

On Characterizing the Trade-off in Invariant Representation Learning

Anonymous authors

Paper under double-blind review

Abstract

Many applications of representation learning, such as privacy-preservation, algorithmic fairness, and domain adaptation, desire explicit control over semantic information being discarded. This goal is formulated as satisfying two objectives: maximizing utility for predicting a target attribute while simultaneously being invariant (independent) to a known semantic attribute. Solutions to invariant representation learning (IRL) problems lead to a trade-off between utility and invariance when they are competing. While existing works study bounds on this trade-off, two questions remain outstanding: 1) *What is the exact trade-off between utility and invariance?* and 2) *What are the encoders (mapping the data to a representation) that achieve the trade-off, and how can we estimate it from training data?* This paper addresses these questions for IRLs in reproducing kernel Hilbert spaces (RKHS)s. We derive a closed-form solution for the global optima of the underlying optimization problem for encoders in RKHSs. This in turn yields closed formulae for the exact trade-off, optimal representation dimensionality, and the corresponding encoder(s). We also numerically quantify the trade-off on representative problems and compare them to those achieved by baseline IRL algorithms.

1 Introduction

Real-world applications of representation learning often have to contend with objectives beyond predictive performance. These include cost functions pertaining to, invariance (e.g., to photometric or geometric variations), semantic independence (e.g., to age or race for face recognition systems), privacy (e.g., mitigating leakage of sensitive information (Roy & Boddeti, 2019)), algorithmic fairness (e.g., demographic parity (Madras et al., 2018)), and generalization across multiple domains (Ganin et al., 2016), to name a few.

At its core, the goal of the aforementioned formulations of representation learning is to satisfy two competing objectives: Extracting as much information necessary to predict a target label Y (e.g., face identity) while *intentionally* and *permanently* suppressing information about a given semantic attribute S (e.g., age or gender). See Figure 1 (a) for illustration. Let Z be a representation of input data from which the target attribute Y can be predicted. When the statistical dependency between Y and S is not negligible, learning a representation Z that is invariant to the semantic attribute S (i.e., $Z \perp\!\!\!\perp S$) will necessarily degrade the performance of the target prediction, i.e., there exists a trade-off between utility and invariance. The existence of a trade-off has been well established, both theoretically and empirically, under various contexts of representation learning such as fairness (Menon & Williamson, 2018; Zhao & Gordon, 2019; Gouic et al., 2020; Zhao, 2021), invariance (Zhao et al., 2020), and domain adaptation (Zhao et al., 2019b). However, much of this body of work only establishes bounds on the trade-off, rather than a *precise* characterization. As such, two aspects of the trade-off in invariant representation learning (IRL) are unknown, including, i) *exact* characterization of the trade-off inherent to IRL and ii) a learning algorithm that achieves the trade-off. This paper establishes the mentioned characteristics restricted to functions in reproducing kernel Hilbert spaces (RKHS)s.

Ideally, the utility-invariance trade-off is defined as a bi-objective optimization problem:

$$\inf_{f \in \mathcal{H}_X, g_Y \in \mathcal{H}_Y} \mathbb{E}_{XY} [L_Y(g_Y(f(X)), Y)] \quad \text{such that} \quad \text{Dep}(f(X), S) \leq \epsilon, \quad (1)$$

where f is the encoder that extracts the representation $Z = f(X)$ from X , g_Y predicts \hat{Y} from the representation Z , \mathcal{H}_X and \mathcal{H}_Y are the corresponding hypothesis classes, and L_Y is the loss function for predicting the target attribute Y .

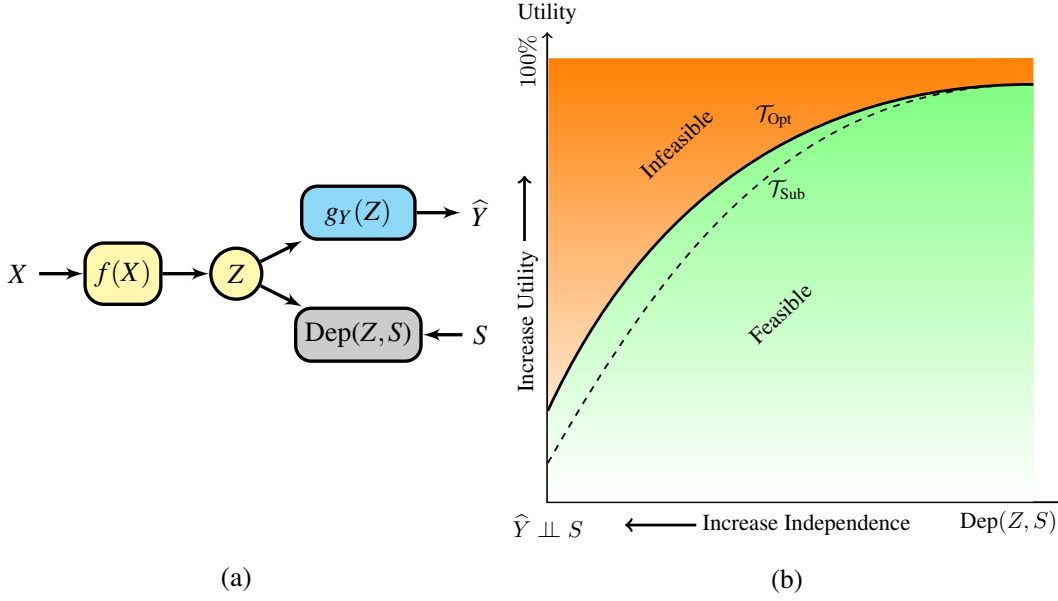


Figure 1: (a): Invariant representation learning seeks a representation $Z = f(X)$ that contains as much information necessary for the downstream target predictor g_Y while being independent of the semantic attribute S . (b): The trade-off (denoted by \mathcal{T}_{Opt}) between utility (target task performance) and invariance (measured by the dependence metric $\text{Dep}(Z, S)$) is induced by a controlled representation learner in the hypothesis class of all Borel functions.

The function $\text{Dep}(\cdot, \cdot) \geq 0$ is a parametric or non-parametric measure of statistical dependence, i.e., $\text{Dep}(Q, U) = 0$ implies Q and U are independent, and $\text{Dep}(Q, U) > 0$ implies Q and U are dependent with larger values indicating greater degrees of dependence. The scalar $\epsilon \geq 0$ is a user-defined parameter that controls the trade-off between the two objectives, with $\epsilon \rightarrow \infty$ being the standard scenario that has no invariance constraints with respect to (w.r.t.) S while $\epsilon \rightarrow 0$ enforces $Z \perp\!\!\!\perp S$ (i.e., total invariance). Involving all Borel functions in \mathcal{H}_X and \mathcal{H}_Y ensures that the best possible trade-off is included within the feasible solution space. For example, when $\epsilon \rightarrow \infty$ and L_Y is MSE loss, the optimal Bayes estimation, $g_Y(f(X)) = \mathbb{E}[Y | X]$ is attainable.

In this paper, we consider the linear combination of utility and invariance in (1) and define the optimal utility-invariance trade-off (denoted by \mathcal{T}_{Opt}) as a single objective optimization problem:

Definition 1.

$$\mathcal{T}_{\text{Opt}} := \inf_{f \in \mathcal{H}_X} \left\{ (1 - \lambda) \inf_{g_Y \in \mathcal{H}_Y} \mathbb{E}_{X,Y} [L_Y(g_Y(f(X)), Y)] + \lambda \text{Dep}(f(X), S) \right\}, \quad 0 \leq \lambda < 1, \quad (2)$$

where λ controls the trade-off between utility and invariance (e.g., $\lambda = 0$ corresponds to ignoring the invariance and only optimizing the utility, while, $\lambda \rightarrow 1$ corresponds to $Z \perp\!\!\!\perp S$).

The motivation behind deploying this single-objective IRL is that any solution to this simplified problem is a solution to the bi-objective problem in (1) and even (2) is challenging to solve, and it has not been fully investigated by existing works. An illustration of the utility-invariance trade-off is illustrated in Figure 1 (b). In this paper, we restrict \mathcal{H}_X to be some RKHSs and $\text{Dep}(Z, S)$ to be a simplified version of the Hilbert-Schmidt Independence Criterion (HSIC) (Gretton et al., 2005a). Further, we replace the target loss function in (2) by $\text{Dep}(Z, Y)$ as presented and justified in Sections 3.2 and 5.2.

Summary of Contributions: i) We design a dependence measure that accounts for all modes of dependence between Z and S (under a mild assumption) while allowing for analytical tractability. ii) We employ functions in RKHSs and obtain closed-form solutions for the IRL optimization problem. Consequently, we exactly characterize \mathcal{T}_{Opt} for encoders restricted to RKHSs. iii) We obtain a closed-form estimator for the encoder that achieves the optimal trade-off, and we establish its numerical convergence. iv) Using random Fourier features (RFF) (Rahimi et al., 2007), we provide a scalable version (in terms of both memory and computation) of our IRL algorithm. v) We numerically quantify our

\mathcal{T}_{Opt} (denoted by $\text{K-}\mathcal{T}_{\text{Opt}}$) on an illustrative problem as well as large scale real-world datasets, Folktables (Ding et al., 2021) and CelebA (Liu et al., 2015), where we compare $\text{K-}\mathcal{T}_{\text{Opt}}$ to those obtained by existing works.

2 Related Work

2.1 Invariant Representation Learning

The basic idea of representation learning that discards unwanted semantic information has been explored under many contexts like invariant, fair, or privacy-preserving learning. In domain adaptation (Ganin & Lempitsky, 2015; Tzeng et al., 2017; Zhao et al., 2018), the goal is to learn features that are independent of the data domain. In fair learning (Dwork et al., 2012; Ruggieri, 2014; Feldman et al., 2015; Calmon et al., 2017; Zemel et al., 2013; Edwards & Storkey, 2015; Beutel et al., 2017; Xie et al., 2017; Zhang et al., 2018; Song et al., 2019; Madras et al., 2018; Bertran et al., 2019; Creager et al., 2019; Locatello et al., 2019; Mary et al., 2019; Martinez et al., 2020; Sadeghi et al., 2019), the goal is to discard the demographic information that leads to unfair outcomes. Similarly, there is growing interest in mitigating unintended leakage of private information from representations (Hamm, 2017; Coavoux et al., 2018; Roy & Boddeti, 2019; Xiao et al., 2020; Dusmanu et al., 2021).

A vast majority of this body of work is empirical in nature. They implicitly look for single or multiple points on the trade-off between utility and semantic information, and do not explicitly seek to characterize the whole trade-off front. Overall, these approaches are not concerned with or aware of the inherent utility-invariance trade-off. In contrast, with the cost of restricting encoders to lie in some RKHSs, we *exactly* characterize the trade-off and propose a practical learning algorithm that achieves this trade-off.

2.2 Adversarial Representation Learning

Most practical approaches for learning fair, invariant, domain adaptive, or privacy-preserving representations discussed above are based on adversarial representation learning (ARL). ARL is typically formulated as

$$\inf_{f \in \mathcal{H}_X} \left\{ (1 - \lambda) \inf_{g_Y \in \mathcal{H}_Y} \mathbb{E}_{X,Y} [L_Y(g_Y(f(X)), Y)] - \lambda \inf_{g_S \in \mathcal{H}_S} \mathbb{E}_{X,S} [L_S(g_S(f(X)), S)] \right\}, \quad (3)$$

where L_S is the loss function of a hypothetical adversary g_S who intends to extract the semantic attribute S through the best estimator within the hypothesis class \mathcal{H}_S and $0 \leq \lambda < 1$ is the utility-invariance trade-off parameter. ARL is a special case of (2) where the negative loss of the adversary, $-\inf_{g_S \in \mathcal{H}_S} \mathbb{E}_{X,S} [L_S(g_S(f(X)), S)]$ plays the role of $\text{Dep}(f(X), S)$. However, this form of adversarial learning suffers from a drawback. The induced independence measure is not guaranteed to account for all modes of non-linear dependence between S and Z if the adversary loss function L_S is not bounded like MSE or cross-entropy (Adeli et al., 2021; Grari et al., 2020). In the case of MSE loss, even if the loss is maximized at a bounded value, where the corresponding representation $Z = f(X)$ is also bounded, still, it is not guaranteed that $Z \perp\!\!\!\perp S$ is attainable (see Appendix H for more details). This implies that, designing the adversary loss in ARL that accounts for all modes of dependence is challenging, and it can be infeasible for some loss functions.

2.3 Trade-Offs in Invariant Representation Learning:

Prior work has established the existence of trade-offs in IRL, both empirically and theoretically. In the following, we categorize them based on properties of interest.

Restricted Class of Attributes: A majority of existing work considers IRL trade-offs under restricted settings, i.e., binary and/or categorical attributes Y and S . For instance, Zhao & Gordon (2019) uses information-theoretic tools and characterizes the utility-fairness trade-off in terms of lower bounds when both Y and S are binary labels. Later McNamara et al. (2019) provided both upper and lower bounds for binary labels. By leveraging Chernoff bound, Dutta et al. (2020) proposed a construction method to generate an ideal representation beyond the input data to achieve perfect fairness while maintaining the best performance on the target task. In the case of categorical features, a lower bound on utility-fairness trade-off has been provided by Zhao et al. (2019a). In contrast to this body of work, our trade-off analysis applies to multidimensional continuous/discrete attributes. To the best of our knowledge, the only

prior works with a general setting are Sadeghi et al. (2019) and Zhao et al. (2020). However, in Zhao et al. (2020), both S and Y are restricted to be continuous/discrete or categorical at the same time (e.g., it is not possible to have Y continuous while S is categorical).

Characterizing Exact vs. Bounds on Trade-Off: To the best of our knowledge, all existing approaches except Sadeghi et al. (2019), which obtains the trade-off for the linear dependence only, characterize the trade-off in terms of upper and/or lower bounds. In contrast, we *exactly* characterize the trade-offs with closed-form expressions for encoders belonging to some RKHSs.

Optimal Encoder and Representation: Another property of practical interest is the optimal encoder that achieves the desired point on the utility-invariance trade-off and the corresponding representation(s). Existing works which only study bounds on the trade-off do not obtain the encoder that achieves those bounds. Sadeghi et al. (2019) do develop a learning algorithm that obtains a globally optimal encoder, but only under a linear dependence measure between Z and S . In contrast, we obtain a closed-form solution for the optimal encoder and its corresponding representation while detecting all modes of dependence between Z and S .

