
On the Fundamental Trade-offs in Learning Invariant Representations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Many applications of representation learning, such as privacy-preservation, al-
2 gorithmic fairness and domain adaptation, desire explicit control over semantic
3 information being discarded. This goal is often formulated as satisfying two po-
4 tentially competing objectives: maximizing utility for predicting a target attribute
5 while simultaneously being independent or invariant with respect to a known seman-
6 tic attribute. In this paper, we *identify and determine* two fundamental trade-offs
7 between utility and semantic dependence induced by the statistical dependencies
8 between the data and its corresponding target and semantic attributes. We derive
9 closed-form solutions for the global optima of the underlying optimization prob-
10 lems under mild assumptions, which in turn yields closed formulae for the exact
11 trade-offs. We also derive empirical estimates of the trade-offs and show their
12 convergence to the corresponding population counterparts. Finally, we numeri-
13 cally quantify the trade-offs on representative problems and compare the solutions
14 achieved by baseline representation learning algorithms.

15 1 Introduction

16 Real-world applications of representation learning algorithms often have to contend with objectives
17 beyond predictive performance. These include cost functions pertaining to, invariance (e.g., to
18 photometric or geometric variations), semantic independence (e.g., w.r.t to age or race for face
19 recognition systems), privacy (e.g., mitigating leakage of sensitive information [1]), algorithmic
20 fairness (e.g., demographic parity [2]), and generalization across multiple domains [3], to name a few.

21 At its core, the underlying goal of the aforementioned formulations of representation learning is to
22 satisfy two competing objectives, extracting as much information necessary to predict a target label
23 \mathbf{y} (e.g., face identity) while *intentionally and permanently* suppressing information pertaining to a
24 desired semantic attribute \mathbf{s} (e.g., age, gender or race). When \mathbf{y} is independent of \mathbf{s} , one can learn a
25 representation that is independent of \mathbf{s} with no loss of performance, i.e., no trade-off exists between
26 the two objectives. However, when the two attributes \mathbf{y} and \mathbf{s} are correlated, attaining semantic
27 independence will necessarily reduce the performance of the target predictor, i.e., there is a trade-off
28 between the two objectives. The trade-off is unknown yet is important for understanding the limits of
29 existing and future representation learning algorithms that involve semantic independence constraints.

30 Let $\mathbf{z} = f(\mathbf{x})$ be a representation of input data \mathbf{x} , and $f(\cdot)$ be the encoder (see Fig 1(a)). Invariant
31 learning requires that prediction of the target label, $\hat{\mathbf{y}} = g_Y(\mathbf{z})$ be independent of a semantic attribute
32 \mathbf{s} i.e., $\hat{\mathbf{y}} \perp\!\!\!\perp \mathbf{s}$ for all possible downstream target predictors $g_Y(\cdot)$. This independence condition is
33 satisfied if and only if (iff), the representation \mathbf{z} is independent of \mathbf{s} i.e., $\mathbf{z} \perp\!\!\!\perp \mathbf{s}$. Therefore, Invariant
34 representation learning (IRL) seeks to optimize two objectives: i) the degree of dependence between
35 data representation \mathbf{z} and semantic attribute \mathbf{s} , and ii) target task utility. These two objectives can be
36 combined into one, with a parameter τ controlling the trade-off.

