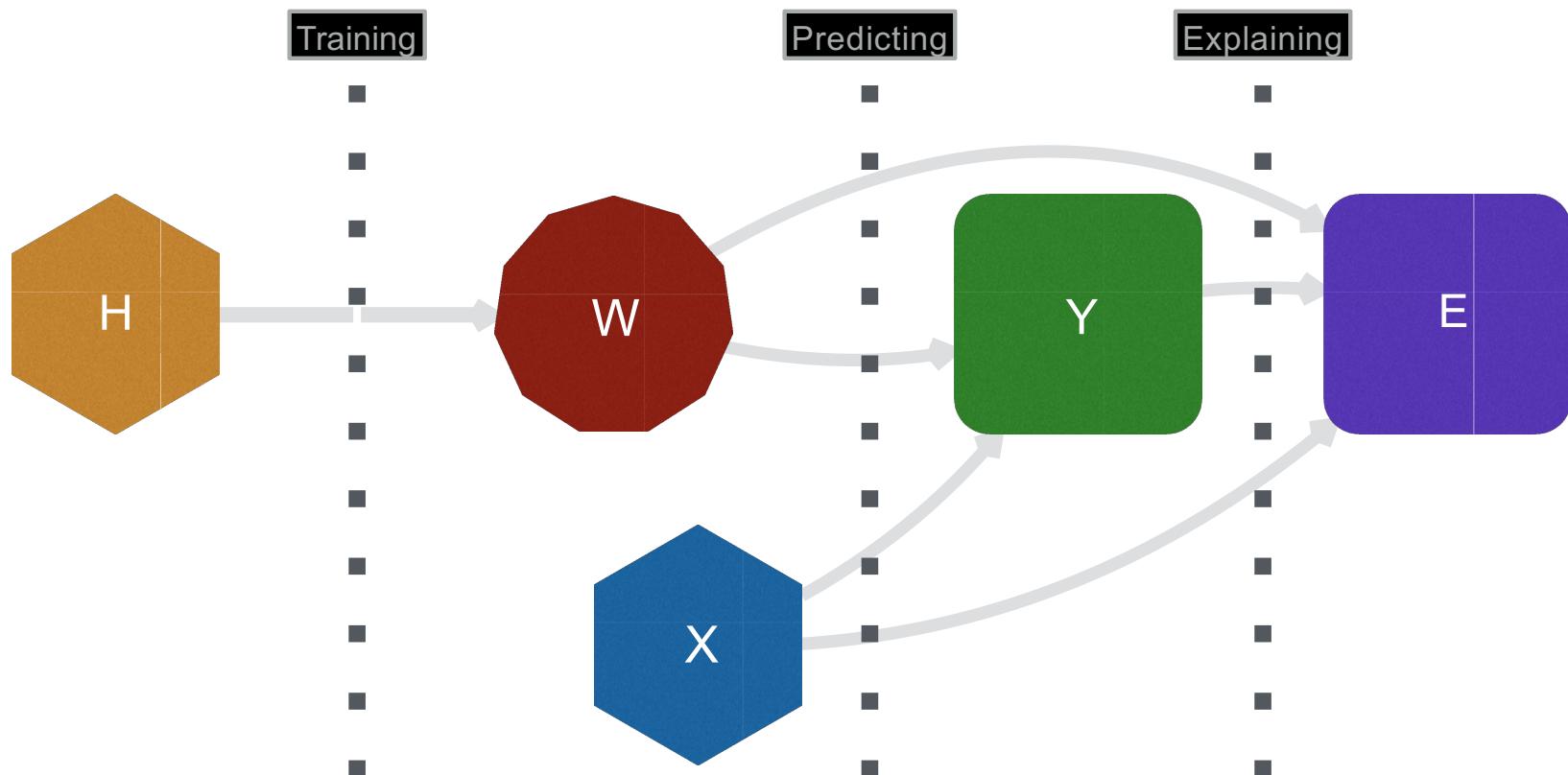


# Motivation: Explanation Generation Process

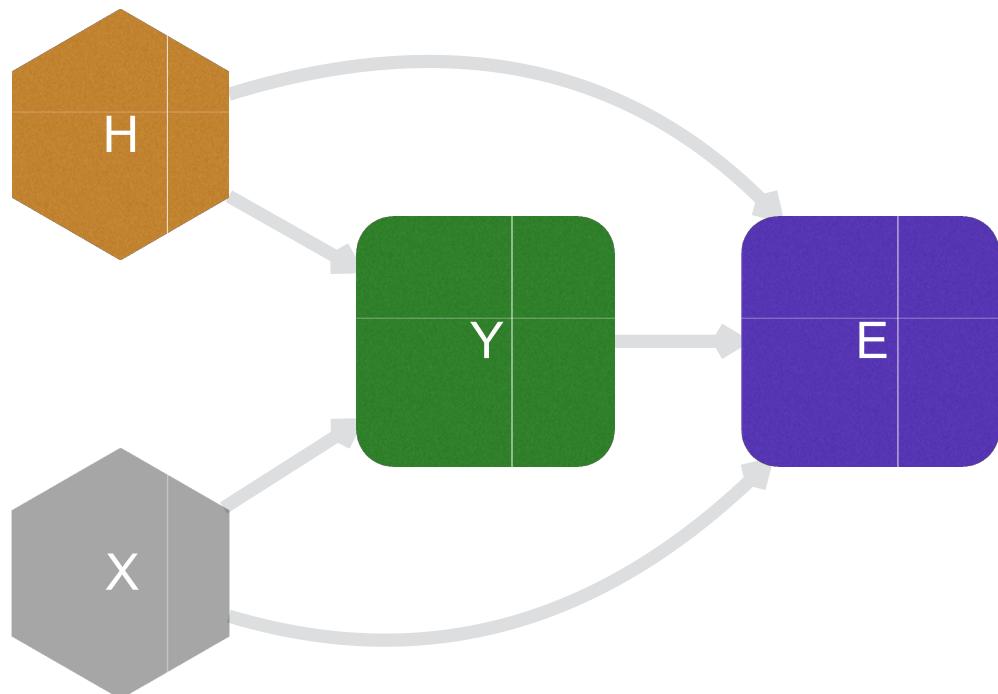


[Karimi et al., ICML 2023]