2.4 Kernel Method

The technical machinery of our kernel method for representation learning is closely related to kernelized component analysis (Schölkopf et al., 1998). Kernel methods have been previously used for fair representation learning by Pérez-Suay et al. (2017) where the Rayleigh quotient is employed to only search for a single point in the utility-invariance trade-off. To find the entire trade-off, Sadeghi et al. (2019) used kernelized ARL with the linear adversary and target estimators. Kernel methods also have been used to measure all modes of dependence between two RVs, pioneered by Bach & Jordan (2002) in kernel canonical correlation (KCC). Building upon KCC, later, Gretton et al. (2005a;b; 2006) have introduced HSIC, constrained covariance (COCO), and maximum mean discrepancy (MMD), to name a few. Inspired by these works, a variation of HSIC is deployed as a measure of dependence in this paper.

3 Problem Setting

3.1 Notations

Scalars are denoted by regular lower case letters, e.g., r, λ . Deterministic vectors are denoted by boldface lower case letters, e.g., \mathbf{x}, \mathbf{s} . We denote both scalar-valued and multidimensional random variables (RV)s by regular upper case letters, e.g., X, S . Deterministic matrices are denoted by boldface upper case letters, e.g., $\mathbf{H}, \mathbf{\Theta}$. The entry at i -th row, j -th column of a matrix \mathbf{M} is denoted by $(\mathbf{M})_{ij}$ or m_{ij} . \mathbf{I}_n or simply \mathbf{I} denotes an $n \times n$ identity matrix, $\mathbf{1}_n$ or $\mathbf{1}$ and $\mathbf{0}_n$ or $\mathbf{0}$ are $n \times 1$ a vector of ones and zeros, respectively. We denote the trace of a square matrix \mathbf{K} by $\text{Tr}[\mathbf{K}]$. The pseudo-inverse of a matrix \mathbf{U} is denoted by \mathbf{U}^\dagger . We denote finite or infinite sets by calligraphy letters, e.g., \mathcal{H}, \mathcal{A} .

3.2 Problem Setup

Consider the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where Ω is the sample space, \mathcal{F} is a σ -algebra on Ω , and \mathbb{P} is a probability measure on \mathcal{F} . We assume that the joint RV, (X, Y, S) containing the input data $X \in \mathbb{R}^{d_X}$, the target label $Y \in \mathbb{R}^{d_Y}$, and the semantic attribute $S \in \mathbb{R}^{d_S}$, is an RV on (Ω, \mathcal{F}) with joint distribution $\mathbf{p}_{X,Y,S}$. Furthermore, Y and S can also belong to any finite set like a categorical set. This setting enables us to work with both classification and multidimensional regression tasks, where the semantic attribute can be either categorical or multidimensional continuous/discrete RV.

Assumption 1. We assume that the encoder consists of r functions from \mathbb{R}^{d_X} to \mathbb{R} in a universal RKHS $(\mathcal{H}_X, k_X(\cdot, \cdot))$ (e.g., RBF Gaussian kernel), where universality ensures that \mathcal{H}_X can approximate any Borel function with arbitrary precision (Sriperumbudur et al., 2011).

Hence, the representation RV Z can be expressed as

$$Z = \mathbf{f}(X) := [Z_1, \dots, Z_r]^T \in \mathbb{R}^r, \quad Z_j = f_j(X), f_j \in \mathcal{H}_X \quad \forall j = 1, \dots, r, \quad (4)$$

where r is the dimensionality of the representation. As we will discuss in Corollary 4.1, unlike common practice where r is chosen on an ad-hoc basis, it is an object of interest for optimization. We consider a general scenario where

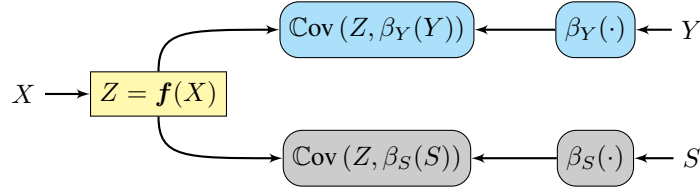


Figure 2: Our IRL model consists of three components: i) An r -dimensional encoder \mathbf{f} belonging to the universal RKHS \mathcal{H}_X . ii) A measure of dependence that accounts for all dependence modes between data representation Z and semantic attribute S induced by the covariance between $Z = \mathbf{f}(X)$ and $\beta_S(S)$ where β_S belongs to a universal RKHS \mathcal{H}_S . iii) A measure of dependency between Z and the target attribute Y defined similarly to that for S .

both Y and S can be continuous/discrete or categorical, or one of Y or S is continuous/discrete while the other is categorical. To accomplish this, we replace the target loss, $\inf_{g_Y \in \mathcal{H}_Y} \mathbb{E}_{X,Y} [L_Y(g_Y(Z), Y)]$ in (2) by the reverse of a non-parametric measure of dependence, i.e., $-\text{Dep}(Z, Y)$. The main reason for this replacement is that maximizing statistical dependency between the representation Z and the target attribute Y can flexibly learn a representation that is applicable for different downstream target tasks, including, regression, classification, clustering, etc (Barshan et al., 2011). Particularly, Theorem 6 in Section 5.2 indicates that with an appropriate choice of involved RKHS for $\text{Dep}(Z, Y)$, we can learn a representation that lends itself to an estimator that performs as optimum as the Bayes estimation, i.e., $\mathbb{E}_X[Y|X]$. Furthermore, in an unsupervised setting, where there is no target attribute Y , the target loss can be replaced with $\text{Dep}(Z, X)$, which implicitly forces the representation Z to be as dependent on the input data X . This scenario is of practical interest when a data producer aims to provide an invariant representation for an unknown downstream target task.

4 Choice of Dependence Measure

We only discuss for $\text{Dep}(Z, S)$ since $\text{Dep}(Z, Y)$ follows similarly. Accounting for all possible non-linear relations between RVs is a key desideratum of dependence measures. A well-known example of such measures is MI (e.g., MINE (Belghazi et al., 2018)). However, calculating MI for multidimensional continuous representation is analytically challenging and computationally intractable. Kernel-based measures are an alternative solution with the attractive properties of being computationally feasible/efficient and analytically tractable (Gretton et al., 2005b).

Definition 2. Let $D = \{(\mathbf{x}_1, \mathbf{y}_1, \mathbf{s}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n, \mathbf{s}_n)\}$ be the training data, containing n i.i.d. samples from the joint distribution $p_{X,Y,S}$. Invoking the representer theorem (Shawe-Taylor & Cristianini, 2004), it follows that for each $f_j \in \mathcal{H}_X$ ($j = 1, \dots, r$) we have $Z_j = f_j(X) = \sum_{i=1}^n \theta_{ji} k_X(\mathbf{x}_i, X)$ where θ_{ji} s are the learnable linear weights. Consequently, it follows that

$$\mathbf{f}(X) = \Theta [k_X(\mathbf{x}_1, X), \dots, k_X(\mathbf{x}_n, X)]^T, \quad (5)$$

where $\Theta \in \mathbb{R}^{r \times n}$ and $(\Theta)_{ji} = \theta_{ji}$.

Principally, $Z \perp\!\!\!\perp S$ if and only if (iff) $\text{Cov}(\alpha(Z), \beta_S(S)) = 0$ for all Borel functions $\alpha : \mathbb{R}^r \rightarrow \mathbb{R}$ and $\beta_S : \mathbb{R}^{d_S} \rightarrow \mathbb{R}$ belonging to the universal RKHSs \mathcal{H}_Z and \mathcal{H}_S , respectively. Alternatively, $Z \perp\!\!\!\perp S$ iff $\text{HSIC}(Z, S) = 0$ for HSIC (Gretton et al., 2005a) being defined as

$$\text{HSIC}(Z, S) := \sum_{\alpha \in \mathcal{U}_Z} \sum_{\beta_S \in \mathcal{U}_S} \text{Cov}^2(\alpha(Z), \beta_S(S)), \quad (6)$$

where \mathcal{U}_Z and \mathcal{U}_S are countable orthonormal basis sets for the separable universal RKHSs \mathcal{H}_Z and \mathcal{H}_S , respectively. However, since $Z = \mathbf{f}(X)$ where \mathbf{f} is defined in (4), calculating $\text{Cov}(\alpha(Z), \beta_S(S))$ necessitates the application of a cascade of kernels, which limits the analytical tractability of our solution. Therefore, we adopt a simplified version of HSIC that considers transformation on S only but affords analytical tractability for solving the IRL optimization problem. We define this measure as

$$\text{Dep}(Z, S) := \sum_{j=1}^r \sum_{\beta_S \in \mathcal{U}_S} \text{Cov}^2(Z_j, \beta_S(S)), \quad (7)$$

where $Z_j = f_j(X)$ for f_j s defined in (4). To guarantee the boundedness of $\text{Dep}(Z, S)$ and $\mathbf{f}(X)$, we consider the following assumption in the remainder of this paper.

Assumption 2. We assume that $(\mathcal{H}_S, k_S(\cdot, \cdot))$ and $(\mathcal{H}_Y, k_Y(\cdot, \cdot))$ are separable¹ and the kernel functions are bounded:

$$\mathbb{E}_U [k_U(U, U)] < \infty, \quad \text{for } U = X, Y, S. \quad (8)$$

The measure $\text{Dep}(Z, S)$ in (7) captures all modes of non-linear dependence under the assumption that the distribution of a low-dimensional projection of high-dimensional data is approximately normal (Diaconis & Freedman, 1984), (Hall & Li, 1993). To see why this reasoning is relevant, we note from (5) that Z can be expressed as $Z = \Theta V$, where $V \in \mathbb{R}^n$ and $\Theta \in \mathbb{R}^{r \times n}$. This indicates that for large n and small r (which is the case for the most real-world datasets), Z is indeed a low-dimensional projection of high-dimensional data. In our numerical experiments in Section 6 we empirically observe that $\text{Dep}(Z, S)$ enjoys a monotonic relation with the underlying invariance measure and captures all modes of dependency in practice, especially as $Z \perp\!\!\!\perp S$.

Lemma 1.² Let $\mathbf{K}_X, \mathbf{K}_S \in \mathbb{R}^{n \times n}$ be the Gram matrices corresponding to \mathcal{H}_X and \mathcal{H}_S , respectively, i.e., $(\mathbf{K}_X)_{ij} = k_X(\mathbf{x}_i, \mathbf{x}_j)$ and $(\mathbf{K}_S)_{ij} = k_S(\mathbf{s}_i, \mathbf{s}_j)$, where covariance is empirically estimated as

$$\text{Cov}(f_j(X), \beta_S(S)) \approx \frac{1}{n} \sum_{i=1}^n f_j(\mathbf{x}_i) \beta_S(\mathbf{s}_i) - \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n f_j(\mathbf{x}_i) \beta_S(\mathbf{s}_k).$$

It follows that, the corresponding empirical estimation for $\text{Dep}(Z, S)$ is

$$\text{Dep}^{\text{emp}}(Z, S) = \frac{1}{n^2} \|\Theta \mathbf{K}_X \mathbf{H} \mathbf{L}_S\|_F^2, \quad (9)$$

where $\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ is the centering matrix, and \mathbf{L}_S is a full column-rank matrix in which $\mathbf{L}_S \mathbf{L}_S^T = \mathbf{K}_S$ (Cholesky factorization). Furthermore, the empirical estimator in (9) has a bias of $\mathcal{O}(n^{-1})$ and a convergence rate of $\mathcal{O}(n^{-1/2})$.

Notice that, the dependence measure between Z and Y can be defined similarly.

5 Exact Kernelized Trade-Off

Consider the optimization problem corresponding to \mathcal{T}_{opt} in (2). Recall that $Z = \mathbf{f}(X)$ is an r -dimensional RV, where the embedding dimensionality r is also a variable to be optimized. A common desideratum of learned representations is that of compactness (Bengio et al., 2013), to avoid learning representations with redundant information where different dimensions are highly correlated to each other. Therefore, going beyond the assumption that each component of \mathbf{f} (i.e., f_j s) belongs to the universal RKHS \mathcal{H}_X , we impose additional constraints on the representation. Specifically, we constrain the search space of the encoder $\mathbf{f}(\cdot)$ to learn a disentangled representation (Bengio et al., 2013) as follows

$$\mathcal{A}_r := \{(f_1, \dots, f_r) \mid f_i, f_j \in \mathcal{H}_X, \text{Cov}(f_i(X), f_j(X)) + \gamma \langle f_i, f_j \rangle_{\mathcal{H}_X} = \delta_{i,j}\}. \quad (10)$$

In the above set, the $\text{Cov}(f_i(X), f_j(X))$ part enforces the covariance matrix of $Z = \mathbf{f}(X)$ to be an identity matrix. This kind of disentanglement is used in the principal component analysis (PCA) and encourages the variance of each entry of Z to be one and different entries of Z to be uncorrelated with each other. The regularization part, $\gamma \langle f_i, f_j \rangle_{\mathcal{H}_X}$ encourages the encoder components to be as orthogonal as possible to each other and to be of unit norm, which aids with numerical stability during empirical estimation (Fukumizu et al., 2007). As the following theorem states formally, such disentanglement is an invertible transformation, and therefore it does not nullify any information.

Theorem 2. Let $Z = \mathbf{f}(X)$ be an arbitrary representation of the input data, where $\mathbf{f} \in \mathcal{H}_X$. Then, there exists an invertible Borel function \mathbf{h} , such that, $\mathbf{h} \circ \mathbf{f}$ belongs to \mathcal{A}_r .

This Theorem implies that the disentanglement preserves the performance of the downstream task, since any target network can revert the disentanglement \mathbf{h} and access to the original representation Z . In addition, any deterministic measurable transformation of Z will not add any information about S that does not already exist in Z .

¹A Hilbert space is separable iff it has a countable orthonormal basis set.

²We differ the proofs of all Lemmas and Theorems to Appendices.

We define our $K\text{-}\mathcal{T}_{\text{Opt}}$ as,

$$\sup_{\mathbf{f} \in \mathcal{A}_r} \{J(\mathbf{f}, \lambda) := (1 - \lambda) \text{Dep}(\mathbf{f}(X), Y) - \lambda \text{Dep}(\mathbf{f}(X), S)\}, \quad 0 \leq \lambda < 1, \quad (11)$$

where λ is the utility-invariance trade-off parameter. Fortunately, the above optimization problem lends itself to a closed-form solution as follows.

