



# Review of “Learning Representations for Counterfactual Inference”

Johansson et al. (2016)

Suyeol Yun

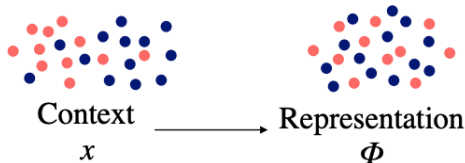
Massachusetts Institute of Technology

June 30, 2022

# What is "Representation Learning"?



- Representation learning is to find effective transformation of input data that makes it easier to perform a task like classification or prediction.
- How to perform representation learning for causal inference (CI) in observational studies to estimate the causal quantity of interest?
- How's the representation learning related to CI?
- Which kind of property is desirable for the representation for CI?
- Balancing property - what else?





- Domain adaptation is a sub-field within machine learning that aims to cope with the situation where the training and the test data fall from different distributions
- Domain adaptation aligns the disparity between domains such that the trained model can be generalized into the domain of interest
- Johansson et al. (2016) view CI as DA and suggests a objective (loss) function to learn effective representation that well generalizes the trained model in factual domain to counterfactual domain

- Binary action set  $\mathcal{T} = \{0, 1\}$
- Causal quantity of interest:  $\text{ITE}(x) = Y_1(x) - Y_0(x)$
- Given  $n$  samples  $\{(x_i, t_i, y_i^F)\}_{i=1}^n$ , where  $y_i^F = t_i \cdot Y_1(x_i) + (1 - t_i) \cdot Y_0(x_i)$ , learn a function  $h: \mathcal{X} \times \mathcal{T} \rightarrow \mathcal{Y}$  such that  $h(x_i, t_i) \approx y_i^F$
- $$\hat{\text{ITE}}(x_i) = \begin{cases} y_i^F - h(x_i, 1 - t_i), & t_i = 1 \\ h(x_i, 1 - t_i) - y_i^F, & t_i = 0 \end{cases}$$
- $\hat{P}^F = \{(x_i, t_i)\}_{i=1}^n$ ,  $\hat{P}^{CF} = \{(x_i, 1 - t_i)\}_{i=1}^n$
- $P^F$  and  $P^{CF}$  need not to be equal- the problem of causal inference requires inference over a different distribution than the one from which samples are given
- In machine learning terms, this means that the feature distribution of the test set differs from that of the train set. This is a case of *covariate shift*, which is a special case of domain adaptation.

# Balancing Objective



- Johansson et al. (2016) proposes a method to jointly learn a representation  $\Phi : \mathcal{X} \rightarrow \mathbb{R}^d$  and  $h : \mathbb{R}^d \times \mathcal{T} \rightarrow \mathbb{R}$  such that the learned representation minimizes the following objective:

$$\begin{aligned} B_{\mathcal{H}, \alpha, \gamma}(\Phi, h) = & \frac{1}{n} \sum_{i=1}^n \left| h(\Phi(x_i), t_i) - y_i^F \right| \\ & + \alpha \operatorname{disc} \mathcal{H} \left( \hat{P}_{\Phi}^F, \hat{P}_{\Phi}^{CF} \right) \\ & + \frac{\gamma}{n} \sum_{i=1}^n \left| h(\Phi(x_i), 1 - t_i) - y_{j(i)}^F \right| \end{aligned}$$

- Let  $j(i) \in \arg \min_{j \in \{1 \dots n\} \text{ s.t. } t_j = 1 - t_i} d(x_j, x_i)$  be the nearest neighbor of  $x_i$  among the group that received

- $\operatorname{disc}_{\mathcal{H}} \left( \hat{P}_{\Phi}^F, \hat{P}_{\Phi}^{CF} \right) =$

$$\max_{\beta, \beta' \in \mathcal{H}} \left| \mathbb{E}_{x \sim \hat{P}_{\Phi}^F} [L(\beta(x), \beta'(x))] - \mathbb{E}_{x \sim \hat{P}_{\Phi}^{CF}} [L(\beta(x), \beta'(x))] \right|$$



