# DiConStruct

## Causal Concept-based Explanations through Black-Box Model Distillation

Ricardo Moreira, Jacopo Bono, Mário Cardoso, Pedro Saleiro, Mário Figueiredo, Pedro Bizarro

**XAI Seminar - Imperial College London**

**20 June 2024**

feedzai    TÉCNICO LISBOA

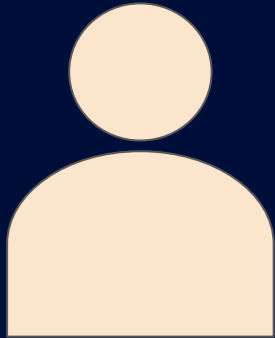# Agenda

Motivation

Methods

Experimental Setup

Results

Conclusions

# **Motivation**

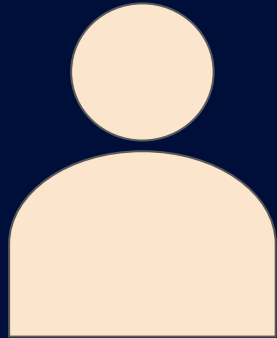# Motivation
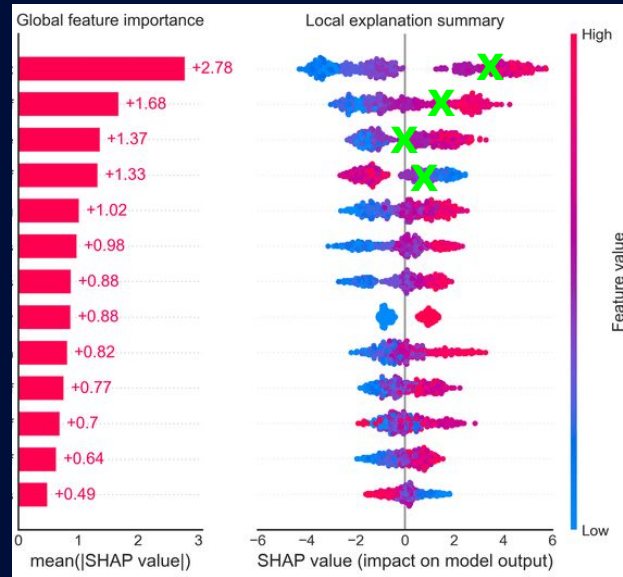
Heart attack?

# Motivation



Heart attack?

...

# Motivation

P(Heart attack) = 0.8



LDL
BMI
glucose
Avg h active / week
...

SHAP: Lundberg et al., NeurIps 2017

# Motivation

1. Feature-based explanations are often **difficult to interpret.**

# Motivation

1. Feature-based explanations are often **difficult to interpret.**

Especially
- if many features
- if semantics / connection to higher level concepts is not obvious

Important when human - AI collaboration is time-sensitive!

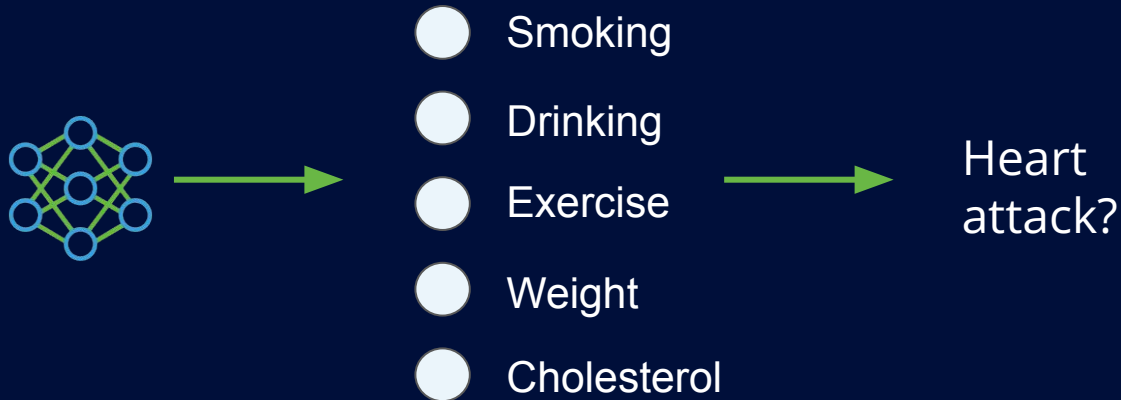# Motivation

1. Feature-based explanations are often **difficult to interpret.**

   Use human understandable concepts instead.