Theorem 3. Consider the operator Σ_{SX} to be induced by the bi-linear functional $\text{Cov}(\alpha(X), \beta_S(S)) = \langle \beta_S, \Sigma_{SX} \alpha \rangle_{\mathcal{H}_S}$ and define Σ_{YX} and Σ_{XX} , similarly. Then, a global optima for the optimization problem in (11) is the eigenfunctions corresponding to the r largest eigenvalues of the following generalized eigenvalue problem

$$((1 - \lambda) \Sigma_{YX}^* \Sigma_{YX} - \lambda \Sigma_{SX}^* \Sigma_{SX}) \mathbf{f} = \tau (\Sigma_{XX} + \gamma I_X) \mathbf{f}, \quad (12)$$

where γ is the disentanglement regularization parameter defined in (10), I_X is the identity operator in \mathcal{H}_X , and Σ^* is the adjoint of Σ .

Remark. If the trade-off parameter $\lambda = 0$ (i.e., no semantic independence constraint is imposed) and $\gamma \rightarrow 0$, the solution in Theorem 3 is equivalent to a supervised kernel-PCA. On the other hand, if $\lambda \rightarrow 1$ (i.e., utility is ignored and only semantic independence is considered), the solution in Theorem 3 is the eigenfunctions corresponding to the r smallest eigenvalues of $\Sigma_{SX}^* \Sigma_{SX}$, which are the directions that are the least explanatory of the semantic attribute S .

Now, consider the empirical counterpart of the optimization problem (11),

$$\sup_{\mathbf{f} \in \mathcal{A}_r} \{J^{\text{emp}}(\mathbf{f}, \lambda) := (1 - \lambda) \text{Dep}^{\text{emp}}(\mathbf{f}(X), Y) - \lambda \text{Dep}^{\text{emp}}(\mathbf{f}(X), S)\}, \quad 0 \leq \lambda < 1 \quad (13)$$

where $\text{Dep}^{\text{emp}}(\mathbf{f}(X), S)$ is given in (9) and $\text{Dep}^{\text{emp}}(\mathbf{f}(X), Y)$ is defined similarly.

Theorem 4. Let the Cholesky factorization of \mathbf{K}_X be $\mathbf{K}_X = \mathbf{L}_X \mathbf{L}_X^T$, where $\mathbf{L}_X \in \mathbb{R}^{n \times d}$ ($d \leq n$) is a full column-rank matrix. Let $r \leq d$, then a solution to (13) is

$$\mathbf{f}^{\text{opt}}(X) = \Theta^{\text{opt}} [k_X(\mathbf{x}_1, X), \dots, k_X(\mathbf{x}_n, X)]^T$$

where $\Theta^{\text{opt}} = \mathbf{U}^T \mathbf{L}_X^\dagger$ and the columns of \mathbf{U} are eigenvectors corresponding to the r largest eigenvalues of the following generalized eigenvalue problem.

$$\mathbf{L}_X^T ((1 - \lambda) \tilde{\mathbf{K}}_Y - \lambda \tilde{\mathbf{K}}_S) \mathbf{L}_X \mathbf{u} = \tau \left(\frac{1}{n} \mathbf{L}_X^T \mathbf{H} \mathbf{L}_X + \gamma \mathbf{I} \right) \mathbf{u}. \quad (14)$$

Further, the objective value of (13) is equal to $\sum_{j=1}^r \tau_j$, where $\{\tau_1, \dots, \tau_r\}$ are the r largest eigenvalues of (14).

Corollary 4.1. Embedding Dimensionality: A useful corollary of Theorem 4 is characterizing optimal embedding dimensionality as a function of the trade-off parameter, λ :

$$r^{\text{Opt}}(\lambda) := \arg \sup_{0 \leq r \leq d} \left\{ \sup_{\mathbf{f} \in \mathcal{A}_r} \{J^{\text{emp}}(\mathbf{f}, \lambda)\} \right\} = \text{number of non-negative eigenvalues of (14)}.$$

To examine these results, consider two extreme cases: i) If there is no semantic independence constraint (i.e., $\lambda = 0$), all eigenvalues of (14) are non-negative since $\tilde{\mathbf{K}}_Y$ is a non-negative definite matrix and $\frac{1}{n} \mathbf{L}_X^T \mathbf{H} \mathbf{L}_X + \gamma \mathbf{I}$ is a positive definite matrix. This indicates that r^{Opt} is equal to the maximum possible value (that is equal to d) and therefore it is not required for Z to nullify any information in X . ii) If we only concern about the semantic independence and ignore the target task utility (i.e., $\lambda \rightarrow 1$), all eigenvalues of (14) are non-positive and therefore r^{Opt} would be the number of zero eigenvalues of (14). This indicates that $\text{Dep}^{\text{emp}}(Z, S)$ in (9) is equal to zero, since $\Theta^{\text{opt}} \mathbf{K}_X$ is zero for zero eigenvalues of (14) when $\lambda \rightarrow 1$. In this case, adding more dimension to Z will necessarily increase $\text{Dep}^{\text{emp}}(Z, S)$.

The following Theorem characterizes the convergence behavior of empirical $K\text{-}\mathcal{T}_{\text{Opt}}$ to its population counterpart.

Theorem 5. Assume that k_S and k_Y are bounded by one and $f_j^2(\mathbf{x}_i) \leq M$ for any $j = 1, \dots, r$ and $i = 1, \dots, n$ for which $\mathbf{f} = (f_1, \dots, f_r) \in \mathcal{A}_r$. Then, for any $n > 1$ and $0 < \delta < 1$, with probability at least $1 - \delta$, we have

$$\left| \sup_{\mathbf{f} \in \mathcal{A}_r} J(\mathbf{f}, \lambda) - \sup_{\mathbf{f} \in \mathcal{A}_r} J^{\text{emp}}(\mathbf{f}, \lambda) \right| \leq rM \sqrt{\frac{\log(6/\delta)}{0.22^2 n}} + \mathcal{O}\left(\frac{1}{n}\right).$$

5.1 Numerical Complexity

Computational Complexity: If \mathbf{L}_X in (14) is provided in the training dataset, then, the computational complexity of obtaining the optimal encoder is $\mathcal{O}(l^3)$, where $l \leq n$ is the numerical rank of the Gram matrix \mathbf{K}_X . However, the dominating part of the computational complexity is due to the Cholesky factorization, $\mathbf{K}_X = \mathbf{L}_X \mathbf{L}_X^T$ which is $\mathcal{O}(n^3)$. Using random Fourier features (RFF) (Rahimi et al., 2007), $k_X(\mathbf{x}, \mathbf{x}')$ can be approximated by $\mathbf{r}_X(\mathbf{x})^T \mathbf{r}_X(\mathbf{x}')$, where $\mathbf{r}_X(\mathbf{x}) \in \mathbb{R}^d$. In this situation, the Cholesky factorization can be directly calculated as

$$\mathbf{L}_X = \begin{bmatrix} \mathbf{r}_X(\mathbf{x}_1)^T \\ \vdots \\ \mathbf{r}_X(\mathbf{x}_n)^T \end{bmatrix} \in \mathbb{R}^{n \times d}. \quad (15)$$

As a result the computational complexity of obtaining the optimal encoder becomes $\mathcal{O}(d^3)$, where the RFF dimension, d can be significantly less than the sample size n with negligible sacrifice on $k_X(\mathbf{x}, \mathbf{x}') \approx \mathbf{r}_X(\mathbf{x})^T \mathbf{r}_X(\mathbf{x}')$ approximation.

Memory Complexity: The memory complexity of (14), if calculated naively, is $\mathcal{O}(n^2)$ since \mathbf{K}_Y and \mathbf{K}_S are n by n matrices. However, using RFF together with Cholesky factorization $\mathbf{K}_Y = \mathbf{L}_Y \mathbf{L}_Y^T$, $\mathbf{K}_S = \mathbf{L}_S \mathbf{L}_S^T$, the left-hand side of (14) can be re-arranged as

$$(1 - \lambda) (\mathbf{L}_X^T \tilde{\mathbf{L}}_Y) (\tilde{\mathbf{L}}_Y^T \mathbf{L}_X) - \lambda (\mathbf{L}_X^T \tilde{\mathbf{L}}_S) (\tilde{\mathbf{L}}_S^T \mathbf{L}_X), \quad (16)$$

where $\tilde{\mathbf{L}}_Y^T = \mathbf{H} \mathbf{L}_Y = \mathbf{L}_Y - \frac{1}{n} \mathbf{1}_n (\mathbf{1}_n^T \mathbf{L}_Y)$ and therefore, the required memory complexity is $\mathcal{O}(nd)$. Note that, $\tilde{\mathbf{L}}_S^T$ and $\mathbf{H} \mathbf{L}_X$ can be calculated similarly.

5.2 Target Task Performance in $\mathbf{K}-\mathcal{T}_{\text{Opt}}$

Assume that the desired target loss function is MSE. In the following Theorem, we show that maximizing $\text{Dep}(\mathbf{f}(X), Y)$ over $\mathbf{f} \in \mathcal{A}_r$ can learn a representation Z that is informative enough for a target predictor on Z to achieve the most optimal estimation, i.e., the Bayes estimation ($\mathbb{E}[Y|X]$).

Theorem 6. Let $\lambda = 0$, $\gamma \rightarrow 0$, \mathcal{H}_Y be a linear RKHS, $r \geq d_Y$, and \mathbf{f}^* be the optimal encoder obtained by (14). Then, there exists a linear regressor on top of the representation $Z^* = \mathbf{f}^*(X)$ that can perform as optimal as the Bayes estimator:

$$\begin{aligned} \min_{\mathbf{W} \in \mathbb{R}^{d_Y \times r}, \mathbf{b} \in \mathbb{R}^{d_Y}} \mathbb{E}_{X,Y} [\|\mathbf{W} Z^* + \mathbf{b} - Y\|^2] &= \inf_{h \text{ is Borel}} \mathbb{E}_{X,Y} [\|h(X) - Y\|^2] \\ &= \mathbb{E}_{X,Y} [\|\mathbb{E}[Y|X] - Y\|^2]. \end{aligned}$$

This Theorem implies that not only $\text{Dep}(\mathbf{f}(X), Y)$ can preserve all the necessary information in Z to optimally predict Y , also, the learned representation is simple enough for a linear regressor to achieve the optimal performance.

6 Experiments

In this section, we numerically quantify our $\mathbf{K}-\mathcal{T}_{\text{Opt}}$ through the closed-form solution for the encoder obtained in Section 5 on an illustrative toy example and two real-world datasets, Folktables and CelebA.

6.1 Baselines

We consider two types of baselines: (1) ARL (the main framework for IRL) with MSE or Cross-Entropy as the adversarial loss. Such methods are expected to fail to learn a fully invariant representation (Adeli et al., 2021; Grari et al., 2020). These include (Xie et al., 2017; Zhang et al., 2018; Madras et al., 2018), and SARL (Sadeghi et al., 2019). (2) HSIC-based adversarial loss that accounts for all modes of dependence, and as such is theoretically expected to learn a fully invariant representation (Quadrianto et al., 2019). Among these baselines, except for SARL, all the others are optimized via iterative minimax optimization which is often unstable and not guaranteed to converge. On the other hand, SARL obtains a closed-form solution for the global-optima of the minimax optimization under a linear dependence measure between Z and, S which may fail to capture all modes of dependence between Z and S .

6.2 Datasets

Gaussian Toy Example: We design an illustrative toy example where X and S are mean independent in some dimensions but not fully independent in those dimensions. Specifically, X and S are 4-dimensional continuous RVs and generated as following

$$\begin{aligned} U &= [U_1, U_2, U_3, U_4] \sim \mathcal{N}(\mathbf{0}_4, \mathbf{I}_4), \quad N \sim \mathcal{N}(\mathbf{0}_4, \mathbf{I}_4), \quad U \perp\!\!\!\perp N \\ X &= \cos\left(\frac{\pi}{6}U\right) + 0.005N, \quad S = \left[\sin\left(\frac{\pi}{6}[U_1, U_2]\right), \cos\left(\frac{\pi}{6}[U_3, U_4]\right)\right], \end{aligned} \quad (17)$$

where $\sin(\cdot)$ and $\cos(\cdot)$ are applied point-wise. To generate the target attribute, we define four binary RVs as follows.

$$Y_i = \mathbf{1}_{\{|U_i| > T\}}(U_i), \quad i = 1, 2, 3, 4,$$

where $\mathbf{1}_B(\cdot)$ is the indicator function, and we set $T = 0.6744$, so it holds that $\mathbb{P}[Y_i = 0] = \mathbb{P}[Y_i = 1] = 0.5$ for $i = 1, 2, 3, 4$. Finally, we define Y as a 16-class categorical RV concatenated by Y_i s. Since S is dependent on X through all the dimensions of X , then, a wholly invariant Z (i.e., $Z \perp\!\!\!\perp S$) should not contain any information about X . However, since $[S_1, S_2]$ is only mean independent of $[X_1, X_2]$ (i.e., $\mathbb{E}[S_1, S_2 | X_1, X_2] = \mathbb{E}[S_1, S_2]$), ARL baselines with MSE as the adversary loss, i.e., Xie et al. (2017); Zhang et al. (2018); Madras et al. (2018) and SARL cannot capture the dependency of Z to $[S_1, S_2]$ and result in a representation that is always dependent on $[S_1, S_2]$ (see Section H for theoretical details). We sample 18,000 instances from $p_{X,Y,S}$, independently, and split these samples equally into training, validation, and testing partitions.

Folktables: We consider a fair representation learning task on Folktables (Ding et al., 2021) dataset (a derivation of the US census data). Particularly, we use 2018-Washington census data where the target attribute Y is the employment status (binary) and the semantic attribute S is age (discrete value between 0 and 95 years). We seek to learn a representation that predicts employment status while being fair in demographic parity (DP) w.r.t. age. DP requires that the prediction \hat{Y} be independent of S which can be achieved by enforcing $Z \perp\!\!\!\perp S$. The data contains 76,225 samples, each constructed from 16 different features. We randomly split the data into training (70%), validation (15%), and testing (15%) partitions. Further, we adopt embeddings for categorical features (learned in a supervised fashion by Y) and normalization for continuous/discrete features (by dividing to the maximum value).

CelebA: CelebA dataset (Liu et al., 2015) contains 202,599 face images of 10,177 different celebrities with a standard training, validation, and testing splits. Each image is annotated with 40 different attributes. We choose the target attribute Y as the high cheekbone attribute (binary) and the semantic attributes S to be the concatenation of gender and age (a 4-class categorical RV). The objective of this experiment is similar to that of Folktables. Since raw image data is not appropriate for kernel methods, we pre-train a ResNet-18 (He et al., 2016) (supervised by Y) on CelebA images and extract features of dimension 256. These features are used as the input data for all methods.