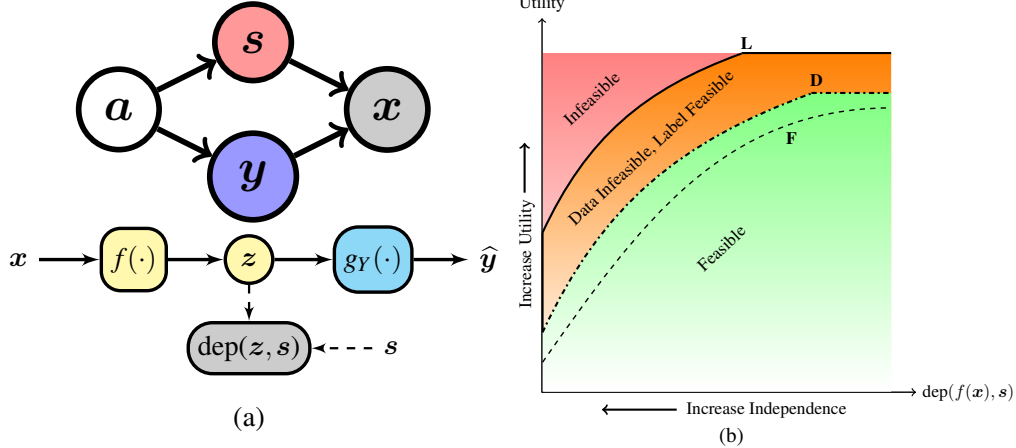


Figure 1: (a): Generic frame work of invariant representation learning (IRL) where attributes s and y are caused by a latent factor a and are not marginally independent. Under this setting, IRL seeks a representation $z = f(x)$ that contains enough information for downstream target predictor $g_Y(\cdot)$ while being independent of the semantic attribute s . Consequently, the prediction $\hat{y} = g_Y(z)$ will also be independent of s for any downstream predictor $g_Y(\cdot)$. (b): We identify and determine two different fundamental trade-offs between utility (i.e., the performance of target task predictor) and dependence measure $\text{dep}(z, s)$ by an optimal learner in the hypothesis class of Borel-measurable functions. Trade-off **L** is induced by the joint distribution of the labels $p_{y,s}$. Trade-off **D** is induced by the joint distribution of the data $p_{x,y,s}$. Trade-off **F** is a relaxed version of trade-off **D** obtained by either using a surrogate measure of dependence, e.g., adversarial learning [3] or from a constrained hypothesis class [4], or from using sub-optimal optimization algorithms.

37 In this paper, we identify and analytically determine two fundamental trade-offs in the invariant
 38 representation learning setting introduced above, namely *Data Space Trade-Off* and *Label Space*
 39 *Trade-Off*. These trade-offs are illustrated in Figure 1 (b) and formally defined next.

40 **Definition 1.** *Data Space Trade-Off* arises from the statistical dependence between the target attribute
 41 y and the semantic attribute s conditioned on the given input data x . When the learner’s hypothesis
 42 class contains all Borel-measurable functions¹ we have:

$$\inf_{f(\cdot) \text{ measurable}} \left\{ (1 - \tau) \inf_{g_Y(\cdot) \text{ measurable}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[\mathcal{L}_Y(g_Y(f(\mathbf{x})), \mathbf{y}) \right] + \tau \text{dep}(f(\mathbf{x}), s) \right\}. \quad (1)$$

43 where $f(\cdot)$ is the encoder that extracts representation z from x , $g_Y(\cdot)$ predicts \hat{y} from the repre-
 44 sentation z , $\mathcal{L}_Y(\cdot, \cdot)$ is the loss for the desired task of predicting the task label y . The function
 45 $\text{dep}(\cdot, \cdot) \geq 0$ is a parametric or non-parametric measure of statistical dependence i.e., $\text{dep}(\mathbf{q}, \mathbf{r}) = 0$
 46 means \mathbf{q} and \mathbf{r} are independent, and $\text{dep}(\mathbf{q}, \mathbf{r}) > 0$ means \mathbf{q} and \mathbf{r} are dependent with larger values
 47 indicating greater degrees of dependence. The scalar $\tau \in [0, 1)$ is a hyper-parameter that controls
 48 the trade-off between the two objectives, with $\tau = 0$ being the standard approach that enforces no
 49 independence to the attribute s , while $\tau \rightarrow 1$ enforces representation z to be independent of s .

50 Including all measurable functions in the hypothesis class of the encoder $f(\cdot)$ and target predic-
 51 tor $g_Y(\cdot)$ ensures that the best possible trade-off is included within the feasible solution space.
 52 For example, when $\tau = 0$ and $\mathcal{L}_Y(\cdot, \cdot)$ is the mean-squared error, the optimal Bayes estimator,
 53 $g_Y(f(\mathbf{x})) = \mathbb{E}_{\mathbf{y}}[y | \mathbf{x}]$ is reachable. This definition corresponds to the trade-off **D** in Figure 1 (b).

54 **Definition 2.** *Label Space Trade-Off* arises by ignoring the data x and is purely determined by the
 55 statistical dependence between the target feature y and the semantic attribute s . Such a trade-off can
 56 be defined as:

$$\inf_{z \in L^2} \left\{ (1 - \tau) \inf_{g_Y(\cdot) \text{ measurable}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[\mathcal{L}_Y(g_Y(z), \mathbf{y}) \right] + \tau \text{dep}(z, s) \right\}, \quad (2)$$

57 where L^2 is the space of all random vectors with finite second-order moment (i.e., $\mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|^2] < \infty$)
 58 on the same probability space in which the joint variable (s, y) comes from.