# Explanation Generation Process

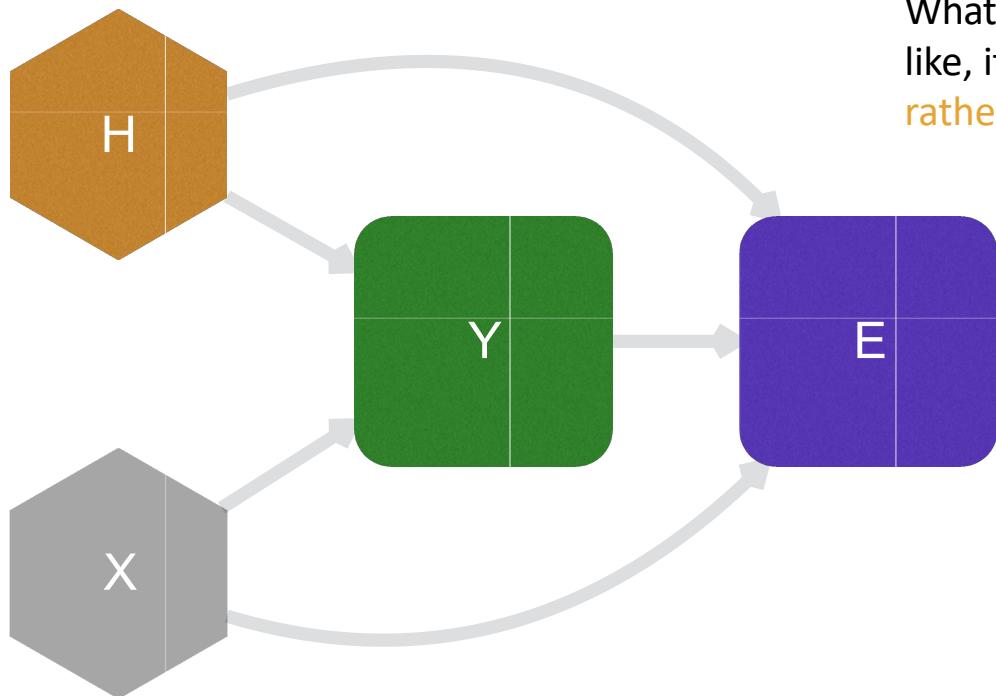


# Hyperparameters as Treatments



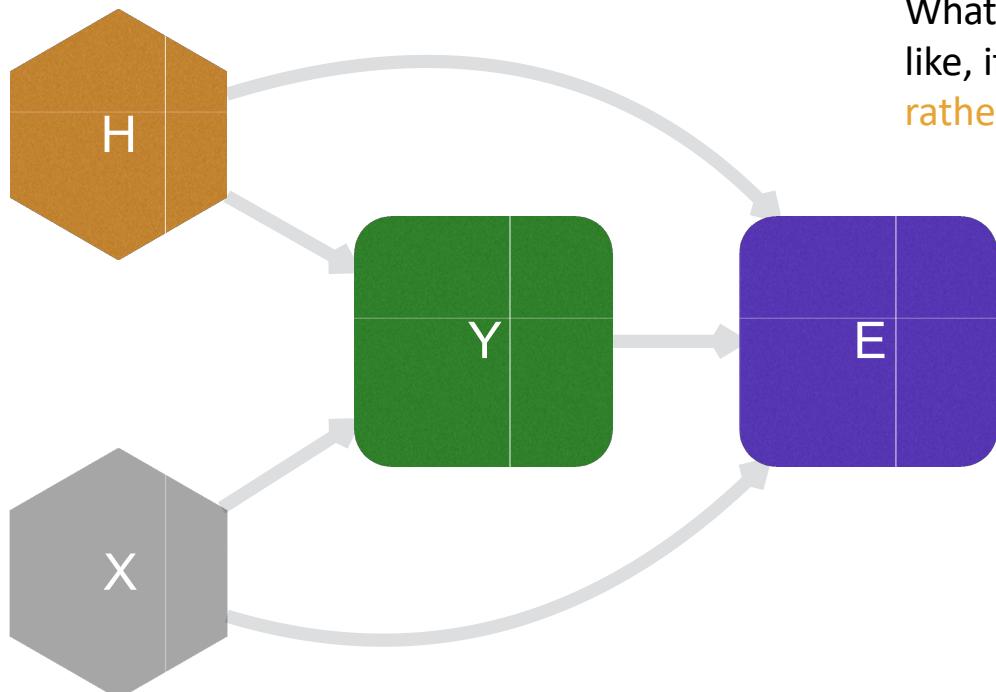
What is the effect of the **hyperparameters** on the resulting prediction/explanation?

# Hyperparameters as Treatments



What does the prediction/explanation for  $X = x$  look like, if the hyperparameters take on value  $H = h$  rather than  $H = h'$ , all else being equal?

# Hyperparameters as Treatments



What does the prediction/explanation for  $X = x$  look like, if the hyperparameters take on value  $H = h$  rather than  $H = h'$ , all else being equal?

$$Y_{h=1} - Y_{h=0}$$

single binary treatment

$$E_{m \neq n} [ Y_{h=n} - Y_{h=m} ]$$

single non-binary treatment

$$E_{h \setminus i} [ E_{m \neq n} [ Y_{h_i=n, h \setminus i} - Y_{h_i=m, h \setminus i} ] ]$$

multiple non-binary treatment

# What remains to be done?

- We need better explainability methods.
- Studied hyperparameters independently; how to study joint effects?
- Analysis limited to the model zoo; how to analyze when zoo not available?
- Assumed linear differences in total and direct effects; what about non-linear effects?
- Extensions beyond saliency map, e.g., SHAP, LIME, CFE, etc.

# References & Reading

- Verma S, Boonsanong V, Hoang M, et al. Counterfactual explanations and algorithmic recourses for machine learning: A review[J]. arXiv
- Lucic A, ter Hoeve M, Tolomei G, et al. CF-GNNExplainer: Counterfactual Explanations for Graph Neural Networks[J]. arXiv 2021.
- Ma, J, et al. “CLEAR: Generative Counterfactual Explanations on Graphs.” *NeurIPS* (2022).
- Karimi A H, Von Kügelgen J, Schölkopf B, et al. Algorithmic recourse under imperfect causal knowledge: a probabilistic approach[J]. *Advances in neural information processing systems*, 2020, 33: 265-277.
- Socially Responsible Machine Learning: A Causal Perspective. KDD 2023 Tutorial.

Thank you!  
Questions?