6.3 Evaluation Metrics

We use the accuracy of the classification tasks (16-class classification for Gaussian toy example, employment prediction for Folktables, and high cheekbone prediction for CelebA) as a utility. For Folktables and CelebA datasets, we define DP violation as

$$\text{DPV}(\hat{Y}, S) := \mathbb{E}_{\hat{Y}} \left[\text{Var}_S \left(\mathbb{P}[\hat{Y} | S] \right) \right] \quad (18)$$

and use it as a metric to measure the variance (unfairness) of the prediction \hat{Y} w.r.t. the semantic attribute S . For the Gaussian toy example, the above metric is challenging to compute because S is a continuous RV. To circumvent this difficulty, we deploy KCC (Bach & Jordan, 2002)

$$\text{KCC}(Z, S) := \sup_{\alpha \in \mathcal{H}_Z, \beta \in \mathcal{H}_S} \frac{\text{Cov}(\alpha(Z), \beta(S))}{\sqrt{\text{Var}(\alpha(Z)) \text{Var}(\beta(S))}}, \quad (19)$$

as a measure of invariance of Z to S , where \mathcal{H}_Z and \mathcal{H}_S are RBF-Gaussian RKHS. The reason for using KCC instead of HSIC is that, unlike HSIC, KCC is normalized, and therefore it is a more readily interpretable measure for comparing the invariance of representations between different methods.

6.4 Choice of (Y, S) Pair

The existence of a utility-invariance trade-off ultimately depends on the statistical dependency between target and semantic attributes. If $\text{Dep}(Z, S)$ is negligible, then there does not exist a trade-off. Keeping this in mind, we first chose the semantic attribute to be a sensitive attribute for Folktables (i.e., age) and CelebA (i.e., concatenation of age and gender) datasets. Then, we calculated the data imbalance (i.e., $|\mathbb{P}[Y = 0] - 0.5|$) and $\text{KCC}(Y, S)$ for all possible Y s. Finally, we chose Y with a small data imbalance and a moderate $\text{KCC}(Y, S)$. For Folktables dataset, $|\mathbb{P}[\text{employment} = 0] - 0.5| = 0.04$ and $\text{KCC}(\text{employment}, \text{age}) = 0.4$. For CelebA dataset, $|\mathbb{P}[\text{high cheekbone} = 0] - 0.5| = 0.05$ and $\text{KCC}(\text{high cheekbone}, [\text{age}, \text{gender}]) = 0.1$.

6.5 Implementation Details

For all methods, we pick different values of λ (100 λ s for Gaussian toy example and 70 λ s for Folktables and CelebA datasets) between zero and one for obtaining the utility-invariance trade-off. We train the baselines that use a neural network for encoder five times with different random seeds. We let the random seed also change the training-validation-testing split for the Folktables dataset (CelebA and Gaussian datasets have fixed splits).

Embedding Dimensionality: None of the baseline methods have any strategy to find the optimum embedding dimensionality (r) and they all set r to be constant w.r.t. λ . Therefore, we set $r = 15$ (i.e., the minimum dimensionality required to linearly classify 16 different categories) for the Gaussian toy example, that is also equal to r^{Opt} when $\lambda = 0$. For $\text{K-}\mathcal{T}_{\text{Opt}}$, we use $r^{\text{Opt}}(\lambda)$ in Corollary 4.1. See Figure 5 (b) for the plot of r^{Opt} vs λ . For Folktables and CelebA datasets, $r^{\text{Opt}}(\lambda = 0)$ is equal to one, and therefore we let $r = 1$ for all methods and all $0 \leq \lambda < 1$.

$\text{K-}\mathcal{T}_{\text{Opt}}$ (Ours): We let \mathcal{H}_X , \mathcal{H}_S , and \mathcal{H}_Y be RBF Gaussian RKHS, where we compute the corresponding band-widths (i.e., σ s) using the median strategy introduced by Gretton et al. (2007). We optimize the regularization parameter γ in the disentanglement set (10) by minimizing the corresponding target losses over γ s in $\{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ on validation sets. RFF (as discussed in Section 5.1) is deployed for all datasets.

SARL (Sadeghi et al., 2019): SARL method is similar to our $\text{K-}\mathcal{T}_{\text{Opt}}$ except that \mathcal{H}_Y and \mathcal{H}_S are linear RKHSs, and therefore we set σ_X and γ similar to that of $\text{K-}\mathcal{T}_{\text{Opt}}$.

ARL (Xie et al., 2017; Zhang et al., 2018; Madras et al., 2018): The representation $Z = f(X)$ is extracted via the encoder f , which is an MLP (4 hidden layers and 15, 15 neurons for Gaussian data; 3 hidden layers and 128, 64 neurons for Folktables and CelebA datasets). Then, Z is fed to a target task predictor g_Y and a proxy adversary g_S networks where both are MLP with (2 hidden layers, and 16 neurons for Gaussian data, 2 hidden layers, and 128 neurons for Folktables and CelebA datasets). All involved networks (f, g_Y, g_S) are trained end-to-end. We use stochastic gradient descent-ascent (SGDA) Xie et al. (2017) with AdamW (Loshchilov & Hutter, 2017) as an optimizer to alternately train the encoder, target predictor, and proxy adversary networks. We choose batch size as 500 for Gaussian data; and 128 for Folktables and CelebA datasets. Then, the corresponding learning rates are optimized over $\{10^{-2}, 10^{-3}, 5 \times 10^{-4}, 10^{-4}, 10^{-5}\}$ by minimizing the target loss on the corresponding validation sets.

HSIC-IRL (Quadrianto et al., 2019): This method can be formulated as (2) where $\text{Dep}(Z, S)$ is replaced by $\text{HSIC}(Z, S)$. The encoder and target predictor networks have the same architecture as that ARL. Therefore, we follow the same optimization procedure as ARL to train the involved neural networks.

6.6 Results

Utility-Invariance Trade-offs: Figures 3 (a, b, c) show the utility-invariance trade-offs for the Gaussian, Folktables, and CelebA datasets, respectively. The invariance measure for the Gaussian toy example is KCC (19), and the invariance measure for Folktables and CelebA datasets is the fairness measure, DPV (18). We make the following observations: 1) $\text{K-}\mathcal{T}_{\text{Opt}}$ is highly stable and almost spans the entire trade-off front. 2) The baseline method HSIC-IRL, despite using a universal dependence measure, leads to a suboptimal trade-off front due to the lack of convergence guarantees to the global optima. 3) The baselines, ARL and SARL span only a small portion of the trade-off front in the Gaussian toy example, since some dimensions of the semantic attribute S in (17) are mean independent (but not fully independent) to some dimensions of X and therefore the adversary does not provide any information to the encoder to discard $[S_1, S_2]$ from the representation. In this dataset, ARL and SARL baselines do not approach $Z \perp\!\!\!\perp S$,

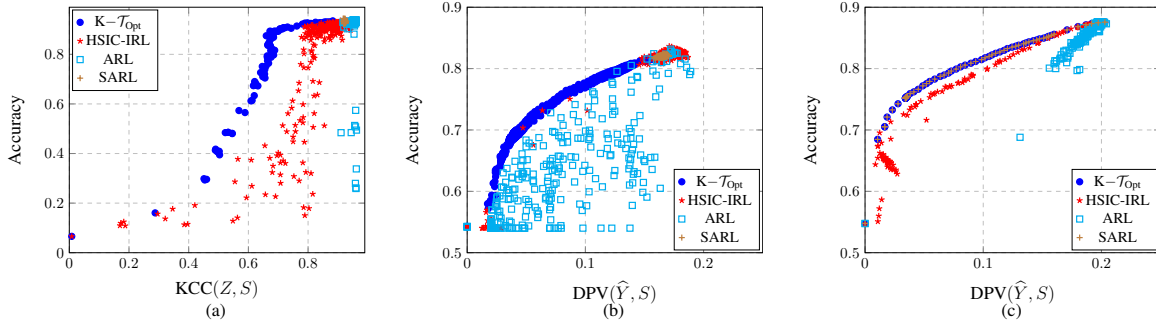


Figure 3: Utility vs. invariance trade-offs obtained by $K-\mathcal{T}_{\text{Opt}}$ and other baselines for (a) Gaussian, (b) Folktables, and (c) CelebA datasets. $K-\mathcal{T}_{\text{Opt}}$ stably spans the entire trade-off front and considerably dominates other methods for all datasets. (a) ARL and SARL span a small portion of the trade-off front since S is mean independent (but not fully independent) of X in some dimensions for the Gaussian toy example. HSIC-IRL, despite using a universal dependence measure, performs sub-optimally due to the lack of convergence guarantees to the global optima.

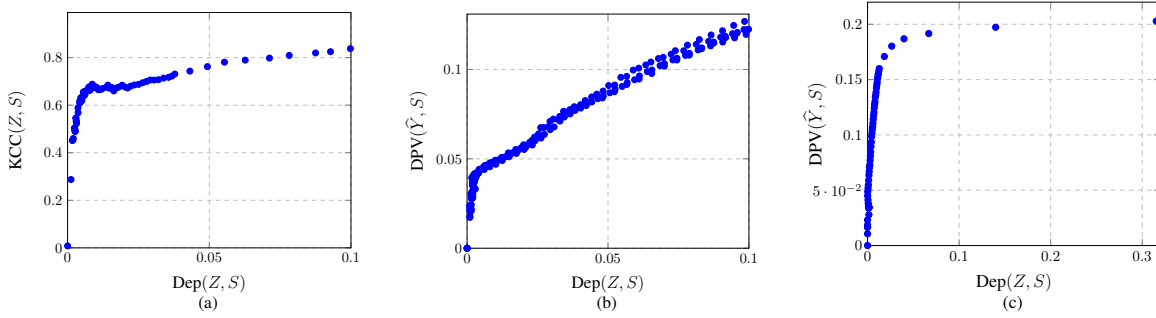


Figure 4: Invariance vs $\text{Dep}(Z, S)$ of $K-\mathcal{T}_{\text{Opt}}$ for (a) Gaussian, (b) Folktables, and (c) CelebA datasets. $\text{Dep}(Z, S)$ enjoys a monotonic relation with the underlying invariance measures.

i.e., $\text{KCC}(Z, S) = 0$ cannot be attained for any value of the trade-off parameter λ . 4) ARL and SARL show high deviation on Folktables dataset due to the unstable nature of the minimax optimization. 5) SARL performs as good as our $K-\mathcal{T}_{\text{Opt}}$ for CelebA dataset. This is because both S and Y are categorical for CelebA dataset, and therefore linear RKHS on one-hot encoded attribute performs just as well as universal RKHSs (Li et al., 2021).

Universality of $\text{Dep}(Z, S)$: We empirically examine the practical validity of our assumption in Section 4 and verify if our dependence measure $\text{Dep}(Z, S)$, defined in (7), can capture all modes of dependency between Z and S . Figure 4 (a) shows the plot of the universal dependence measure $\text{KCC}(Z, S)$ versus $\text{Dep}(Z, S)$ for the Gaussian dataset and Figures 4 (b, c) illustrate the relation between $\text{DPV}(\hat{Y}, S)$ and $\text{Dep}(Z, S)$ for Folktables and CelebA datasets, respectively. We observe that there is a non-decreasing relation between the corresponding invariance measures and $\text{Dep}(Z, S)$. More importantly, as $\text{KCC}(Z, S) \rightarrow 0$ (or $\text{DPV}(\hat{Y}, S) \rightarrow 0$) so does $\text{dep}(Z, S)$. These observations verify that $\text{Dep}(Z, S)$ accounts for all modes of dependence between Z and S .

6.7 Ablation Study

Effect of Embedding Dimensionality: In this experiment, we examine the significance of the embedding dimensionality, $r^{\text{Opt}}(\lambda)$, discussed in Corollary 4.1. We obtain the utility-invariance trade-off when the embedding dimensionality is fixed to be $r = r^{\text{Opt}}(\lambda = 0) = 15$. A comparison plot between the utility-invariance trade-off induced by $r^{\text{Opt}}(\lambda)$ and the fixed $r = 15$ is illustrated in Figure 5 (a). We can observe that not only the utility-invariance trade-off for fixed r is dominated by that of $r^{\text{Opt}}(\lambda)$, but also, using fixed r is unable to achieve the total invariance representation, i.e., $Z \perp\!\!\!\perp S$. Further, $r^{\text{Opt}}(\lambda)$ and some of the largest eigenvalues of (14) vs the invariance trade-off parameter λ are plotted in Figures 5 (b, c), respectively. We recall from Corollary 4.1 that, for any given λ , r^{Opt} is the number of non-negative eigenvalues of (14).

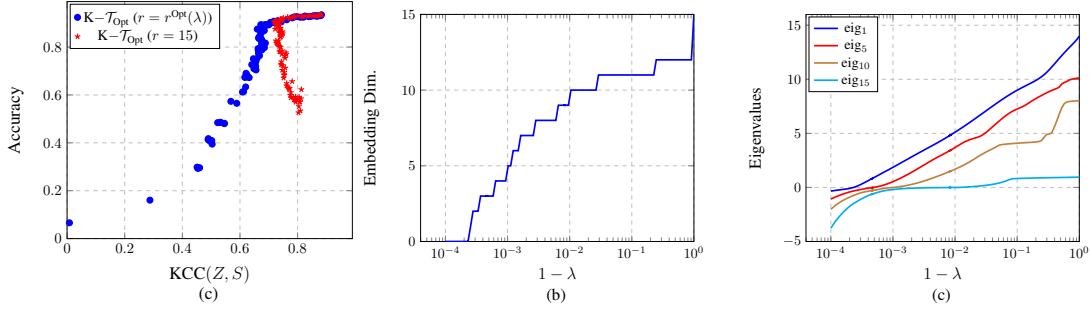


Figure 5: (a) Comparison between the utility-invariance trade-offs induced by the optimal embedding dimensionality $r^{Opt}(\lambda)$ and that of fixed $r = 15$. Fixed $r = 15$ is significantly dominated by that of $r^{Opt}(\lambda)$ and fails to attain $Z \perp\!\!\!\perp S$. (b) The plot of $r^{Opt}(\lambda)$ vs the dependence trade-off parameter $1 - \lambda$. There is a non-decreasing relation between $r^{Opt}(\lambda)$ and $1 - \lambda$. (c) The first, fifth, tenth, and fifteenth largest eigenvalues in (14) vs $1 - \lambda$. Given λ , r^{Opt} is equal to the number of non-negative eigenvalues. As $1 - \lambda$ decreases, the largest eigenvalues approach to negative numbers.

