

# Learning conditionally independent representations with kernel regularizers

**Roman Pogodin\***

Gatsby, UCL → Mila

**Namrata Deka\***

UBC → CMU

**Yazhe Li\***

Gatsby, UCL + DeepMind

**Danica J. Sutherland**

UBC + Amii

**Victor Veitch**

UChicago + Google

**Arthur Gretton**

Gatsby, UCL

LIKE-23, June 2023

based on arXiv:2212.08645 (ICLR 2023, "notable: top 5%")

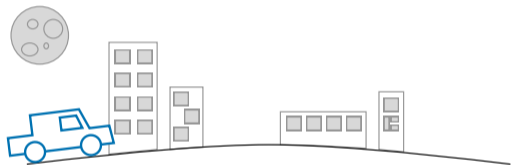


## Intro: conditionally invariant representations

- ▶ Self-driving car tries to predict its location
- ▶ Starts in the morning
- ▶ Finishes in the evening
- ▶ ...learns to predict **location** from **time of day**

**Distribution shift:** car starts in the afternoon

- ▶ ...and makes lots of errors



## Intro: conditionally invariant representations

Idealized solution to this **distribution shift** problem:

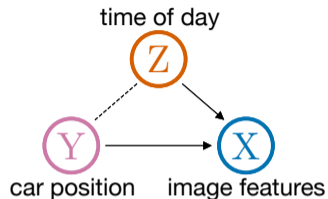
- ▶ **predictions** should be **conditionally independent** of **time** given the **car position**:  $X \perp\!\!\!\perp Z \mid Y$

Same form as a common **domain invariance** objective:

$$\text{features} \perp\!\!\!\perp \text{domain ID} \mid \text{true label}$$

Same form as common **fairness** criterion (equalized odds):

$$\text{predictions} \perp\!\!\!\perp \text{protected attribute} \mid \text{true label}$$



Problem: *conditional* dependence is **hard to measure!**

- ▶ Discrete  $Y$ : check dependence of  $X$  and  $Z$  for *each*  $Y$  value
  - ▶ On *each minibatch* during training...
- ▶ Continuous  $Y$ : prior work runs regression on each minibatch

## Warmup: detecting **unconditional** dependence

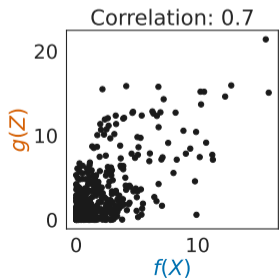
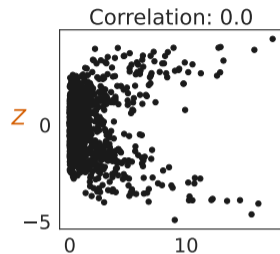
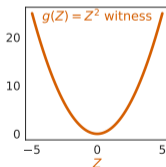
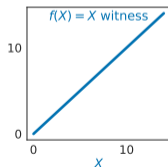
$$Y \sim \mathcal{N}(0, 1)$$

$\xi_1, \xi_2 \sim \mathcal{N}(0, 1)$  i.i.d. noise

$$X = (Y + \xi_1)^2$$

$$Z = Y + \xi_1 + \xi_2$$

►  $X$  and  $Z$  are **uncorrelated**



$X \perp\!\!\!\perp Z$  if and only if **all** square-integrable functions  $f(X)$  and  $g(Z)$  are uncorrelated