# Motivation

For example, concept bottleneck models (CBM)*

Smoking

Drinking

Exercise

Heart attack?

Weight

Cholesterol

* Koh et al., ICML 2020

# Motivation

1. Feature-based explanations are often **difficult to interpret.**
2. CBMs are self-explainable, but **trade-off** between the main task and the explanation task.

# Motivation

1. Feature-based explanations are often **difficult to interpret.**

2. CBMs are self-explainable, but **trade-off** between the main task and the explanation task.

3. No rigorous **counterfactual reasoning** possible.

   ("What if I would stop smoking?")

# Motivation

1. Feature-based explanations are often **difficult to interpret.**
2. CBMs are self-explainable, but **trade-off** between the main task and the explanation task.
3. No rigorous counterfactual reasoning possible.

-> Use **post hoc** explanations

-> Incorporate **causal** principles

# Motivation

## Causal diagrams

- Causal relations are represented in a **DAG**.
- **Nodes** represent (endogenous) **variables**.
- **Directed edges** represent **causal relationships**.

# Motivation

**Structural Causal models (SCMs)**

- **Exogenous** variables: effects from **outside** the model
- **Endogenous** variables: determined **within** the model
- **Structural equations**: express the relationship between the variables mathematically.
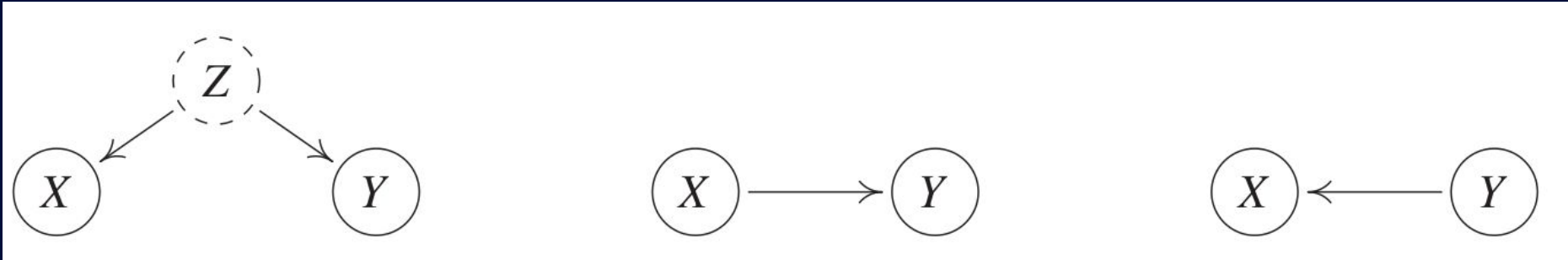


$$X^i := N_X^{\;i}$$

$$Y^i := \alpha X^i + N_Y^{\;i}$$

# Motivation

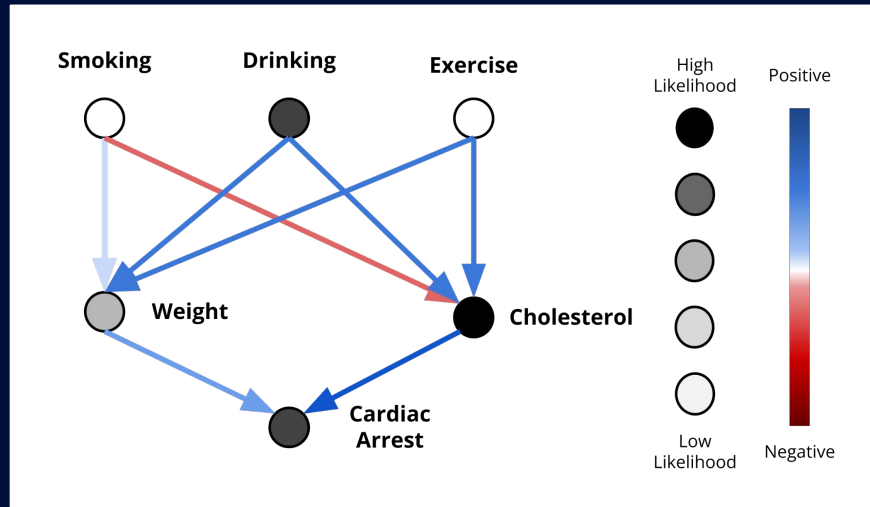**Common Cause Principle:** If two random variables are statistically dependent, then there exists a variable causally influencing both.