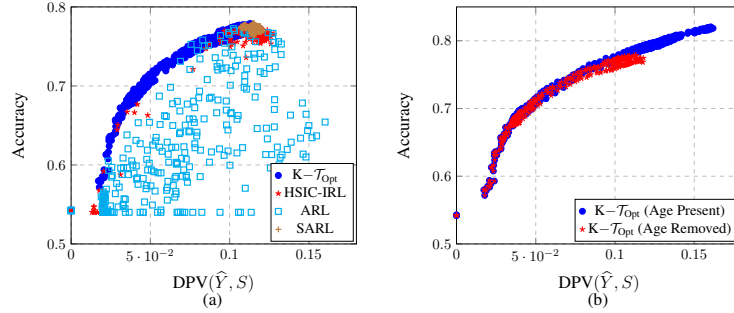


Figure 6: Utility vs. invariance trade-offs obtained by $K-\mathcal{T}_{Opt}$ for Folktables dataset when age is discarded from the input data and: (a) other baselines when age is removed from the input data, (b) $K-\mathcal{T}_{Opt}$ when age is present in the input data. Removing the age attribute slightly degrades the trade-off due to information discarding.

Effect of Semantic Attribute Removal: In this experiment, we examine the effect of removing S (i.e., age) from the input data in the Folktables dataset and examine whether this removal helps the utility-invariance trade-off. Figure 6 (a) shows the utility-invariance trade-off resulting from all methods and Figure 6 (b) compares between removing and keeping the age information from the input data for $K-\mathcal{T}_{Opt}$. Observe that: 1) There is almost the same trend in both keeping and removing the age attribute from the input data for all methods. 2) Removing the age attribute from input data slightly degrades the utility-invariance trade-off due to the lesser information contained in the input data.

7 Conclusion

Invariant representation learning (IRL) often involves a trade-off between utility and invariance. While the existence of such trade-off and its bounds have been studied, its *exact* characterization has not been investigated. This paper takes some steps to address this problem by, i) establishing the *exact* kernelized trade-off (denoted by $K-\mathcal{T}_{Opt}$), ii) determining the optimal dimensionality of the data representation necessary to achieve a desired optimal trade-off point, and iii) developing a scalable learning algorithm for encoders in some RKHSs to achieve $K-\mathcal{T}_{Opt}$. Numerical results on both an illustrative example and two real-world datasets show that commonly used adversarial representation learning-based techniques are unable to attain the optimal trade-off estimated by our solution.

Our theoretical results and empirical solutions shed light on the utility-invariance trade-off for various settings such as algorithmic fairness and privacy-preserving learning under scalarization of the bi-objective trade-off formulation. The trade-off in IRL is also a function of the involved dependence measure that quantifies the dependence of learned representations on the semantic attribute. As such, the trade-off obtained in this paper is optimal for HSIC-like dependence measures. Studying the bi-objective trade-off (rather than the scalarization) and other universal measures are interesting directions for future work.

References

- Ehsan Adeli, Qingyu Zhao, Adolf Pfefferbaum, Edith V Sullivan, Li Fei-Fei, Juan Carlos Niebles, and Kilian M Pohl. Representation learning with statistical independence to mitigate bias. *IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2513–2523, 2021.
- Francis R Bach and Michael I Jordan. Kernel independent component analysis. *Journal of Machine Learning Research*, 3(6):1–48, 2002.
- Elnaz Barshan, Ali Ghodsi, Zohreh Azimifar, and Mansoor Zolghadri Jahromi. Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds. *Pattern Recognition*, 44(7): 1357–1371, 2011.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. *International Conference on Machine Learning*, pp. 531–540, 2018.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- Martin Bertran, Natalia Martinez, Afroditi Papadaki, Qiang Qiu, Miguel Rodrigues, Galen Reeves, and Guillermo Sapiro. Adversarially learned representations for information obfuscation and inference. *International Conference on Machine Learning*, 2019.
- Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi. Data decisions and theoretical implications when adversarially learning fair representations. *arXiv preprint arXiv:1707.00075*, 2017.
- Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. Optimized pre-processing for discrimination prevention. *Advances in Neural Information Processing Systems*, pp. 3992–4001, 2017.
- Maximin Coavoux, Shashi Narayan, and Shay B Cohen. Privacy-preserving neural representations of text. *arXiv preprint arXiv:1808.09408*, 2018.
- Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa A Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. *International Conference on Machine Learning*, 2019.
- Persi Diaconis and David Freedman. Asymptotics of graphical projection pursuit. *The Annals of Statistics*, pp. 793–815, 1984.
- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. *arXiv preprint arXiv:2108.04884*, 2021.
- Mihai Dusmanu, Johannes L Schönberger, Sudipta N Sinha, and Marc Pollefeys. Privacy-preserving visual feature descriptors through adversarial affine subspace embedding. *IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- Sanghamitra Dutta, Dennis Wei, Hazar Yueksel, Pin-Yu Chen, Sijia Liu, and Kush Varshney. Is there a trade-off between fairness and accuracy? A perspective using mismatched hypothesis testing. *International Conference on Machine Learning*, pp. 2803–2813, 2020.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. *Innovations in Theoretical Computer Science Conference*, pp. 214–226, 2012.
- Harrison Edwards and Amos Storkey. Censoring representations with an adversary. *arXiv preprint arXiv:1511.05897*, 2015.
- Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 259–268, 2015.

- Kenji Fukumizu, Francis R Bach, and Arthur Gretton. Statistical consistency of kernel canonical correlation analysis. *Journal of Machine Learning Research*, 8(2), 2007.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. *International Conference on Machine Learning*, 2015.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Thibaut Le Gouic, Jean-Michel Loubes, and Philippe Rigollet. Projection to fairness in statistical learning. *arXiv preprint arXiv:2005.11720*, 2020.
- Vincent Grari, Oualid El Hajouji, Sylvain Lamprier, and Marcin Detyniecki. Learning unbiased representations via Rényi minimization. *arXiv preprint arXiv:2009.03183*, 2020.
- Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with Hilbert-Schmidt norms. *International Conference on Algorithmic Learning Theory*, pp. 63–77, 2005a.
- Arthur Gretton, Ralf Herbrich, Alexander Smola, Olivier Bousquet, and Bernhard Schölkopf. Kernel methods for measuring independence. *Journal of Machine Learning Research*, 6(12):2075–2129, 2005b.
- Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alexander J. Smola. A kernel method for the two-sample-problem. *Advances in Neural Information Processing Systems*, 19, 2006.
- Arthur Gretton, Kenji Fukumizu, Choon Hui Teo, Le Song, Bernhard Schölkopf, and Alexander J. Smola. A kernel statistical test of independence. *Advances in Neural Information Processing Systems*, 20:585–592, 2007.
- Peter Hall and Ker-Chau Li. On almost linearity of low dimensional projections from high dimensional data. *The Annals of Statistics*, pp. 867–889, 1993.
- Jihun Hamm. Minimax filter: Learning to preserve privacy from inference attacks. *Journal of Machine Learning Research*, 18(1):4704–4734, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. In *The Collected Works of Wassily Hoeffding*, pp. 409–426. Springer, 1994.
- Jean Jacod and Philip Protter. *Probability essentials*. Springer Science & Business Media, 2012.
- Effrosini Kokiopoulou, Jie Chen, and Yousef Saad. Trace optimization and eigenproblems in dimension reduction methods. *Numerical Linear Algebra with Applications*, 18(3):565–602, 2011.
- Yazhe Li, Roman Pogodin, Danica J Sutherland, and Arthur Gretton. Self-supervised learning with kernel dependence maximization. *arXiv preprint arXiv:2106.08320*, 2021.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. *IEEE International Conference on Computer Vision*, 2015.
- Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and Olivier Bachem. On the fairness of disentangled representations. *Advances in Neural Information Processing Systems*, pp. 14611–14624, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and transferable representations. *arXiv preprint arXiv:1802.06309*, 2018.
- Natalia Martinez, Martin Bertran, and Guillermo Sapiro. Minimax pareto fairness: A multi objective perspective. *International Conference on Machine Learning*, pp. 6755–6764, 2020.

- Han Zhao, Shanghang Zhang, Guanhong Wu, José MF Moura, Joao P Costeira, and Geoffrey J Gordon. Adversarial multiple source domain adaptation. *Advances in Neural Information Processing Systems*, 31:8559–8570, 2018.
- Han Zhao, Jianfeng Chi, Yuan Tian, and Geoffrey J Gordon. Trade-offs and guarantees of adversarial representation learning for information obfuscation. *arXiv preprint arXiv:1906.07902*, 2019a.
- Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. *International Conference on Machine Learning*, pp. 7523–7532, 2019b.
- Han Zhao, Chen Dan, Bryon Aragam, Tommi S Jaakkola, Geoffrey J Gordon, and Pradeep Ravikumar. Fundamental limits and tradeoffs in invariant representation learning. *arXiv preprint arXiv:2012.10713*, 2020.

A A Population Expression for Definition in (7)

A population expression for $\text{Dep}(Z, S)$ in (7) is given in the following.

$$\begin{aligned} \text{Dep}(Z, S) = & \sum_{j=1}^r \left\{ \mathbb{E}_{X,S,X',S'} [f_j(X) f_j(X') k_S(X, X')] + \mathbb{E}_X [f_j(X)] \mathbb{E}_{X'} [f_j(X')] \mathbb{E}_{S,S'} [k_S(X, S')] \right. \\ & \left. - 2 \mathbb{E}_{X,S} [f_j(X) \mathbb{E}_{X'} [f_j(X')] \mathbb{E}_{S'} [k_S(S, X')]] \right\} \end{aligned}$$

where (X', S') is independent of (X, S) with the same distribution as \mathbf{p}_{XS} .

Proof. We first note that this population expression is inspired by that of HSIC (Gretton et al., 2005a).

Consider the operator Σ_{SX} induced by the linear functional $\text{Cov}(\alpha(X), \beta_S(S)) = \langle \beta_S, \Sigma_{SX} \alpha \rangle_{\mathcal{H}_S}$. Then, it follows that

$$\begin{aligned} \text{Dep}(Z, S) &= \sum_{j=1}^r \sum_{\beta_S \in \mathcal{U}_S} \text{Cov}^2(f_j(X), \beta_S(S)) \\ &= \sum_{j=1}^r \sum_{\beta_S \in \mathcal{U}_S} \langle \beta_S, \Sigma_{SX} f_j \rangle_{\mathcal{H}_S}^2 \\ &= \sum_{j=1}^r \sum_{\beta_S \in \mathcal{U}_S} \langle \beta_S, \Sigma_{SX} f_j \rangle_{\mathcal{H}_S}^2 \\ &\stackrel{(a)}{=} \sum_{j=1}^r \|\Sigma_{SX} f_j\|_{\mathcal{H}_S}^2 \\ &= \sum_{j=1}^r \langle \Sigma_{SX} f_j, \Sigma_{SX} f_j \rangle_{\mathcal{H}_S} \\ &\stackrel{(b)}{=} \sum_{j=1}^r \text{Cov}(f_j(X), (\Sigma_{SX} f_j)(S)) \\ &= \sum_{j=1}^r \text{Cov}\left(f_j(X), \langle k_S(\cdot, S), \Sigma_{SX} f_j \rangle_{\mathcal{H}_S}\right) \\ &= \sum_{j=1}^r \text{Cov}(f_j(X), \text{Cov}(f_j(X'), k_S(S', S))) \\ &= \sum_{j=1}^r \text{Cov}(f_j(X), \mathbb{E}_{X',S'} [f_j(X') k_S(S, S')] - \mathbb{E}_{X'} [f_j(X')] \mathbb{E}_{S'} [k_S(S, S')]) \\ &= \sum_{j=1}^r \left\{ \mathbb{E}_{X,S,X',S'} [f_j(X) f_j(X') k_S(S, S')] + \mathbb{E}_X [f_j(X)] \mathbb{E}_{X'} [f_j(X')] \mathbb{E}_{S,S'} [k_S(S, S')] \right. \\ &\quad \left. - 2 \mathbb{E}_{X,S} [f_j(X) \mathbb{E}_{X'} [f_j(X')] \mathbb{E}_{S'} [k_S(S, S')]] \right\} \end{aligned}$$

where (a) is due to Parseval relation for orthonormal basis and (b) is from the definition of Σ_{SX} . \square

B Proof of Lemma 1

Lemma 1. Let $\mathbf{K}_X, \mathbf{K}_S \in \mathbb{R}^{n \times n}$ be Gram matrices corresponding to \mathcal{H}_X and \mathcal{H}_S , respectively, i.e., $(\mathbf{K}_X)_{ij} = k_X(\mathbf{x}_i, \mathbf{x}_j)$ and $(\mathbf{K}_S)_{ij} = k_S(\mathbf{s}_i, \mathbf{s}_j)$, where covariance is empirically estimated as

$$\text{Cov}(f_j(X), \beta_S(S)) \approx \frac{1}{n} \sum_{i=1}^n f_j(\mathbf{x}_i) \beta_S(\mathbf{s}_i) - \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n f_j(\mathbf{x}_i) \beta_S(\mathbf{s}_k).$$

It follows that, the corresponding empirical estimation for $\text{Dep}(Z, S)$ is

$$\text{Dep}^{\text{emp}}(Z, S) := \frac{1}{n^2} \|\Theta K_X H L_S\|_F^2, \quad (20)$$

where $H = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ is the centering matrix, and L_S is a full column-rank matrix in which $L_S L_S^T = K_S$ (Cholesky factorization). Furthermore, the empirical estimator in (9) has a bias of $\mathcal{O}(n^{-1})$ and a convergence rate of $\mathcal{O}(n^{-1/2})$.