## Warmup: detecting **unconditional** dependence

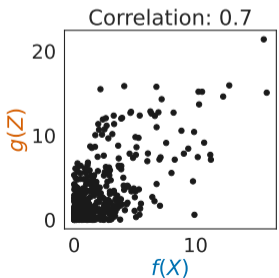
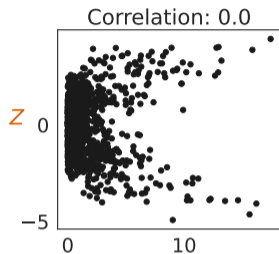
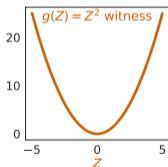
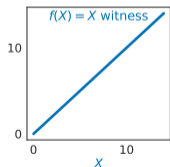
$$Y \sim \mathcal{N}(0, 1)$$

$$\xi_1, \xi_2 \sim \mathcal{N}(0, 1) \text{ i.i.d. noise}$$

$$X = (Y + \xi_1)^2$$

$$Z = Y + \xi_1 + \xi_2$$

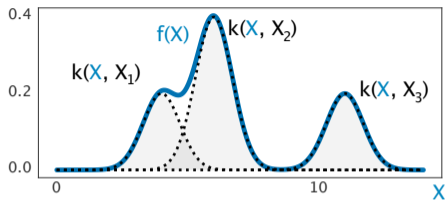
- ▶  $X$  and  $Z$  are **uncorrelated**
- ▶ One way to detect dependence: we can find correlated **nonlinear** functions  $f(X)$  and  $g(Z)$



$X \perp\!\!\!\perp Z$  if and only if **all** square-integrable functions  $f(X)$  and  $g(Z)$  are uncorrelated

## Warmup: detecting **unconditional** dependence

- ▶ If there aren't *any* correlated  $f(X)$  and  $g(Z)$ , then  $X$  and  $Z$  are independent
- ▶ How to check *all enough* nonlinear functions?
- ▶ Check  $f(X)$  and  $g(Z)$  from **kernel spaces** (RKHSes):  $f(X) = \sum_i \alpha_i k(X, X_i)$



- ▶ From RKHS properties:  $\text{Cov}(f(X), g(Z)) = \langle f, C_{XZ} g \rangle$  for the linear operator

$$C_{XZ} = \mathbb{E}[k(X, \cdot) \otimes k(Z, \cdot)] - \mathbb{E}[k(X, \cdot)] \otimes \mathbb{E}[k(Z, \cdot)]$$

- ▶ With linear kernels,  $C_{XZ}$  is just the cross-covariance matrix  $\mathbb{E}[XZ^\top] - \mathbb{E}[X]\mathbb{E}[Z]^\top$
- ▶ If  $C_{XZ} = 0$ , all  $f(X)$  and  $g(Z)$  in the RKHSes are uncorrelated
- ▶ If our kernels are "rich enough" (Gaussian is enough), this implies independence
- ▶ Hilbert-Schmidt Independence Criterion:  $\text{HSIC}(X, Z) = \|C_{XZ}\|_{\text{HS}}^2 = 0$  iff  $C_{XZ} = 0$ 
  - ▶ Can estimate with  $\widehat{\text{HSIC}}(X, Z) = \frac{1}{B^2} \mathbf{1}^\top (HK_{XX}H \odot K_{ZZ}) \mathbf{1}$ , where  $H$  is "centring matrix"
- ▶ Deep nets with features  $X_\theta \sim$ independent of  $Z$ :  $\min_\phi \text{loss}(\phi(X), Y) + \gamma \widehat{\text{HSIC}}(\phi(X), Z)$

## Detecting **conditional** dependence

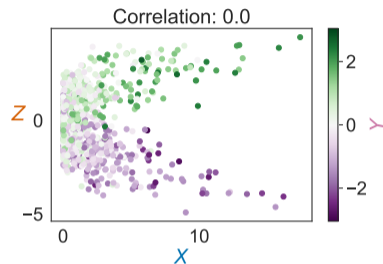
$$Y \sim \mathcal{N}(0, 1)$$

$$\xi_1, \xi_2 \sim \mathcal{N}(0, 1) \text{ i.i.d. noise}$$

$$X = (Y + \xi_1)^2$$

$$Z = Y + \xi_1 + \xi_2$$

- ▶  $X$  and  $Z$  are **dependent**
- ▶  $X$  and  $Z$  are **conditionally dependent** given  $Y$  (through  $\xi_1$ )



