

Exploiting Causal Chains for Domain Generalization

Olawale Salaudeen, Oluwasanmi Koyejo

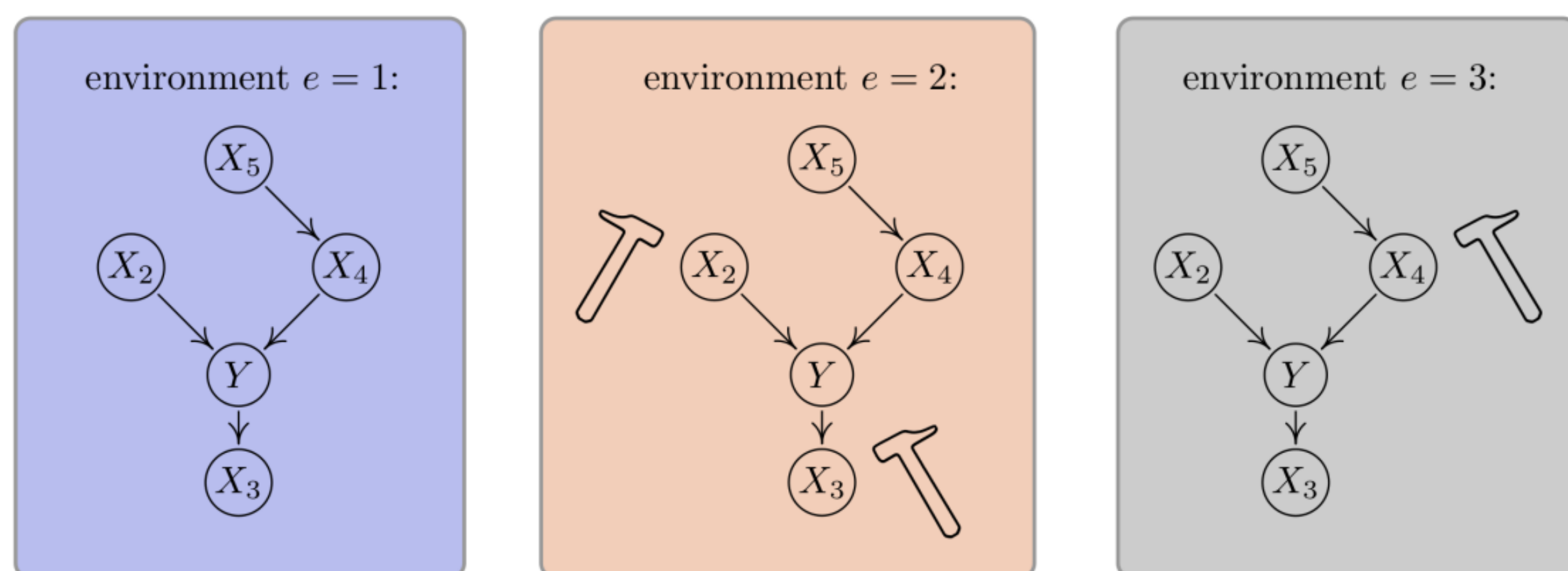
Department of Computer Science, College of Engineering, University of Illinois at Urbana-Champaign

Background

We often encounter distribution shifts from train to test time; how can one learn a predictor that generalizes well to test distributions in this setting?

Invariant Causal Predictors

Assume that datasets belong to distinct domains identified by distinct interventions on the same shared **causal mechanisms**, resulting in different distributions



Invariant Mechanisms: Across domains, $\mu(Y | \text{parents}(Y))$ does not change, though the data distribution can be arbitrarily different. Consequently, identifying parents(Y) is sufficient to learn a predictor that is robust to distribution changes.

Representations that induce Invariant Predictors

When features are latent, a similarly strategy is to learn an embedding function Φ that induces invariance across training distributions, i.e.,

$$\mu(Y | \Phi(X)).$$

Invariant Risk Minimization (IRM): Learn a representation Φ such that the optimal predictor w , across domains \mathcal{E}_{tr} , is the same:

$$\min_{\substack{\Phi: \mathcal{X} \rightarrow \mathcal{H} \\ w: \mathcal{H} \rightarrow \mathcal{Y}}} \sum_{e \in \mathcal{E}_{tr}} R^e(w \circ \Phi)$$

subject to $w \in \arg \min_{\tilde{w}: \mathcal{H} \rightarrow \mathcal{Y}} R^e(\tilde{w} \circ \Phi)$, for all $e \in \mathcal{E}_{tr}$

Let w^* be the optimal invariant predictor that we aim to identify.