¹More specifically, we consider square-integrable Borel-measurable functions for boundedness.

59 This definition corresponds to the optimal trade-off obtained by an *ideal* representation \mathbf{z} that is not
60 constrained by the learnability of the encoder $f(\cdot)$. For example, if $\tau = 0$, the ideal representation
61 \mathbf{z} is perfectly aligned with the target label \mathbf{y} i.e., $\mathbf{z} = \mathbf{y}$ and $g_Y(\cdot)$ is the identity function, perfect
62 prediction of target attribute is feasible. Therefore, this trade-off corresponds to the best trade-off that
63 any combination of data \mathbf{x} and learnable encoder $f(\cdot)$ can aspire to. This definition corresponds to
64 the trade-off \mathbf{L} in Figure 1 (b), and it necessarily dominates the *Data Space Trade-Off* \mathbf{D} .

65 **Contributions:** i) Identify two fundamental trade-offs in invariant representation learning. ii) Obtain
66 closed-form solution for the corresponding optimization problems, and consequently determine the
67 trade-offs exactly. iii) Provide consistent empirical closed-form solution for the representations that
68 achieve optimal trade-offs. iv) Numerically quantify the trade-offs defined here and compare them to
69 those obtained by existing solutions.

70 **Implications:** i) Our closed-form empirical estimators for the optimal representations lend themselves
71 to practical invariant representation learning algorithms. ii) Theoretically elucidating and empirically
72 quantifying the intrinsic limits of invariant representations will enable researchers and practitioners
73 alike to identify the feasible and infeasible solution space for the trade-offs and lead to informed
74 development and deployment of optimal IRL methods. iii) Our theoretical analysis sheds light on the
75 utility-semantic independence trade-off, the role of statistical dependency between target label \mathbf{y} , the
76 semantic attribute \mathbf{s} , and the input data \mathbf{x} , and the hypothesis class adopted for the learners.

77 2 Related Work

78 **Trade-Offs in Representation Learning:** While there are abundant empirical approaches for the
79 representation learning applications considered in this paper, to the best of our knowledge, there
80 is no prior work that *exactly* characterizes and empirically quantifies the trade-offs inherent to
81 representation learning with semantic independence constraints.

82 Prior work primarily sought to either obtain lower or upper bounds or characterize the extreme
83 points of the trade-off in specific contexts such as fair representation learning. For instance, [5]
84 uses information theoretic tools and characterizes the utility-fairness trade-off in terms of a lower
85 bounds when both \mathbf{y} and \mathbf{s} are binary labels. Later [6] provided both upper and lower bound for the
86 binary labels. By leveraging Chernoff bound [7] proposed a construction method to generate an ideal
87 representation beyond input data to achieve perfect fairness while maintaining the best performance
88 on target task for equalized odds. In the case of categorical features, a lower bound on utility-fairness
89 trade-off has been provided by [8]. The notion of Pareto optimality was used by [9] to minimize
90 the maximum possible error among sensitive attributes where both target and sensitive features are
91 categorical. In contrast to this body of work, our trade-off analysis is applicable to multi-dimensional
92 discrete and/or continuous attributes where we find the exact optimal trade-offs.

93 The only prior work that investigates fundamental trade-offs in a general setting where both \mathbf{y} and \mathbf{s}
94 can be continuous or discrete features, are [4] and [10]. [4] considers only linear dependence between
95 the representation and semantic attribute and proposed a closed-form solution for the utility-fairness
96 trade-off. Even though [10] considers non-linear dependencies, optimal losses have been derived only
97 for the extremes of the trade-off (i.e., $\tau \rightarrow 0$ and $\tau \rightarrow 1$). In a more general setting where $0 < \tau < 1$,
98 [10] only provides a lower bound on utility-invariance trade-off through information plane analysis.
99 In contrast to the foregoing, we take a functional analysis approach and utilize covariance operator
100 based measures of dependence that account for all non-linear dependence relations. We exactly
101 characterize and quantify the utility-invariance trade-offs, while also providing a means to empirically
102 estimate the encoder that achieves said optimal trade-off. Lastly, in addition to the *Data Space*
103 *Trade-Off*, we also introduce and determine the *Label Space Trade-Off* which is the ideal trade-off
104 that any unrestricted learning algorithm can aspire to.