Elements of Causal Inference; Peters, Janzing, Schölkopf; 2017

# Motivation

Ideally, explanations are in the form of an SCM connecting concepts,

# Methods

# Methods

**Goals:**

1. Post hoc
2. Causal
3. Concept-based

# Methods

1. Post hoc: train surrogate model

# Methods

2. Causal: SCM inductive bias

# Methods

## 3. Concept-based

# Methods

# Methods



- Concepts and DAG are known a priori.
- We assume:
  - Exogenous independence
  - Concept completeness

# Methods



- **Exogenous model** (outputs exogenous variables $u_k$ for each concept $c_k$).
  - L common neural network layer blocks
  - N concept-specific neural network layer blocks.

# Methods



- **SCM**
  - **Edges:**
    global: $m_{j,k}(\hat{c}_j) = \sigma^{-1}(\hat{c}_j) w_{j,k}$
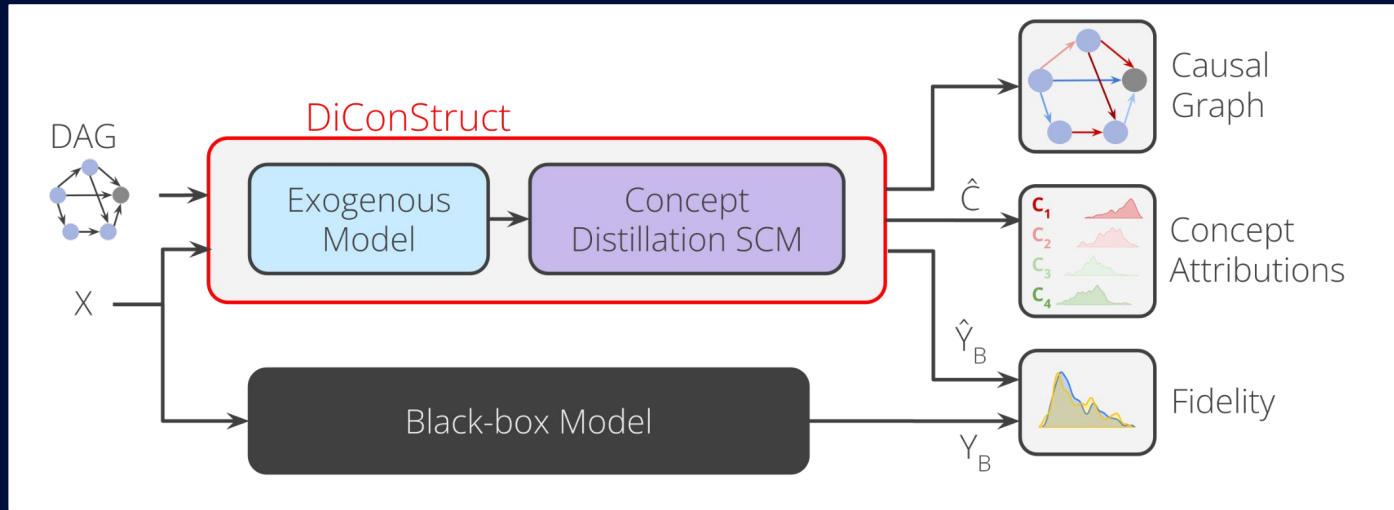    
    local: $m_{j,k}(\hat{c}_j, \boldsymbol{u}, y_B) = \sigma^{-1}(\hat{c}_j) W_{j,k}(\boldsymbol{u}, y_B)$
  - **Concepts:**
    $$\hat{c}_k = \sigma\left(b_k + \sigma^{-1}(u_k) + \sum_{j \in \mathbf{PA}_k} m_{j,k}\right)$$

$b_k, w_{j,k}$ $\quad W_{j,k}$ : **learnable** biases, global weights, local weight functions, respectively.

# Methods



$$\mathcal{L} = \gamma \mathcal{L}_E + \beta \mathcal{L}_C + \mathcal{L}_D$$

- **Objectives:**

$$\mathcal{L}_E = \mathbb{E}_{q(\boldsymbol{u})}\left[\log \frac{f_I(\boldsymbol{u})}{1 - f_I(\boldsymbol{u})}\right] \qquad \mathcal{L}_C = \sum_{i=1}^{M}\sum_{k=1}^{K}\sum_{c_k} c_k^{(i)} \log_2(\hat{c}_k^{(i)}), \qquad \mathcal{L}_D = \sum_{i=1}^{M}\sum_{y_B} y_B^{(i)} \log_2(\hat{y}_B^{(i)}),$$

Exogenous independence loss, concept loss and distillation loss, respectively.

$f_I$ discriminates between the joint distribution of exogenous variables and the product of marginals obtained by randomly shuffling the exogenous variables.