## How do we characterize conditional (in)dependence?

- ▶ Start by just conditioning everything on  $Y$ :  $X \perp\!\!\!\perp Z \mid Y$  iff for all  $f_Y \in L^2_X$  and  $g_Y \in L^2_Z$ ,

$$\mathbb{E}_{XZ}[f_Y(X) g_Y(Z) \mid Y] = \mathbb{E}_X[f_Y(X) \mid Y] \mathbb{E}_Z[g_Y(Z) \mid Y] \quad Y\text{-a.s.}$$

- ▶ Equivalent:  $X \perp\!\!\!\perp Z \mid Y$  iff for all  $f \in L^2_{XY}$  and  $g \in L^2_{ZY}$ ,

$$\mathbb{E}_{XZ}[f(X, Y) g(Z, Y) \mid Y] = \mathbb{E}_X[f(X, Y) \mid Y] \mathbb{E}_Z[g(Z, Y) \mid Y] \quad Y\text{-a.s.}$$

- ▶ Equivalent (Daudin 1980):  $X \perp\!\!\!\perp Z \mid Y$  iff for all  $\tilde{f} \in L^2_{XY}$  such that  $\mathbb{E}_X[\tilde{f}(X, Y) \mid Y] = 0 \quad Y\text{-a.s.}$  and all  $\tilde{g} \in L^2_{ZY}$  such that  $\mathbb{E}_Z[\tilde{g}(Z, Y) \mid Y] = 0 \quad Y\text{-a.s.},$

▶ proof

$$\mathbb{E}[\tilde{f}(X, Y) \tilde{g}(Z, Y)] = 0$$

- ▶ Equivalent:  $X \perp\!\!\!\perp Z \mid Y$  iff for all  $f \in L^2_X$ ,  $g \in L^2_{ZY}$ ,

$$\mathbb{E}\left[f(X) (g(Z, Y) - \mathbb{E}_Z[g(Z, Y) \mid Y])\right] = 0$$



# Detecting **conditional** dependence

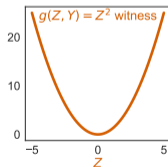
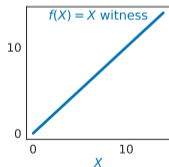
$$Y \sim \mathcal{N}(0, 1)$$

$$\xi_1, \xi_2 \sim \mathcal{N}(0, 1) \text{ i.i.d. noise}$$

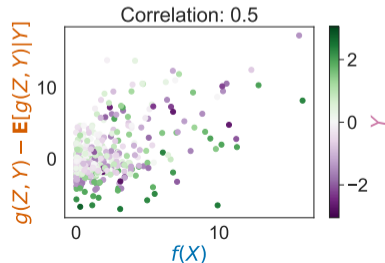
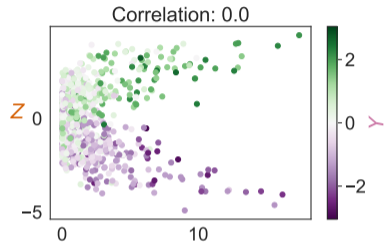
$$X = (Y + \xi_1)^2$$

$$Z = Y + \xi_1 + \xi_2$$

- ▶  $X$  and  $Z$  are **dependent**
- ▶  $X$  and  $Z$  are **conditionally dependent** given  $Y$  (through  $\xi_1$ )



$X \perp\!\!\!\perp Z \mid Y$  if and only if **all**  $f(X)$  are uncorrelated with **all**  $g(Z, Y) - \mathbb{E}[g(Z, Y) \mid Y]$  [Daudin 1980]