Proof. Firstly, let us reconstruct the orthonormal set \mathcal{U}_S when n i.i.d. observations $\{\mathbf{s}_j\}_{j=1}^n$ are given. Invoking representer theorem, for two arbitrary elements β_i and β_m in \mathcal{U}_S , we have

$$\begin{aligned} \langle \beta_i, \beta_m \rangle_{\mathcal{H}_S} &= \left\langle \sum_{j=1}^n \alpha_j k_S(\mathbf{s}_j, \cdot), \sum_{l=1}^n \eta_l k_S(\mathbf{s}_l, \cdot) \right\rangle_{\mathcal{H}_S} \\ &= \sum_{j=1}^n \sum_{l=1}^n \alpha_j \eta_l k_S(\mathbf{s}_j, \mathbf{s}_l) \\ &= \boldsymbol{\alpha}^T K_S \boldsymbol{\eta} \\ &= \langle L_S^T \boldsymbol{\alpha}, L_S^T \boldsymbol{\eta} \rangle_{\mathbb{R}^q} \end{aligned}$$

where $L_S \in \mathbb{R}^{n \times q}$ is a full column-rank matrix and $K_S = L_S L_S^T$ is the Cholesky factorization of K_S . As a result, searching for $\beta_i \in \mathcal{U}_S$ is equivalent to searching for $L_S^T \boldsymbol{\alpha} \in \mathcal{U}_q$ where \mathcal{U}_q is any complete orthonormal set for \mathbb{R}^q . Using empirical expression for covariance, we get

$$\begin{aligned} \text{Dep}^{\text{emp}}(Z, S) &:= \sum_{\beta_S \in \mathcal{U}_S} \sum_{j=1}^r \left\{ \frac{1}{n} \sum_{i=1}^n f_j(\mathbf{x}_i) \beta_S(\mathbf{s}_i) - \frac{1}{n^2} \sum_{i=1}^n f_j(\mathbf{x}_i) \sum_{k=1}^n \beta_S(\mathbf{s}_k) \right\}^2 \\ &= \sum_{L_S^T \boldsymbol{\alpha} \in \mathcal{U}_q} \sum_{j=1}^r \left\{ \frac{1}{n} \boldsymbol{\theta}_j^T K_X K_S \boldsymbol{\alpha} - \frac{1}{n^2} \boldsymbol{\theta}_j^T K_X \mathbf{1}_n \mathbf{1}_n^T K_S \boldsymbol{\alpha} \right\}^2 \\ &= \sum_{L_S^T \boldsymbol{\alpha} \in \mathcal{U}_q} \sum_{j=1}^r \left\{ \frac{1}{n} \boldsymbol{\theta}_j^T K_X H K_S \boldsymbol{\alpha} \right\}^2 \\ &= \sum_{L_S^T \boldsymbol{\alpha} \in \mathcal{U}_q} \sum_{j=1}^r \left\{ \frac{1}{n} \boldsymbol{\theta}_j^T K_X H L_S L_S^T \boldsymbol{\alpha} \right\}^2 \\ &= \sum_{\boldsymbol{\zeta} \in \mathcal{U}_q} \sum_{j=1}^r \left\{ \frac{1}{n} \boldsymbol{\theta}_j^T K_X H L_S \boldsymbol{\zeta} \right\}^2 \\ &= \sum_{\boldsymbol{\zeta} \in \mathcal{U}_q} \frac{1}{n^2} \|\Theta K_X H L_S \boldsymbol{\zeta}\|_2^2 \\ &= \frac{1}{n^2} \|\Theta K_X H L_S\|_F^2, \end{aligned}$$

where $\mathbf{f}(X) = \Theta [k_X(\mathbf{x}_1, X), \dots, k_X(\mathbf{x}_n, X)]^T$ and $\Theta := [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r]^T$.

We now show that the bias of $\text{Dep}^{\text{epm}}(Z, S)$ for estimating $\text{Dep}(Z, S)$ in (9) is $\mathcal{O}(\frac{1}{n})$. To achieve this, we split $\text{Dep}^{\text{epm}}(Z, S)$ into three terms as,

$$\begin{aligned}
\frac{1}{n^2} \|\Theta K_X H L_S\|_F^2 &= \frac{1}{n^2} \text{Tr} \{ \Theta K_X H K_S H K_X \Theta^T \} \\
&= \frac{1}{n^2} \text{Tr} \left\{ \Theta K_X \left(I - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) K_S \left(I - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) K_X \Theta^T \right\} \\
&= \frac{1}{n^2} \underbrace{\text{Tr} \{ K_X \Theta^T \Theta K_X K_S \}}_{\text{I}} - \frac{2}{n^3} \underbrace{\text{Tr} \{ \mathbf{1}^T K_X \Theta^T \Theta K_X K_S \mathbf{1} \}}_{\text{II}} \\
&\quad + \frac{1}{n^4} \underbrace{\text{Tr} \{ \mathbf{1}^T K_X \Theta^T \Theta K_X \mathbf{1} \mathbf{1}^T K_S \mathbf{1} \}}_{\text{III}}
\end{aligned} \tag{21}$$

Let \mathbf{c}_p^n denote the set of all p -tuples drawn without replacement from $\{1, \dots, n\}$. Moreover, let $\Theta = [\theta_1, \dots, \theta_r]^T \in \mathbb{R}^{r \times n}$ and $(A)_{ij}$ denote the element of an arbitrary matrix A at i -th row and j -th column. Then, it follows that

(I):

$$\begin{aligned}
\mathbb{E} [\text{Tr} \{ K_X \Theta^T \Theta K_X K_S \}] &= \sum_{k=1}^r \mathbb{E} \left[\text{Tr} \left\{ \underbrace{K_X \theta_k \theta_k^T}_{\alpha_k} K_X K_S \right\} \right] \\
&= \sum_{k=1}^r \mathbb{E} [\text{Tr} \{ \alpha_k \alpha_k^T K_S \}] \\
&= \sum_{k=1}^r \mathbb{E} \left[\sum_i (\alpha_k \alpha_k^T)_{ii} (K_S)_{ii} + \sum_{(i,j) \in \mathbf{c}_2^n} (\alpha_k \alpha_k^T)_{ij} (K_S)_{ji} \right] \\
&= n \sum_{k=1}^r \mathbb{E}_{X,S} [f_k^2(X) k_S(S, S)] + \frac{n!}{(n-2)!} \sum_{k=1}^r \mathbb{E}_{X,S,X',S'} [f_k(X) f_k(X') k_S(S, S')] \\
&= \mathcal{O}(n) + \frac{n!}{(n-2)!} \sum_{k=1}^r \mathbb{E}_{X,S,X',S'} [f_k(X) f_k(X') k_S(S, S')]
\end{aligned} \tag{22}$$

where (X, S) and (X', S') are independently drawn from the joint distribution p_{XS} .

(II):

$$\begin{aligned}
\mathbb{E} [\mathbf{1}^T K_X \Theta^T \Theta K_X K_S \mathbf{1}] &= \sum_{k=1}^r \mathbb{E} \left[\mathbf{1}^T \underbrace{K_X \theta_k \theta_k^T}_{\alpha_k} K_X K_S \mathbf{1} \right] \\
&= \sum_{k=1}^r \mathbb{E} [\mathbf{1}^T \alpha_k \alpha_k^T K_S \mathbf{1}] \\
&= \sum_{k=1}^r \mathbb{E} \left[\sum_{m=1}^n \sum_{i=1}^n \sum_{j=1}^n (\alpha_k \alpha_k^T)_{mi} (K_S)_{mj} \right] \\
&= \sum_{k=1}^r \mathbb{E} \left[\sum_i (\alpha_k \alpha_k^T)_{ii} (K_S)_{ii} + \sum_{(m,j) \in \mathbf{c}_2^n} (\alpha_k \alpha_k^T)_{mm} (K_S)_{mj} \right] \\
&\quad + \sum_{k=1}^r \mathbb{E} \left[\sum_{(m,i) \in \mathbf{c}_2^n} (\alpha_k \alpha_k^T)_{mi} (K_S)_{mm} + \sum_{(m,j) \in \mathbf{c}_2^n} (\alpha_k \alpha_k^T)_{mj} (K_S)_{mj} \right]
\end{aligned}$$

$$\begin{aligned}
& + \sum_{k=1}^r \mathbb{E} \left[\sum_{(m,i,j) \in \mathbf{c}_3^n} (\boldsymbol{\alpha}_k \boldsymbol{\alpha}_k^T)_{mi} (\mathbf{K}_S)_{mj} \right] \\
& = n \sum_{k=1}^r \mathbb{E}_{X,S} [f_k^2(X) k_S(S, S)] + \frac{n!}{(n-2)!} \sum_{k=1}^r \mathbb{E}_{X,S,S'} [f_k^2(X) k_S(S, S')] \\
& + \frac{n!}{(n-2)!} \sum_{k=1}^r \mathbb{E}_{X,S,X'} [f_k(X) f_k(X') k_S(S, S)] \\
& + \frac{n!}{(n-2)!} \sum_{k=1}^r \mathbb{E}_{X,S,X',S'} [f_k(X) f_k(X') k_S(S, S')] \\
& + \frac{n!}{(n-3)!} \sum_{k=1}^r \mathbb{E}_{X,S} [f_k(X) \mathbb{E}_{X'} [f_k(X')] \mathbb{E}_{S'} [k_S(S, S')]] \\
& = \mathcal{O}(n^2) + \frac{n!}{(n-3)!} \sum_{k=1}^r \mathbb{E}_{X,S} [f_k(X) \mathbb{E}_{X'} [f_k(X')] \mathbb{E}_{S'} [k_S(S, S')]] . \tag{23}
\end{aligned}$$

(III):

$$\begin{aligned}
\mathbb{E} [\mathbf{1}^T \mathbf{K}_X \boldsymbol{\Theta}^T \boldsymbol{\Theta} \mathbf{K}_X \mathbf{1} \mathbf{1}^T \mathbf{K}_S \mathbf{1}] & = \sum_{k=1}^r \mathbb{E} \left[\mathbf{1}^T \underbrace{\mathbf{K}_X \boldsymbol{\theta}_k}_{\boldsymbol{\alpha}_k} \boldsymbol{\theta}_k^T \mathbf{K}_X \mathbf{1} \mathbf{1}^T \mathbf{K}_S \mathbf{1} \right] \\
& = \sum_{k=1}^r \mathbb{E} [\mathbf{1}^T \boldsymbol{\alpha}_k \boldsymbol{\alpha}_k^T \mathbf{1} \mathbf{1}^T \mathbf{K}_S \mathbf{1}] \\
& = \sum_{k=1}^r \mathbb{E} \left[\sum_{i,j,m,l} (\boldsymbol{\alpha}_k \boldsymbol{\alpha}_k^T)_{ij} (\mathbf{K}_S)_{ml} \right] \\
& = \mathcal{O}(n^3) + \sum_{k=1}^r \mathbb{E} \left[\sum_{(i,j,m,l) \in \mathbf{c}_4^n} (\boldsymbol{\alpha}_k \boldsymbol{\alpha}_k^T)_{ij} (\mathbf{K}_S)_{ml} \right] \\
& = \mathcal{O}(n^3) + \frac{n!}{(n-4)!} \sum_{k=1}^r \mathbb{E}_X [f_k(X)] \mathbb{E}_{X'} [f_k(X')] \mathbb{E}_{S,S'} [k_S(S, S')] \\
& \tag{24}
\end{aligned}$$

Using above calculations together with Lemma 2 lead to

$$\text{Dep}(Z, S) = \mathbb{E} [\text{Dep}^{\text{emp}}(Z, S)] + \mathcal{O} \left(\frac{1}{n} \right).$$

We now obtain the convergence of $\text{dep}^{\text{emp}}(Z, S)$. Consider the decomposition in (21) together with (22), (23), and (24). Let $\boldsymbol{\alpha}_k := \mathbf{K}_X \boldsymbol{\theta}_k$, then it follows that

$$\begin{aligned}
& \mathbb{P} \{ \text{Dep}(Z, S) - \text{Dep}^{\text{emp}}(Z, S) \geq t \} \\
& \leq \mathbb{P} \left\{ \sum_{k=1}^r \mathbb{E}_{X,S,X',S'} [f_k(X) f_k(X') k_S(S, S')] - \frac{(n-2)!}{n!} \sum_{k=1}^r \sum_{(i,j) \in \mathbf{c}_2^n} (\boldsymbol{\alpha}_k \boldsymbol{\alpha}_k^T)_{ij} (\mathbf{K}_S)_{ji} + \mathcal{O} \left(\frac{1}{n} \right) \geq at \right\} \\
& + \mathbb{P} \left\{ \sum_{k=1}^r \mathbb{E}_{X,S} [f_k(X) \mathbb{E}_{X'} [f_k(X')] \mathbb{E}_{S'} [k_S(S, S')]] - \frac{(n-3)!}{n!} \sum_{k=1}^r \sum_{(i,j,m) \in \mathbf{c}_3^n} (\boldsymbol{\alpha}_k \boldsymbol{\alpha}_k^T)_{mi} (\mathbf{K}_S)_{mj} + \mathcal{O} \left(\frac{1}{n} \right) \geq bt \right\} \\
& + \mathbb{P} \left\{ \sum_{k=1}^r \mathbb{E}_X [f_k(X)] \mathbb{E}_{X'} [f_k(X')] \mathbb{E}_{S,S'} [k_S(S, S')] \right\}
\end{aligned}$$

$$-\frac{(n-4)!}{n!} \sum_{k=1}^r \sum_{(i,j,m,l) \in \mathcal{C}_4^n} (\alpha_k \alpha_k^T)_{ij} (\mathbf{K}_S)_{ml} + \mathcal{O}\left(\frac{1}{n}\right) \geq (1-a-b)t \Big\},$$

where $a, b > 0$ and $a + b < 1$. For convenience, we omit the term $\mathcal{O}\left(\frac{1}{n}\right)$ and add it back in the last stage.

Define $\zeta := (X, S)$ and consider the following U-statistics (Hoeffding, 1994)

$$\begin{aligned} u_1(\zeta_i, \zeta_j) &= \frac{(n-2)!}{n!} \sum_{(i,j) \in \mathcal{C}_2^n} \sum_{k=1}^r (\alpha_k \alpha_k^T)_{ij} (\mathbf{K}_S)_{ij} \\ u_2(\zeta_i, \zeta_j, \zeta_m) &= \frac{(n-3)!}{n!} \sum_{(i,j,m) \in \mathcal{C}_3^n} \sum_{k=1}^r (\alpha_k \alpha_k^T)_{mi} (\mathbf{K}_S)_{mj} \\ u_3(\zeta_i, \zeta_j, \zeta_m, \zeta_l) &= \frac{(n-4)!}{n!} \sum_{(i,j,m,l) \in \mathcal{C}_4^n} \sum_{k=1}^r (\alpha_k \alpha_k^T)_{ij} (\mathbf{K}_S)_{ml} \end{aligned}$$

Then, from Hoeffding's inequality (Hoeffding, 1994) it follows that

$$\mathbb{P}\{\text{Dep}(Z, S) - \text{Dep}^{\text{emp}}(Z, S) \geq t\} \leq e^{\frac{-2a^2 t^2}{2r^2 M^2} n} + e^{\frac{-2b^2 t^2}{3r^2 M^2} n} + e^{\frac{-2(1-a-b)^2 t^2}{4r^2 M^2} n},$$

where we assumed that $k_S(\cdot, \cdot)$ is bounded by one and $f_k^2(X_i)$ is bounded by M for any $k = 1, \dots, r$ and $i = 1, \dots, n$.