# Methods



- **Concept attributions:**

$$\mathrm{CA}^{(i)}(C_k) = \sum_{a \in \{0,1\}} |\mathbb{P}_{do(C_k := a)}(Y_B^{(i)} = 1) - \mathbb{P}(Y_B^{(i)} = 1)|$$

# Experimental Setup

# Experimental setup

## 1.   Datasets

- **CUB-200-2011**\*: Bird classification (binarized for the purpose of our work)

    Data comes with annotated concepts such as "eye color", "back color", etc.

- **Merchant fraud detection**

    Data manually (and partially) annotated by in-house analysts.
    Remaining data was pseudo-labeled by training "concept teacher" models.

    Concept examples are "high speed ordering", "suspicious device", etc.

\*Wah et al., 2011

# Experimental setup

## 2. Evaluation metrics

- **Main Task Performance**: Given the class imbalance, we chose to use the metric true positive rate (TPR) evaluated at a fixed false positive rate (FPR), which we set to be 5%.

- **Fidelity**: We use the 1 - MAE (mean absolute error).

- **Concept Performance**: Average accuracy over the K concepts.

# Experimental setup

### 3.   Baselines

- **CBM\*:** Concept bottleneck model

- **Distillation CBM**: variation of the above, trained using the same distillation setup as DiConStruct.

- Various ablation studies on the DiConStruct components.

\* Koh et al., ICML 2020

# Experimental setup

## 4. DAGs

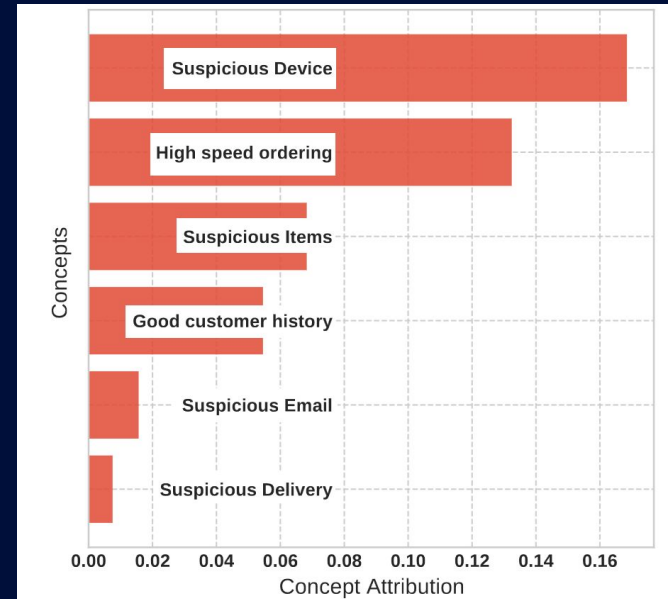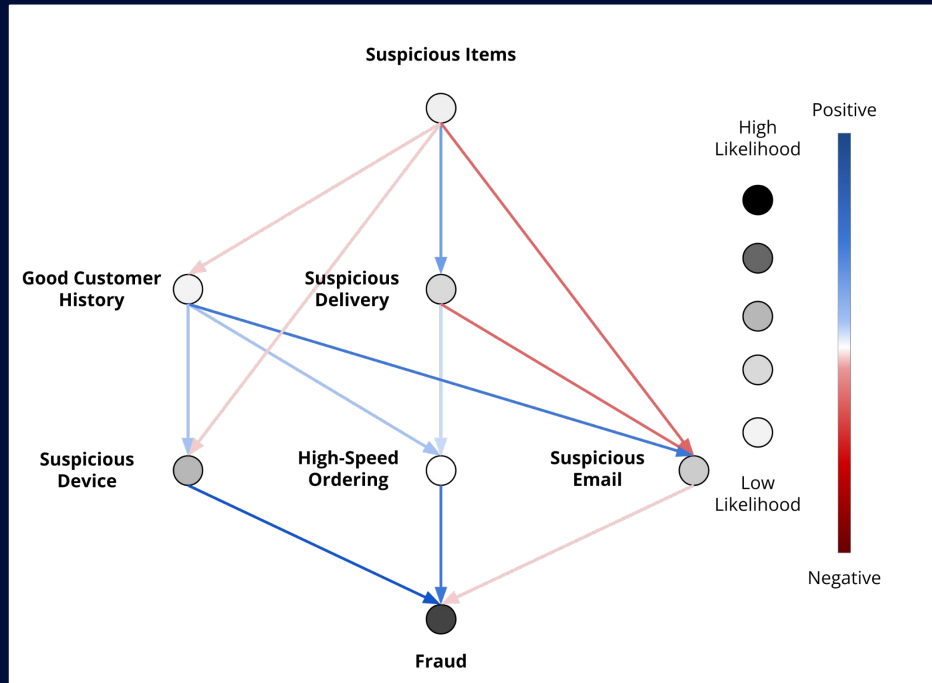We obtain the causal DAG using three causal discovery methods