The Risks of Invariant Risk Minimization: Let E be the number of distinct training domains and d_e the dimensionality of non-invariant features):

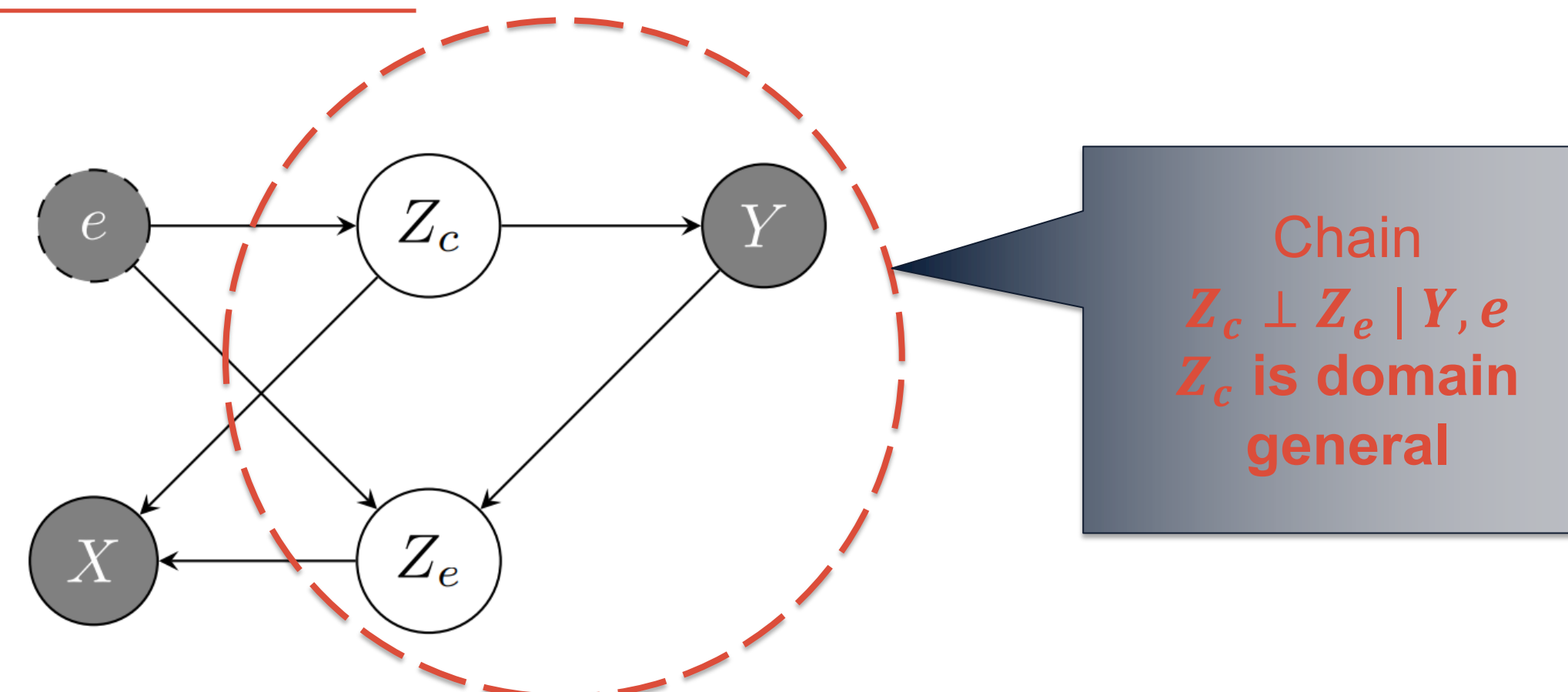
Linear case

- If $E \leq d_e$: There exists a feasible linear Φ which depends on non-invariant features and obtains a lower training risk than the optimal invariant predictor w^*
- If $E > d_e$: The optimal invariant predictor w^* achieves the lowest training risk

Nonlinear case

- IRM behaves just like ERM at test time

Contribution



Shaded nodes are observed while others are latent. e indicates the domain and X a function of the two latents Z_c, Z_e .

Under this generative process in the graph above, we propose to enforce a different Markov property than ICP, namely

$$Z_c \perp Z_e | Y, e. \text{ (TCRI)}$$

\perp means indicates independence.