Further, if $0.22 \leq a < 1$, it holds that

$$e^{\frac{-2a^2 t^2}{2r^2 M^2} n} + e^{\frac{-2b^2 t^2}{3r^2 M^2} n} + e^{\frac{-2(1-a-b)^2 t^2}{4r^2 M^2} n} \leq 3e^{\frac{-a^2 t^2}{r^2 M^2} n}.$$

Consequently, we have

$$\mathbb{P}\{|\text{Dep}(Z, S) - \text{Dep}^{\text{emp}}(Z, S)| \geq t\} \leq 6e^{\frac{-a^2 t^2}{r^2 M^2} n}.$$

Therefore, with probability at least $1 - \delta$, it holds

$$|\text{Dep}(Z, S) - \text{Dep}^{\text{emp}}(Z, S)| \leq \sqrt{\frac{r^2 M^2 \log(6/\delta)}{\alpha^2 n}} + \mathcal{O}\left(\frac{1}{n}\right). \quad (25)$$

□

C Proof of Theorem 2

Theorem 2. Let $Z = \mathbf{f}(X)$ be an arbitrary representation of the input data, where $\mathbf{f} \in \mathcal{H}_X$. Then, there exist an invertible Borel function \mathbf{h} , such that, $\mathbf{h} \circ \mathbf{f}$ belongs to \mathcal{A}_r .

Proof. Recall that the space of disentangled representation is

$$\mathcal{A}_r := \left\{ (f_1, \dots, f_r) \mid f_i, f_j \in \mathcal{H}_X, \mathbb{Cov}(f_i(X), f_j(X)) + \gamma \langle f_i, f_j \rangle_{\mathcal{H}_X} = \delta_{i,j} \right\},$$

where $\gamma > 0$. Let I_X denote the identity operator from \mathcal{H}_X to \mathcal{H}_X . We claim that $\mathbf{h} := [h_1, \dots, h_r]$, where

$$\begin{aligned} \mathbf{G}_0 &= \begin{bmatrix} \langle f_1, f_1 \rangle_{\mathcal{H}_X} & \cdots & \langle f_1, f_r \rangle_{\mathcal{H}_X} \\ \vdots & \ddots & \vdots \\ \langle f_r, f_1 \rangle_{\mathcal{H}_X} & \cdots & \langle f_r, f_r \rangle_{\mathcal{H}_X} \end{bmatrix} \\ \mathbf{G} &= \mathbf{G}_0^{-1/2} \\ h_j \circ \mathbf{f} &= \sum_{m=1}^r g_{jm} (\Sigma_{XX} + \gamma I_X)^{-1/2} f_j, \quad \forall j = 1, \dots, r \end{aligned}$$

is the desired invertible transformation. To see this, construct

$$\begin{aligned}
& \text{Cov}(h_i(\mathbf{f}(X)), h_j(\mathbf{f}(X))) + \gamma \langle h_i \circ \mathbf{f}, h_j \circ \mathbf{f} \rangle_{\mathcal{H}_X} \\
&= \langle h_i \circ \mathbf{f}, (\Sigma_{XX} + \gamma I_X) h_j \circ \mathbf{f} \rangle_{\mathcal{H}_X} \\
&= \left\langle \sum_{m=1}^r g_{im} (\Sigma_{XX} + \gamma I_X)^{-1/2} f_i, \sum_{k=1}^r g_{jk} (\Sigma_{XX} + \gamma I_X) (\Sigma_{XX} + \gamma I_X)^{-1/2} f_j \right\rangle_{\mathcal{H}_X} \\
&= \sum_{m=1}^r \sum_{k=1}^r g_{im} g_{jk} \langle f_i, f_j \rangle_{\mathcal{H}_X} = (\mathbf{G} \mathbf{G}_0 \mathbf{G})_{ij} = \delta_{i,j}
\end{aligned}$$

The inverse of \mathbf{h} is $\mathbf{h}' := [h'_1, \dots, h'_r]$ where

$$\begin{aligned}
\mathbf{H} &= \mathbf{G}_0^{1/2} \\
h'_j \circ \mathbf{h} &= \sum_{m=1}^r h_{jm} (\Sigma_{XX} + \gamma I_X)^{1/2} h_j, \quad \forall j = 1, \dots, r.
\end{aligned}$$

□

D Proof of Theorem 3

Theorem 3. Consider the operator Σ_{SX} induced by the bi-linear functional $\text{Cov}(\alpha(X), \beta_S(S)) = \langle \beta_S, \Sigma_{SX} \alpha \rangle_{\mathcal{H}_S}$ and define Σ_{YX} and Σ_{XX} , similarly. Then, a global optima for the optimization problem in (11) is the eigenfunctions corresponding to r largest eigenvalues of the following generalized problem

$$((1 - \lambda) \Sigma_{YX}^* \Sigma_{YX} - \lambda \Sigma_{SX}^* \Sigma_{SX}) f = \lambda (\Sigma_{XX} + \gamma I_X) f, \quad (26)$$

where γ is the disentanglement regularization parameter defined in (10), I_X is the identity operator in \mathcal{H}_X , and Σ^* is the adjoint of Σ .

Proof. Consider $\text{Dep}(Z, S)$ in (7):

$$\begin{aligned}
\text{Dep}(Z, S) &= \sum_{\beta_S \in \mathcal{U}_S} \sum_{j=1}^r \text{Cov}^2(f_j(X), \beta_S(S)) \\
&= \sum_{j=1}^r \sum_{\beta_S \in \mathcal{U}_S} \langle \beta_S, \Sigma_{SX} f_j \rangle_{\mathcal{H}_S}^2 \\
&= \sum_{j=1}^r \|\Sigma_{SX} f_j\|_{\mathcal{H}_S}^2,
\end{aligned}$$

where the last step is due to Parseval's identity for orthonormal basis set. Similarly, we have $\text{dep}(Z, Y) = \sum_{j=1}^r \|\Sigma_{YX} f_j\|_{\mathcal{H}_Y}^2$. Recall that $Z = \mathbf{f}(X) = [(f_1(X), \dots, f_r(X))]$, then it follows that

$$\begin{aligned}
J(\mathbf{f}(X)) &= (1 - \lambda) \sum_{j=1}^r \|\Sigma_{YX} f_j\|_{\mathcal{H}_Y}^2 - \lambda \sum_{j=1}^r \|\Sigma_{SX} f_j\|_{\mathcal{H}_S}^2 \\
&= (1 - \lambda) \sum_{j=1}^r \langle \Sigma_{YX} f_j, \Sigma_{YX} f_j \rangle_{\mathcal{H}_Y} - \lambda \sum_{j=1}^r \langle \Sigma_{SX} f_j, \Sigma_{SX} f_j \rangle_{\mathcal{H}_S} \\
&= \sum_{j=1}^r \langle f_j, ((1 - \lambda) \Sigma_{YX}^* \Sigma_{YX} - \lambda \Sigma_{SX}^* \Sigma_{SX}) f_j \rangle_{\mathcal{H}_X},
\end{aligned}$$

where Σ^* is the adjoint operator of Σ . Further, note that $\text{Cov}(f_i(X), f_j(X))$ is equal to $\langle f_i, \Sigma_{XX} f_j \rangle_{\mathcal{H}_X}$. As a result, the optimization problem in (12) can be restated as

$$\sup_{\langle f_i, (\Sigma_{XX} + \gamma I_X) f_k \rangle_{\mathcal{H}_X} = \delta_{i,k}} \sum_{j=1}^r \langle f_j, ((1-\lambda)\Sigma_{YX}^* \Sigma_{YX} - \lambda \Sigma_{SX}^* \Sigma_{SX}) f_j \rangle_{\mathcal{H}_X}, \quad 1 \leq i, k \leq r$$

where I_X denotes identity operator from \mathcal{H}_X to \mathcal{H}_X . This optimization problem is known as generalized Rayleigh quotient (Strawderman, 1999) and a possible solution to it is given by the eigenfunctions corresponding to the r largest eigenvalues of the following generalized problem

$$((1-\lambda)\Sigma_{YX}^* \Sigma_{YX} - \lambda \Sigma_{SX}^* \Sigma_{SX}) f = \lambda (\Sigma_{XX} + \gamma I_X) f.$$

□

E Proofs of Theorem 4 and Corollary 4.1

Theorem 4. Let the Cholesky factorization of \mathbf{K}_X be $\mathbf{K}_X = \mathbf{L}_X \mathbf{L}_X^T$, where $\mathbf{L}_X \in \mathbb{R}^{n \times d}$ ($d \leq n$) is a full column-rank matrix. Let $r \leq d$, then a solution to (13) is

$$\mathbf{f}^{\text{opt}}(X) = \boldsymbol{\Theta}^{\text{opt}} [k_X(\mathbf{x}_1, X), \dots, k_X(\mathbf{x}_n, X)]^T$$

where $\boldsymbol{\Theta}^{\text{opt}} = \mathbf{U}^T \mathbf{L}_X^\dagger$ and the columns of \mathbf{U} are eigenvectors corresponding to the r largest eigenvalues of the following generalized eigenvalue problem.

$$\mathbf{L}_X^T ((1-\lambda)\tilde{\mathbf{K}}_Y - \lambda \tilde{\mathbf{K}}_S) \mathbf{L}_X \mathbf{u} = \tau \left(\frac{1}{n} \mathbf{L}_X^T \mathbf{H} \mathbf{L}_X + \gamma \mathbf{I} \right) \mathbf{u}. \quad (27)$$

Further, the supremum value of (13) is equal to $\sum_{j=1}^r \tau_j$, where $\{\tau_1, \dots, \tau_r\}$ are r largest eigenvalues of (14).

Proof. Consider the Cholesky factorization, $\mathbf{K}_x = \mathbf{L}_x \mathbf{L}_x^T$ where \mathbf{L}_x is a full column-rank matrix. Using the representer theorem, the disentanglement property in (10) can be expressed as

$$\begin{aligned} & \text{Cov}(f_i(X), f_j(X)) + \gamma \langle f_i, f_j \rangle_{\mathcal{H}_X} \\ &= \frac{1}{n} \sum_{k=1}^n f_i(\mathbf{x}_k) f_j(\mathbf{x}_k) - \frac{1}{n^2} \sum_{k=1}^n f_i(\mathbf{x}_k) \sum_{m=1}^n f_j(\mathbf{x}_m) + \gamma \langle f_i, f_j \rangle_{\mathcal{H}_X} \\ &= \frac{1}{n} \sum_{k=1}^n \sum_{t=1}^n \mathbf{K}_X(\mathbf{x}_k, \mathbf{x}_t) \theta_{it} \sum_{m=1}^n \mathbf{K}_X(\mathbf{x}_k, \mathbf{x}_m) \theta_{jm} - \frac{1}{n^2} \boldsymbol{\theta}_i^T \mathbf{K}_X \mathbf{1}_n \mathbf{1}_n^T \mathbf{K}_X \boldsymbol{\theta}_j + \gamma \langle f_i, f_j \rangle_{\mathcal{H}_X} \\ &= \frac{1}{n} (\mathbf{K}_X \boldsymbol{\theta}_i)^T (\mathbf{K}_X \boldsymbol{\theta}_j) - \frac{1}{n^2} \boldsymbol{\theta}_i^T \mathbf{K}_X \mathbf{1}_n \mathbf{1}_n^T \mathbf{K}_X \boldsymbol{\theta}_j + \gamma \left\langle \sum_{k=1}^n \theta_{ik} k_X(\cdot, \mathbf{x}_k), \sum_{t=1}^n \theta_{jt} k_X(\cdot, \mathbf{x}_t) \right\rangle_{\mathcal{H}_X} \\ &= \frac{1}{n} \boldsymbol{\theta}_i^T \mathbf{K}_X \mathbf{H} \mathbf{K}_X \boldsymbol{\theta}_j + \gamma \boldsymbol{\theta}_i^T \mathbf{K}_X \boldsymbol{\theta}_j \\ &= \frac{1}{n} \boldsymbol{\theta}_i^T \mathbf{L}_X (\mathbf{L}_X^T \mathbf{H} \mathbf{L}_X + n\gamma \mathbf{I}) \mathbf{L}_X^T \boldsymbol{\theta}_j \\ &= \delta_{i,j}. \end{aligned}$$

As a result, $\mathbf{f} \in \mathcal{A}_r$ is equivalent to

$$\boldsymbol{\Theta} \mathbf{L}_X \underbrace{\left(\frac{1}{n} \mathbf{L}_X^T \mathbf{H} \mathbf{L}_X + \gamma \mathbf{I} \right)}_{:= \mathbf{C}} \mathbf{L}_X^T \boldsymbol{\Theta}^T = \mathbf{I}_r,$$

where $\boldsymbol{\Theta} := [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r]^T \in \mathbb{R}^{r \times n}$.