- (1) Enabling low-error prediction of the observed outcomes over the factual representation
- (2) Make the sampling distributions of factual and counterfactual to be similar
- (3) Enabling low-error prediction of unobserved counterfactuals by taking into account relevant factual outcomes

$$B_{\mathcal{H},\alpha,\gamma}(\Phi, h) = \frac{1}{n} \sum_{i=1}^n \left| h(\Phi(x_i), t_i) - y_i^F \right| \quad (1)$$

$$+ \alpha \text{disc } \mathcal{H} \left( \hat{P}_{\Phi}^F, \hat{P}_{\Phi}^{CF} \right) \quad (2)$$

$$+ \frac{\gamma}{n} \sum_{i=1}^n \left| h(\Phi(x_i), 1 - t_i) - y_{j(i)}^F \right| \quad (3)$$

$$\bullet \frac{\lambda}{\mu r} \left( \mathcal{L}_{PCF} \left( \hat{\beta}^F(\Phi) \right) - \mathcal{L}_{PCF} \left( \hat{\beta}^{CF}(\Phi) \right) \right)^2 \leq$$

$$\min_{h \in \mathcal{H}_I} \frac{1}{n} \sum_{i=1}^n \left( \left| \hat{y}_i^F(\Phi, h) - y_i^F \right| + \left| \hat{y}_i^{CF}(\Phi, h) - y_{j(i)}^F \right| \right) \quad (4)$$

$$+ \text{disc}_{\mathcal{H}_I} \left( \hat{\rho}_{\Phi}^F, \hat{\rho}_{\Phi}^{CF} \right) \quad (5)$$

$$+ \frac{K_0}{n} \sum_{i:t_i=1} d_{i,j(i)} + \frac{K_1}{n} \sum_{i:t_i=0} d_{i,j(i)} \quad (6)$$

- $\mathcal{H}_I \subset \mathbb{R}^{d+1}$  be the space of linear functions
- $\mathcal{L}_{PCF}(\beta) = \mathbb{E}_{(x,t,y) \sim PCF} [L(\beta(x, t), y)]$  be the expected loss of  $\beta$  over distribution  $P^{CF}$ .
- $\hat{\beta}^F(\Phi) = \arg \min_{\beta \in \mathcal{H}_I} \mathcal{L}_{\hat{\rho}_{\Phi}^F}(\beta) + \lambda \|\beta\|_2^2,$   
 $\hat{\beta}^{CF}(\Phi) = \arg \min_{\beta \in \mathcal{H}_I} \mathcal{L}_{\hat{\rho}_{\Phi}^{CF}}(\beta) + \lambda \|\beta\|_2^2$

$$\bullet \frac{\lambda}{\mu r} \left( \mathcal{L}_{PCF} \left( \hat{\beta}^F(\Phi) \right) - \mathcal{L}_{PCF} \left( \hat{\beta}^{CF}(\Phi) \right) \right)^2 \leq$$

$$\min_{h \in \mathcal{H}_I} \frac{1}{n} \sum_{i=1}^n \left( \left| \hat{y}_i^F(\Phi, h) - y_i^F \right| + \left| \hat{y}_i^{CF}(\Phi, h) - y_{j(i)}^F \right| \right) \quad (7)$$

$$+ \text{disc}_{\mathcal{H}_I} \left( \hat{P}_{\Phi}^F, \hat{P}_{\Phi}^{CF} \right) \quad (8)$$

$$+ \frac{K_0}{n} \sum_{i:t_i=1} d_{i,j(i)} + \frac{K_1}{n} \sum_{i:t_i=0} d_{i,j(i)} \quad (9)$$

- Minimizing the bound make the estimator fit on factual distribution to generalize better over counterfactual distribution
- It's important to find  $\Phi$  such that minimizes the bound
- $\Phi$  such that attains low prediction error, less discrepancy between representation space of factual and counterfactual
- The GB holds regardless of how  $\Phi$  is obtained, e.g. if  $\Phi$  is a neural net - it still holds.





- Johansson et al., Learning Representations for Counterfactual Inference (2016)
- Mansour et al., Domain adaptation: Learning bounds and algorithms (2009)