## CIRCE: Conditional Independence Regression Covariance

- ▶ Want to check covariance of  $f(X)$  and  $g^c(Z, Y) = g(Z, Y) - \mathbb{E}[g(Z, Y) | Y]$ 
  - ▶  $g^c(Z, Y)$  has mean zero, so they're uncorrelated iff  $\mathbb{E}[f(X) g^c(Z, Y)] = 0$
- ▶ The **CIRCE operator** gives  $\langle f, C_{XZ|Y}^c g \rangle = \mathbb{E}[f(X) g^c(Z, Y)]$ , using

$$C_{XZ|Y}^c = \mathbb{E} \left[ k(X, \cdot) \otimes (k((Z, Y), \cdot) - \mathbb{E}[k((Z', Y), \cdot) | Y]) \right]$$

- ▶  $\text{CIRCE}(X, Z | Y) = \|C_{XZ|Y}^c\|_{\text{HS}}^2 = 0$  iff  $X \perp\!\!\!\perp Z | Y$ , if kernels are “rich enough”
- ▶ Special case: if  $k((Z, Y), (Z', Y')) = k(Z, Z') k(Y, Y')$ , we get

$$C_{XZ|Y}^c = \mathbb{E}[k(X, \cdot) \otimes k(Y, \cdot) \otimes (k(Z, \cdot) - \mu_{Z|Y}(Y))]$$

where  $\mu_{Z|Y}$  is the **conditional mean embedding** of  $Z$  given  $Y$

## CIRCE estimator

- ▶ Want squared norm of  $C_{XZ|Y}^c = \mathbb{E}[k(X, \cdot) \otimes k(Y, \cdot) \otimes (k(Z, \cdot) - \mu_{Z|Y}(Y))]$
- ▶ First, estimate conditional mean embedding  $\mu_{Z|Y}$  on a dataset  $\{(Z_i, Y_i)\}_{i=1}^M$ 
  - ▶ Use kernel ridge regression: inputs  $Y$ , RKHS-valued labels  $k(Z, \cdot)$
  - ▶ Use this to estimate the conditionally-centred kernel function

$$\begin{aligned}\hat{k}^c((Z, Y), (Z', Y')) &= \langle k(Z, \cdot) - \hat{\mu}_{Z|Y}(Y), k(Z', \cdot) - \hat{\mu}_{Z|Y}(Y') \rangle \\ &\approx k(Z, Z') - \mathbb{E}[k(Z, Z') | Y] - \mathbb{E}[k(Z, Z') | Y'] + \mathbb{E}[k(Z, Z') | Y, Y']\end{aligned}$$

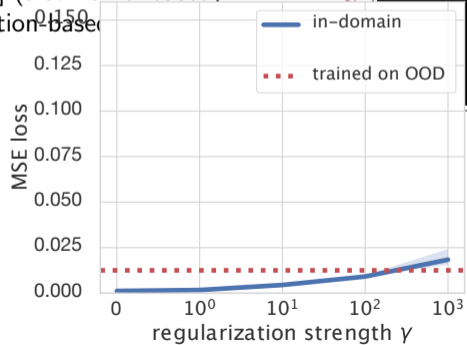
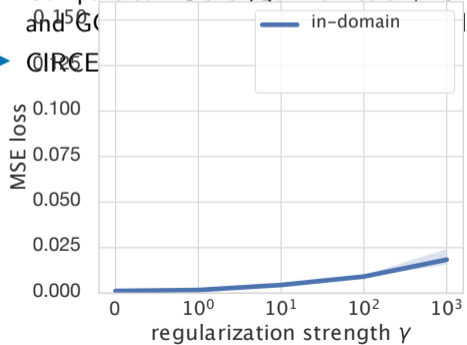
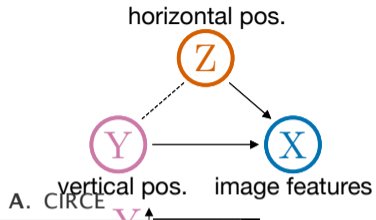
- ▶ While training  $\phi(X)$ , for each batch  $\{(\phi(X_i), Z_i, Y_i)\}_{i=1}^B$ :
  - ▶ Get  $(K_{XX})_{ij} = k(\phi(X_i), \phi(X_j))$ ,  $(K_{YY})_{ij} = k(Y_i, Y_j)$ ,  $(\hat{K}_{ZZ}^c)_{ij} = \hat{k}^c((Z, Y), (Z', Y'))$
  - ▶ Regularize with  $\widehat{\text{CIRCE}} = \frac{1}{B(B-1)} \mathbf{1}^\top (K_{XX} \odot K_{YY} \odot \hat{K}_{ZZ}^c) \mathbf{1}$