105 **Invariant, Fair, Privacy-Preserving Representation Learning:** The basic idea of representation
106 learning that discards unwanted semantic information has been explored under different contexts like
107 invariant, fair, or privacy-preserving learning. In domain adaptation [11, 12, 13], the goal is to learn
108 features that are independent of the data domain. In fair learning [14, 15, 16, 17, 18, 19, 20, 21, 22, 23,
109 2, 24, 25, 26, 27, 4], the goal is to discard the demographic information that leads to unfair outcomes.
110 Similarly, there is a growing interest in mitigating unintended leakage of private information from
111 data representations [28, 29, 1, 30, 31]. A vast majority of this body of work is empirical in nature.
112 These methods implicitly look for a single or more points in the trade-off between utility and fairness

113 and do not explicitly seek to characterize the whole trade-off front. Overall, these approaches are
 114 not concerned (or aware) about the feasibility and limitations on the utility-invariance trade-off. In
 115 contrast, this paper determines the fundamental theoretical limits of controlling independence to
 116 semantic attributes, and proposes practical learning algorithms that achieve this limit.

117 **Adversarial Representation Learning:** Most practical approaches for learning fair, invariant, do-
 118 main adaptive or privacy-preserving representations discussed above are based on adversarial repre-
 119 sentation learning (ARL). This learning problem is typically formulated as,

$$\inf_{f \in \mathcal{H}_x} \left\{ (1 - \tau) \inf_{g_Y \in \mathcal{H}_y} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[\mathcal{L}_Y \left(g_Y(f(\mathbf{x})), \mathbf{y} \right) \right] - \tau \inf_{g_S \in \mathcal{H}_s} \mathbb{E}_{\mathbf{x}, \mathbf{s}} \left[\mathcal{L}_S \left(g_S(f(\mathbf{x})), \mathbf{s} \right) \right] \right\}, \quad (3)$$

120 where $\mathcal{L}_S(\cdot, \cdot)$ is the loss function of a hypothetical adversary $g_S(\cdot)$ who intends to extract the semantic
 121 attribute \mathbf{s} through the best predictor within the hypothesis class \mathcal{H}_s . ARL is a special case of the *Data*
 122 *Space Trade-Off* in (1) where the negative loss of the adversary, $-\inf_{g_S \in \mathcal{H}_s} \mathbb{E}_{\mathbf{x}, \mathbf{s}} \left[\mathcal{L}_S \left(g_S(f(\mathbf{x})), \mathbf{s} \right) \right]$
 123 plays the role of $\text{dep}(f(\mathbf{x}), \mathbf{s})$. However, this form of adversarial learning suffers from a fundamental
 124 drawback as also noted in [32, 33]. The measure of dependence induced by ARL does not account
 125 for all modes of non-linear dependence between \mathbf{s} and the representation \mathbf{z} . The next theorem states
 126 this observation precisely,

127 **Theorem 1.**² Let \mathcal{H}_s contain all Borel-measurable functions and $\mathcal{L}_S(\cdot, \cdot)$ be mean squared error
 128 (MSE) loss. Then,

$$\mathbf{z} \in \arg \sup \left\{ \inf_{g_S \in \mathcal{H}_s} \mathbb{E}_{\mathbf{x}, \mathbf{s}} \left[\mathcal{L}_S \left(g_S(\mathbf{z}), \mathbf{s} \right) \right] \right\} \Leftrightarrow \mathbb{E}[\mathbf{s} | \mathbf{z}] = \mathbb{E}[\mathbf{s}].$$

129 This theorem implies that an optimal adversary does not necessarily lead to a representation \mathbf{z} that
 130 is statistically independent of \mathbf{s} (i.e., $p(\mathbf{s} | \mathbf{z}) = p(\mathbf{s})$), but rather leads to \mathbf{s} being mean independent
 131 of representation \mathbf{z} i.e., independence with respect to first order moment only. In other words,
 132 adversarially learned measure of dependence is not a complete measure of dependence and hence
 133 does not account for all modes of non-linear dependence between two random variables. As such, ARL
 134 is inherently incapable of attaining the trade-offs achievable by complete measures of dependence.