Let $\mathbf{V} = \mathbf{L}_X^T \boldsymbol{\Theta}^T$ and consider the optimization problem in (13):

$$\begin{aligned}
& \sup_{\mathbf{f} \in \mathcal{A}_r} \{(1 - \lambda) \text{Dep}^{\text{emp}}(\mathbf{f}(X), Y) - \lambda \text{Dep}^{\text{emp}}(\mathbf{f}(X), S)\} \\
&= \sup_{\mathbf{f} \in \mathcal{A}_r} \frac{1}{n^2} \left\{ (1 - \lambda) \|\boldsymbol{\Theta} \mathbf{K}_X \mathbf{H} \mathbf{L}_Y\|_F^2 - \lambda \|\boldsymbol{\Theta} \mathbf{K}_X \mathbf{H} \mathbf{L}_S\|_F^2 \right\} \\
&= \sup_{\mathbf{f} \in \mathcal{A}_r} \frac{1}{n^2} \left\{ (1 - \lambda) \text{Tr} \{ \boldsymbol{\Theta} \mathbf{K}_X \mathbf{H} \mathbf{K}_Y \mathbf{H} \mathbf{K}_X \boldsymbol{\Theta}^T \} - \lambda \text{Tr} \{ \boldsymbol{\Theta} \mathbf{K}_X \mathbf{H} \mathbf{K}_S \mathbf{H} \mathbf{K}_X \boldsymbol{\Theta}^T \} \right\} \\
&= \max_{\mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{I}_r} \frac{1}{n^2} \text{Tr} \{ \boldsymbol{\Theta} \mathbf{L}_X \mathbf{B} \mathbf{L}_X^T \boldsymbol{\Theta}^T \} \\
&= \max_{\mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{I}_r} \frac{1}{n^2} \text{Tr} \{ \mathbf{V}^T \mathbf{B} \mathbf{V} \}
\end{aligned} \tag{28}$$

where the second step is due to (9) and

$$\begin{aligned}
\mathbf{B} &:= \mathbf{L}_X^T ((1 - \lambda) \mathbf{H} \mathbf{K}_Y \mathbf{H} - \lambda \mathbf{H} \mathbf{K}_S \mathbf{H}) \mathbf{L}_X \\
&= \mathbf{L}_X^T ((1 - \lambda) \tilde{\mathbf{K}}_Y - \lambda \tilde{\mathbf{K}}_S) \mathbf{L}_X.
\end{aligned}$$

It is shown in Kokiopoulou et al. (2011) that an³ optimizer of (28) is any matrix \mathbf{U} whose columns are eigenvectors corresponding to r largest eigenvalues of generalized problem

$$\mathbf{B} \mathbf{u} = \tau \mathbf{C} \mathbf{u} \tag{29}$$

and the maximum value is the summation of r largest eigenvalues. Once \mathbf{U} is determined, then, any $\boldsymbol{\Theta}$ in which $\mathbf{L}_X^T \boldsymbol{\Theta}^T = \mathbf{U}$ is optimal $\boldsymbol{\Theta}$ (denoted by $\boldsymbol{\Theta}^{\text{opt}}$). Note that $\boldsymbol{\Theta}^{\text{opt}}$ is not unique and has a general form of

$$\boldsymbol{\Theta}^T = (\mathbf{L}_X^T)^\dagger \mathbf{U} + \boldsymbol{\Lambda}_0, \quad \mathcal{R}(\boldsymbol{\Lambda}_0) \subseteq \mathcal{N}(\mathbf{L}_X^T).$$

However, setting $\boldsymbol{\Lambda}_0$ to zero would lead to minimum norm for $\boldsymbol{\Theta}$. Therefore, we opt $\boldsymbol{\Theta}^{\text{opt}} = \mathbf{U}^T \mathbf{L}_X^\dagger$. \square

Corollary 4.1. Embedding Dimensionality: A useful corollary of Theorem 4 is characterizing optimal embedding dimensionality as a function of trade-off parameter, λ :

$$r^{\text{Opt}}(\lambda) := \arg \sup_{0 \leq r \leq l} \left\{ \sup_{\mathbf{f} \in \mathcal{A}_r} \{J^{\text{emp}}(\mathbf{f}, \lambda)\} \right\} = \text{number of non-negative eigenvalues of (14)}$$

Proof. From proof of Theorem 4, we know that

$$\sup_{\mathbf{f} \in \mathcal{A}_r} \{(1 - \lambda) \text{Dep}^{\text{emp}}(\mathbf{f}(X), Y) - \lambda \text{Dep}^{\text{emp}}(\mathbf{f}(X), S)\} = \sum_{j=1}^r \tau_j,$$

where $\{\tau_1, \dots, \tau_n\}$ are eigenvalues of the generalized problem in (14) in decreasing order. It follows immediately that

$$\arg \sup_r \left\{ \sum_{j=1}^r \tau_j \right\} = \text{number of non-negative elements of } \{\tau_1, \dots, \tau_l\}.$$

\square

F Proof of Theorem 5

Theorem 5. Assume that k_S and k_Y are bounded by one and $f_j^2(\mathbf{x}_i) \leq M$ for any $j = 1, \dots, r$ and $i = 1, \dots, n$ for which $\mathbf{f} = (f_1, \dots, f_r) \in \mathcal{A}_r$. Then, for any $n > 1$ and $0 < \delta < 1$, with probability at least $1 - \delta$, we have

$$\left| \sup_{\mathbf{f} \in \mathcal{A}_r} J(\mathbf{f}, \lambda) - \sup_{\mathbf{f} \in \mathcal{A}_r} J^{\text{emp}}(\mathbf{f}, \lambda) \right| \leq rM \sqrt{\frac{\log(6/\delta)}{0.22^2 n}} + \mathcal{O}\left(\frac{1}{n}\right).$$

³Optimal \mathbf{V} is not unique.

Proof. Recall that in the proof of Lemma 1 we have shown that with probability at least $1 - \delta$, the following inequality holds

$$|\text{Dep}(Z, S) - \text{Dep}^{\text{emp}}(Z, S)| \leq \sqrt{\frac{r^2 M^2 \log(6/\sigma)}{0.22^2 n}} + \mathcal{O}\left(\frac{1}{n}\right).$$

Using the same reasoning for $\text{dep}(Z, Y)$, with probability at least $1 - \delta$, we have

$$|\text{Dep}(Z, Y) - \text{Dep}^{\text{emp}}(Z, Y)| \leq \sqrt{\frac{r^2 M^2 \log(6/\sigma)}{0.22^2 n}} + \mathcal{O}\left(\frac{1}{n}\right).$$

Since $J(\mathbf{f}(X)) = (1 - \lambda) \text{dep}(Z, Y) - \lambda \text{dep}(Z, S)$ and $J^{\text{emp}}(\mathbf{f}(X)) := (1 - \lambda) \text{Dep}^{\text{emp}}(Z, Y) - \lambda \text{Dep}^{\text{emp}}(Z, S)$, it follows that with probability at least $1 - \delta$,

$$|J(\mathbf{f}, \lambda) - J^{\text{emp}}(\mathbf{f}, \lambda)| \leq rM \sqrt{\frac{\log(6/\sigma)}{0.22^2 n}} + \mathcal{O}\left(\frac{1}{n}\right).$$

We complete the proof by noting that, the following inequality holds for any bounded J and J^{emp} :

$$\left| \sup_{\mathbf{f} \in \mathcal{A}_r} J(\mathbf{f}, \lambda) - \sup_{\mathbf{f} \in \mathcal{A}_r} J^{\text{emp}}(\mathbf{f}, \lambda) \right| \leq \sup_{\mathbf{f} \in \mathcal{A}_r} |J(\mathbf{f}, \lambda) - J^{\text{emp}}(\mathbf{f}, \lambda)|.$$

□

G Optimality of Target Task Performance in $\mathbf{K}-\mathcal{T}_{\text{Opt}}$

We show that maximizing $\text{dep}(\mathbf{f}(X), Y)$ can lead to a representation Z that is sufficient to result in the optimal Bayes prediction of Y .

Theorem 6. Let $\lambda = 0$, $\gamma \rightarrow 0$, \mathcal{H}_Y be a linear RKHS, $r \geq d_Y$, and \mathbf{f}^* be the optimal encoder obtained by (14). Then, there exists a linear regressor on top of the representation $Z^* = \mathbf{f}^*(X)$ that can perform as optimal as Bayes estimator:

$$\begin{aligned} \min_{\mathbf{W} \in \mathbb{R}^{d_Y \times r}, \mathbf{b} \in \mathbb{R}^{d_Y}} \mathbb{E}_{X,Y} [\|\mathbf{W}Z^* + \mathbf{b} - Y\|^2] &= \inf_{h \text{ is Borel}} \mathbb{E}_{X,Y} [\|h(X) - Y\|^2] \\ &= \mathbb{E}_{X,Y} [\|\mathbb{E}[Y|X] - Y\|^2]. \end{aligned}$$

Proof. We only prove this theorem for the empirical version due to its convergence to the population counterpart. The optimal Bayes estimator can be the composition of the kernelized encoder $Z = \mathbf{f}(X)$ and a linear regressor on top of it. More specifically, $\hat{Y} = \mathbf{W}\mathbf{f}(X) + \mathbf{b}$ can approach to $\mathbb{E}[Y|X]$ if we optimize \mathbf{f} , \mathbf{W} , and \mathbf{b} all together. This is because $\mathbf{f} \in \mathcal{H}_X$ can approximate any Borel function (due to the universality of \mathcal{H}_X) and, since $r \geq d_Y$, \mathbf{W} can be surjective. Let $\mathbf{Z} := [\mathbf{z}_1, \dots, \mathbf{z}_n] \in \mathbb{R}^{r \times n}$ and $\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{d_Y \times n}$. Further, let $\tilde{\mathbf{Z}}$ and $\tilde{\mathbf{y}}$ be the centered (i.e., mean subtracted) version of \mathbf{Z} and \mathbf{Y} , respectively. We firstly optimize \mathbf{b} for any given \mathbf{f} , r , and \mathbf{W} :

$$\begin{aligned} \mathbf{b}_{\text{opt}} &:= \arg \min_{\mathbf{b}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{W}\mathbf{z}_i + \mathbf{b} - \mathbf{y}_i\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i - \mathbf{W} \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i. \end{aligned}$$

Then, optimizing over \mathbf{W} would lead to

$$\begin{aligned} \min_{\mathbf{W}} \frac{1}{n} \|\mathbf{W}\tilde{\mathbf{Z}} - \tilde{\mathbf{Y}}\|_F^2 &= \frac{1}{n} \min_{\mathbf{W}} \|\tilde{\mathbf{Z}}^T \mathbf{W}^T - \tilde{\mathbf{Y}}^T\|_F^2 \\ &= \min_{\mathbf{W}} \frac{1}{n} \|\tilde{\mathbf{Z}}^T \mathbf{W}^T - P_{\tilde{\mathbf{Z}}} \tilde{\mathbf{Y}}^T\|_F^2 + \frac{1}{n} \|P_{\tilde{\mathbf{Z}}^\perp} \tilde{\mathbf{Y}}^T\|_F^2 \\ &= \frac{1}{n} \|P_{\tilde{\mathbf{Z}}^\perp} \tilde{\mathbf{Y}}^T\|_F^2 = \frac{1}{n} \|\tilde{\mathbf{Y}}\|_F^2 - \frac{1}{n} \|P_{\tilde{\mathbf{Z}}} \tilde{\mathbf{Y}}^T\|_F^2, \end{aligned}$$

where $P_{\tilde{Z}}$ denotes the orthogonal projector onto the column space of \tilde{Z}^T and a possible minimizer is $\mathbf{W}_{\text{opt}}^T = (\tilde{Z}^T)^\dagger \tilde{Y}^T$ or equivalently $\mathbf{W}_{\text{opt}} = \tilde{Y}(\tilde{Z})^\dagger$. Since the MSE loss is a function of the range (column space) of \tilde{Z}^T , we can consider only \tilde{Z}^T with orthonormal columns or equivalently $\frac{1}{n} \tilde{Z} \tilde{Z}^T = \mathbf{I}_r$. In this setting, it holds $P_{\tilde{Z}} = \frac{1}{n} \tilde{Z} \tilde{Z}^T$. Now, consider optimizing $f(X) = \Theta [k_X(\mathbf{x}_1, X), \dots, k_X(\mathbf{x}_n, X)]^T$. We have, $\tilde{Z} = \Theta K_X H$ where H is the centering matrix. Let $V = L_X^T \Theta^T$ and $C = \frac{1}{n} L_X^T H L_X$, then it follows that

$$\begin{aligned}
\min_{\Theta K_X H K_X \Theta^T = n \mathbf{I}_r} \frac{1}{n} \left\{ \|\tilde{Y}\|_F^2 - \|P_{\tilde{Z}} \tilde{Y}^T\|_F^2 \right\} &= \frac{1}{n} \|\tilde{Y}\|_F^2 - \max_{\Theta K_X H K_X \Theta^T = n \mathbf{I}_r} \frac{1}{n} \|P_{\tilde{Z}} \tilde{Y}^T\|_F^2 \\
&= \frac{1}{n} \|\tilde{Y}\|_F^2 - \max_{V^T C V = \mathbf{I}_r} \frac{1}{n^2} \text{Tr} [\tilde{Y} H K_X \Theta^T \Theta K_X H \tilde{Y}^T] \\
&= \frac{1}{n^2} \|\tilde{Y}\|_F^2 - \max_{V^T C V = \mathbf{I}_r} \frac{1}{n^2} \text{Tr} [\Theta K_X H \tilde{Y}^T \tilde{Y} H K_X \Theta^T] \\
&= \|\tilde{Y}\|_F^2 - \max_{V^T C V = \mathbf{I}_r} \frac{1}{n^2} \text{Tr} [V^T L_X^T \tilde{Y}^T \tilde{Y} L_X V] \\
&= \frac{1}{n} \|\tilde{Y}\|_F^2 - \frac{1}{n^2} \sum_{j=1}^r \lambda_j,
\end{aligned}$$

where $\lambda_1, \dots, \lambda_r$ are r largest eigenvalues of the following generalized problem

$$B_0 u = \lambda C u$$

and $B_0 := L_X^T \tilde{Y}^T \tilde{Y} L_X$. This resembles the eigenvalue problem in Section E, equation (29) where $\lambda = 0$, \mathcal{H}_Y is a linear RKHS and $\gamma \rightarrow 0$. \square

H Deficiency of Mean-Squared Error as A Measure of Dependence

Theorem. Let \mathcal{H}_S contain all Borel functions, S be a d_S -dimensional RV, and $L_S(\cdot, \cdot)$ be MSE loss. Then,

$$Z \in \arg \sup \left\{ \inf_{g_S \in \mathcal{H}_S} \mathbb{E}_{X,S} [L_S(g_S(Z), S)] \right\} \Leftrightarrow \mathbb{E}[S | Z] = \mathbb{E}[S].$$

Proof. Let S_i , $(g_S(Z))_i$, and $(\mathbb{E}[S | Z])_i$ denote the i -th entries of S , $g_S(Z)$, and $\mathbb{E}[S | Z]$, respectively. Then, it follows that

$$\begin{aligned}
\inf_{g_S \in \mathcal{H}_S} \mathbb{E}_{X,S} [L_S(g_S(Z), S)] &= \inf_{g_S \in \mathcal{H}_S} \sum_{i=1}^{d_S} \mathbb{E}_{X,S} [((g_S(Z))_i - S_i)^2] \\
&= \sum_{i=1}^{d_S} \mathbb{E}_{X,S} [((\mathbb{E}[S | Z])_i - S_i)^2] \\
&\leq \sum_{i=1}^{d_S} \mathbb{E}_S [((\mathbb{E}[S])_i - S_i)^2] = \sum_{i=1}^{d_S} \text{Var}[S_i],
\end{aligned}$$

where the second step is due to the optimality of conditional mean (i.e., Bayes estimation) for MSE (Jacod & Protter, 2012) and the last step is because independence between Z and S leads to an upper bound on MSE. Therefore, if $Z \in \arg \sup \{ \inf_{g_S \in \mathcal{H}_S} \mathbb{E}_{X,S} [L_S(g_S(Z), S)] \}$, then $\mathbb{E}[S | Z] = \mathbb{E}[S]$. On the other hand, if $\mathbb{E}[S | Z] = \mathbb{E}[S]$, then it follows immediately that $Z \in \arg \sup \{ \inf_{g_S \in \mathcal{H}_S} \mathbb{E}_{X,S} [L_S(g_S(Z), S)] \}$. \square

This theorem implies that an optimal adversary does not necessarily lead to a representation Z that is statistically independent of S , but rather leads to S being mean independent of the representation Z i.e., independence with respect to first order moment only.