### Benefits of CIRCE:

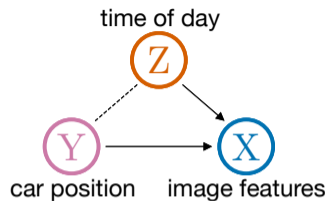
- ▶ As  $B, M \rightarrow \infty$ ,  $\widehat{\text{CIRCE}} \rightarrow 0$  iff  $\phi(X) \perp\!\!\!\perp Z | Y$ ; rate is known (see paper)
- ▶  $K_{YY}$  and  $\hat{K}_{ZZ}^c$  don't depend on  $\phi$ :
  - ▶ Can precompute them, so only need  $k(\phi(X_i), \phi(X_j))$  for each new  $\phi$
  - ▶ Separates (small) batch size  $B$  and (big) regression training size  $M$ : better convergence

# Experiments

- ▶ dSprites dataset [Matthey et al., 2017]:  
2D shapes in different locations
- ▶ Task: predict vertical position  $Y$   
But be invariant to horizontal position  $Z$   
 $Z$  and  $Y$  have strong dependence in training
- ▶ Compare to HSCIC [Quinzan et al., 2022] (also kernel-based) and G
- ▶ CIRCE



## Discussion



- ▶ **CIRCE**: a measure of conditional independence for feature learning
- ▶ It works with continuous variables and in deep learning settings
- ▶ Applications: domain shift invariance, fairness
- ▶ Ongoing extensions:
  - ▶ Learn kernels on **Y** (straightforward) and **Z** (harder)
  - ▶ Testing whether  $\text{CIRCE}(X, Z | Y) = 0$



(code link inside)

# Characterizing conditional (in)dependence – proof sketch

▶ back

$$\forall f \in L^2_{XY}, g \in L^2_{ZY}, \quad \mathbb{E}[fg | Y] = \mathbb{E}[f | Y]\mathbb{E}[g | Y] \quad (\text{A})$$

⇕ [Daudin 1980]

$$\forall \tilde{f} \in L^2_{XY}, \tilde{g} \in L^2_{ZY} \text{ s.t. } \mathbb{E}[\tilde{f} | Y] = 0 = \mathbb{E}[\tilde{g} | Y], \quad \mathbb{E}[\tilde{f}\tilde{g}] = 0 \quad (\text{B})$$

- ▶ (A)  $\implies$  (B), (C): Just apply (A) to  $\tilde{f}$  and  $\tilde{g}$ , RHS becomes 0
- ▶ (B)  $\implies$  (A):
  - ▶ Choose  $\tilde{f}(X, Y) = f(X, Y) - \mathbb{E}[f(X, Y) | Y]$  and  $\tilde{g}(Z, Y) = g(Z, Y) - \mathbb{E}[g(Z, Y) | Y]$ .
  - ▶  $0 = \mathbb{E}[\tilde{f}\tilde{g}] = \mathbb{E}_Y [\mathbb{E}[\tilde{f}\tilde{g} | Y]] = \mathbb{E}_Y [\mathbb{E}[fg | Y] - \mathbb{E}[f | Y]\mathbb{E}[g | Y]]$
  - ▶ Letting  $g$  include an indicator on sets of  $Y$ , implies must hold almost surely in  $Y$
- ▶ Same basic idea works for uncentred  $f \in L^2_{XY}$  and centred  $g \in L^2_{ZY}$
- ▶ Slightly more argument to drop the  $Y$  in  $f$