# Debiasing Concept-based Explanations with Causal Analysis

Taha Bahadori & David Heckerman May 3, 2021

Amazon

# Incorporating Domain Knowledge via Concepts



- Concepts allows analysis of models on complex manifolds.
- Concepts bring domain knowledge into the explanation process.

# Confounding and Noise Bias in Concept-Based Explanations

Spearman correlation coefficients ( $\rho$ ) of the predictors of the concepts given features  $\hat{c}(x)$  and labels  $\hat{c}(y)$ .



 $\widehat{c}(x)$  captures spurious correlation and isn't just predictive of y.

# Causal Prior Graph for the Case without Confounding



#### Generative model:

- Generate labels y randomly.
- Generate concepts  $\mathbf{c} \sim p(\mathbf{c}|\mathbf{y})$ .
- Generate images from concepts  $\mathbf{x} \sim p(\mathbf{x}|\mathbf{c})$ .

#### Drawbacks:

- No shared context information between  $\boldsymbol{x}$  and  $\boldsymbol{c}.$
- Concept completeness

## A More Realistic Causal Prior Graph



- Latent variable u represents the shared context between x and c.
- Discriminative concept vector d
- Direct  $x \leftarrow y$  to capture the residual correlation between x and y

## A Technique from Instrumental Variables



– An estimate for the discriminative concepts  $\widehat{d}=\textit{E}[c|y].$  –  $u\perp\!\!\!\perp y\Longrightarrow u\perp\!\!\!\perp \widehat{d}$ 

Our prediction of concepts  $\widehat{c}(x)$  can be uncertain. Incorporate uncertainty into  $\widehat{y}(x) = \widehat{y}(\widehat{c}(x))$ .

To incorporate the uncertainty in our estimation:

$$E[\mathbf{y}|\mathbf{x}] = E[g_{\theta}(\widehat{\mathbf{d}})|\mathbf{x}] = \int g_{\theta}(\mathbf{d}) \mathrm{d}p_{\phi}(\mathbf{d} = \mathbf{d}|\mathbf{x}),$$

The integral is computed using Monte Carlo method.

### Our causal prior graph with linear transformations:

- Generate *n* vector pairs  $\mathbf{y}_i, \mathbf{u}_i \in \mathbb{R}^{100}$  with elements  $\sim \mathcal{N}(0, 1)$ .
- Generate *n* noise vector pairs  $\varepsilon_{c,i}, \varepsilon_{x,i} \in \mathbb{R}^{100}$  with elements  $\sim \mathcal{N}(0, \sigma = 0.02)$ .
- Generate matrices  $W_{y \to d}$ ,  $W_{u \to c}$ ,  $W_{d \to x}$ ,  $W_{u \to x} \in \mathbb{R}^{100 \times 100}$  with elements  $\sim \mathcal{N}(0, \sigma = 0.1)$ .
- Compute  $d_i = W_{y \rightarrow d}y_i + \varepsilon_{d,i}$  for  $i = 1, \dots, n$ .
- Compute  $c_i = d_i + W_{u \to c} u_i$  for  $i = 1, \dots, n$ .
- Compute  $\mathbf{x}_i = \mathbf{W}_{\mathbf{d} \to \mathbf{x}} \mathbf{d}_i + \mathbf{W}_{\mathbf{u} \to \mathbf{x}} \mathbf{u}_i + \boldsymbol{\varepsilon}_{\mathbf{x}, i}$  for  $i = 1, \dots, n$ .



orthogonal  $(W_{y \to d} \perp W_{u \to c}, W_{u \to x})$ 

Regular design

## CUB-200-2011 Data



winter wren

downy woodpecker

bohemian waxwing

- 11788 pictures (in 5994/5794 train/test partitions)
- 200 different types of birds
- Annotations for each picture: bird type and 312 different concepts.
- Randomly choose 15% of the training set as the validation set.

Examples of concepts:

- has\_bill\_shape::dagger, has\_bill\_shape::needle
- has\_wing\_color::purple, has\_wing\_color::blue
- has\_breast\_pattern::solid, has\_breast\_pattern::spotted

## **ROAR Evaluation**

Top-5 accuracy of label prediction improved from 39.5% to 49.3%.



#### ROAR: RemOve And Retrain

– Mask bottom x% of concepts and retrain the  $\mathbf{c} \rightarrow \mathbf{y}$  predictor.  $_{10/10}$