- **PC algorithm** (Sprites et al., Causation, prediction, and search, 2000)

- **ICA-LiNGAM** (Shimizu et al., JMLR,2006)

- **NO TEARS** (Zheng et al., NeurIPS, 2018)

# Results

# Results

| | Model | Variant | Validation | | | Test | | |
|---|---|---|---|---|---|---|---|---|
| | | | Task Perf. (%) | Concept Perf. (%) | Fidelity (%) | Task Perf. (%) | Concept Perf. (%) | Fidelity (%) |
| **CUB-200-2011** | Ours | Global | 77.85 | $75.58 \pm 0.65$ | $93.52 \pm 0.77$ | 79.05 | $75.44 \pm 0.44$ | $94.3 \pm 0.8$ |
| | | Global w/ Ind. | | $75.61 \pm 0.69$ | $93.16 \pm 0.75$ | | $75.43 \pm 0.65$ | $93.83 \pm 0.64$ |
| | | Local | | $75.25 \pm 0.76$ | $98.72 \pm 0.79$ | | $75.11 \pm 0.63$ | $98.79 \pm 0.74$ |
| | | Local w/ Ind. | | $75.05 \pm 1.08$ | $98.78 \pm 0.86$ | | $74.89 \pm 1.19$ | $98.83 \pm 0.8$ |
| | Baselines | Joint CBM ($\lambda = 1$) | $79.25 \pm 0.98$ | $75.57 \pm 0.46$ | - | $67.33 \pm 2.13$ | $75.76 \pm 0.55$ | - |
| | | Distill. Joint CBM ($\lambda = 1$) | 77.85 | $75.48 \pm 0.53$ | $93.1 \pm 0.52$ | 79.05 | $75.52 \pm 0.59$ | $93.9 \pm 0.56$ |
| | | Single task - Task Perf. | 77.85 | - | - | 79.05 | - | - |
| | | Single task - Concept Perf. | - | $76.11 \pm 0.21$ | - | - | $76.07 \pm 0.26$ | - |
| | | Single task - Fidelity | - | - | $96.07 \pm 0.49$ | - | - | $96.33 \pm 0.26$ |
| **Merchant Fraud - NN** | Ours | Global | 74.67 | $82.64 \pm 0.14$ | $97.12 \pm 0.29$ | 63.35 | $82.58 \pm 0.12$ | $96.62 \pm 0.28$ |
| | | Global w/ Ind. | | $82.6 \pm 0.11$ | $96.96 \pm 0.13$ | | $82.55 \pm 0.09$ | $96.45 \pm 0.24$ |
| | | Local | | $82.5 \pm 0.14$ | $99.39 \pm 0.37$ | | $82.45 \pm 0.13$ | $99.27 \pm 0.42$ |
| | | Local w/ Ind. | | $82.47 \pm 0.13$ | $99.34 \pm 0.41$ | | $82.42 \pm 0.12$ | $99.23 \pm 0.49$ |
| | Baselines | Joint CBM ($\lambda = 1$) | $48.42 \pm 0.31$ | $82.49 \pm 0.14$ | - | $47.47 \pm 3.64$ | $82.34 \pm 0.08$ | - |
| | | Distill. Joint CBM ($\lambda = 1$) | 74.67 | $82.62 \pm 0.13$ | $96.87 \pm 0.18$ | 63.35 | $82.57 \pm 0.12$ | $96.19 \pm 0.29$ |
| | | Single task - Task Perf. | 74.67 | - | - | 63.35 | - | - |
| | | Single task - Concept Perf. | - | $82.25 \pm 0.19$ | - | - | $82.25 \pm 0.19$ | - |
| | | Single task - Fidelity | - | - | $98.13 \pm 0.22$ | - | - | $97.86 \pm 0.23$ |

# Results

# Conclusions

# Conclusions

**Key Takeaways**:

We propose a **novel explainer** that is (1) **concept-based** and **causal**, (2) a **surrogate model** not affecting the predictive performance of the ML model.

**Limitations and future work**:

- Concept completeness assumption
- Multi-class version
- Learning of the DAG