We empirically show that under the chain generative model, the target conditioned representation invariance (TCRI) constraint yields a predictor that generalizes better than ERM and IRM to test distributions.

Target Conditioned Representation Invariance (TCRI)

Define two feature embedding functions Φ, Ψ – domain general and domain specific, respectively. These two embeddings are related by TCRI:

$$\Phi(X) \perp \Psi(X) | Y, e.$$

$$\min_{\Phi, \Psi, \theta_c, \theta_1, \theta_2, \dots, \theta_E} \sum_{e \in \mathcal{E}} \left[\mathcal{R}^e(\theta_c \circ \Phi) + \beta \mathcal{R}^e(\theta_e \circ \Psi) + \rho \widehat{TCRI}(\Phi^e, \Psi^e) \right]$$

(1) θ_c – optimal linear predictor on Φ across environments.

(2) $\theta_1, \theta_2, \dots, \theta_E$ – optimal linear predictors on Ψ for each environment.

(3) The TCRI captures the constraint $\Phi(X) \perp \Psi(X) | Y \forall e$.

One option is the V-statistic-based Hilbert-Schmidt Independence Criterion (HSIC) estimate:

$$\widehat{HSIC}(X, Y) = \frac{1}{n^2} \text{trace}(\mathbf{K}_{XX'} \mathbf{H}_n \mathbf{K}_{YY'} \mathbf{H}_n),$$

where $\mathbf{K}_{XX'}, \mathbf{K}_{YY'} \in R^{n \times n}$ are Gram matrices, $\mathbf{K}_{XX'}^{ij} = \phi(X_i, X_j)$, $\mathbf{K}_{YY'}^{ij} = \psi(X_i, X_j)$,

$\mathbf{H}_n = \frac{1}{n} \mathbf{I}_n \mathbf{I}_n'$ is a centering matrix, \mathbf{I}_n is the $n \times n$ dimensional identity matrix.

Another is the norm of the conditional cross-covariance:

$$\Sigma_{X_\Phi X_\Psi | Y} = \Sigma_{X_\Phi X_\Psi} - \Sigma_{X_\Phi Y} \Sigma_{Y Y}^{-1} \Sigma_{Y X_\Psi}.$$

Experiments

We evaluate the following linear-Gaussian structural equation model (SEM):

$$SCM(e) = \begin{cases} z_c^e \sim \mathcal{N}(0, (\sigma_c^2)^e I_{d_c}) \\ y^e = z_c^e \alpha + \varepsilon \\ z_e^e = y^e \gamma + \eta \end{cases},$$

where d_c, d_e are the dimensions of z_c, z_e , respectively, $\varepsilon \sim \mathcal{N}(0, (\sigma_\varepsilon^2)^e)$, and $\eta \sim \mathcal{N}(0, (\sigma_\eta^2)^e I_{d_e})$.

$$\Phi, \Psi: \mathbb{R}^{d_c + d_e} \rightarrow \mathbb{R}, \quad \theta_i: \mathbb{R} \rightarrow \mathbb{R} \forall i.$$

- Loss function: Mean Squared Error; TCRI: HSIC
- $\theta_c = 1.0$, is a *dummy* predictor; $x^e = \text{concatenation of } z_c^e, z_e^e$.

Results

Below are relative mean squared errors, $\frac{\text{Algorithm}}{\text{ERM}}$. Errors are computed for each domain independently – average is across all test environment errors and worst case is the worst error achieved on a distinct test domain.

| Algorithm | Average | | Worst Case | |
|-----------------|----------|------|------------|------|
| | Train | Test | Train | Test |
| ERM | baseline | | | |
| IRM | 1.08 | 1.24 | 1.92 | 2.16 |
| TCRI (ours) | 1.38 | 1.20 | 0.16 | 0.11 |
| Causal (Oracle) | 1.29 | 1.13 | 0.11 | 0.06 |

Summary of Results

- ERM performs best on average across train and test domains
 - Does not consider how large a distinct domain's error can be
 - Utilizing non-general features yields lower errors during training
- IRM is worse than ERM in this setting
 - We observe that IRM relies on domain-specific features more than ERM
- TCRI outperforms ERM and IRM in the worst case

Ongoing Work

- Experiments on real-world datasets using TCRI
 - Comparison to more SOA models on benchmark datasets