135 3 Theoretical Results

136 3.1 Problem Setting

137 Consider the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where Ω is the sample space, \mathcal{F} is a σ -algebra on Ω , and
 138 \mathbb{P} is a probability measure on \mathcal{F} . We assume that the joint random vector $(\mathbf{x}, \mathbf{y}, \mathbf{s})$, containing the
 139 input data $\mathbf{x} \in \mathbb{R}^{d_x}$, the target label $\mathbf{y} \in \mathbb{R}^{d_y}$ and the sensitive attribute $\mathbf{s} \in \mathbb{R}^{d_s}$, is a random vector
 140 on (Ω, \mathcal{F}) with joint distribution $\mathbf{p}_{\mathbf{x}\mathbf{y}\mathbf{s}}$.

141 **Assumption 1.** We assume that the encoder consists of r functions in an L_2 -universal RKHS
 142 $(\mathcal{H}_x, k_x(\cdot, \cdot))$ (e.g., Gaussian kernel), where L_2 -universality guarantees that \mathcal{H}_x can approximate
 143 any Borel-measurable function with arbitrary precision [34].

144 Now, the representation vector \mathbf{z} can be expressed as

$$\mathbf{z} = \mathbf{f}(\mathbf{x}) := \left[f_1(\mathbf{x}), \dots, f_r(\mathbf{x}) \right]^T \in \mathbb{R}^r, \quad f_j(\cdot) \in \mathcal{H}_x \quad \forall j = 1, \dots, r. \quad (4)$$

145 where r is the dimensionality of the representation \mathbf{z} . As discussed in Corollary 5.1, unlike common
 146 practice where it is chosen arbitrarily, r itself is an object of interest for optimization. We consider a
 147 general scenario where both \mathbf{y} and \mathbf{s} can be continuous or discrete, or one of \mathbf{y} or \mathbf{s} is continuous
 148 while the other is discrete. To do this, we substitute³ the target loss, $\inf_{g_Y} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[\mathcal{L}_Y(g_Y(\mathbf{z}), \mathbf{y}) \right]$ in (1)
 149 with the negative of a non-parametric measure of dependence i.e., $-\text{dep}(\mathbf{z}, \mathbf{y})$. Furthermore, in

²We defer the proofs of all lemmas, theorems and corollaries to the supplementary material.

³Many standard loss functions can be written in term of dependence measures [35] that capture all non-linear dependencies i.e., $\mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[\mathcal{L}_Y(f_T(\mathbf{f}(\mathbf{x})), \mathbf{y}) \right] \propto -\text{dep}(f(\mathbf{x}), \mathbf{y})$. For example, the mean squared error is proportional to $1 - \rho(f(\mathbf{x}), \mathbf{y})$, where ρ is the Pearson correlation coefficient, a plausible dependence measure.

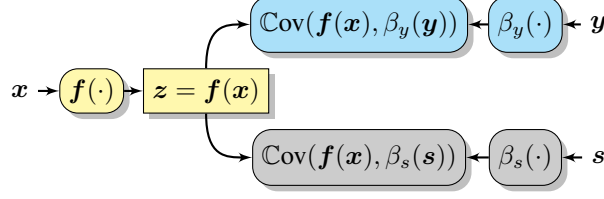


Figure 2: Our IRL model consists of three components: i) An r -dimensional encoder $f(\cdot)$ in a RKHS \mathcal{H}_x . ii) A measure of dependence that accounts for all kinds of linear or non-linear dependencies between the representation z and the semantic attribute s via the covariance between $f(x)$ and $\beta_s(s)$ where x is the input data and $\beta_s(\cdot)$ belongs to RKHS \mathcal{H}_s . iii) A measure of dependency between $f(x)$ and the target attribute y defined similar to the one for s .

150 unsupervised settings, when there is no target attribute y , the target dependence $\text{dep}(z, y)$ can be
 151 replaced with $\text{dep}(z, x)$, which implicitly forces the representation z to retain as much information
 152 as is necessary for reconstructing the input data x . This scenario is of practical interest when a data
 153 producer aims to provide a representation of data that is independent of a desired semantic attribute
 154 for any arbitrary downstream task.

155 We start by designing $\text{dep}(z, s)$, and $\text{dep}(z, y)$ follows similarly. A key desiderata of dependence
 156 measures is that they should be able to account for all possible non-linear dependence relations
 157 between the random variables (or vectors). Examples of such measures include information theoretic
 158 measures such as mutual information (e.g., MINE [36]) or covariance operator based measures such
 159 as Hilbert-Schmidt Independence Criterion [37], Constrained Covariance [38] and Kernel Canonical
 160 Correlation [39]. The underlying principle behind the latter class of dependence measures is that
 161 finite dimensional spaces with non-linear dependencies behave as linearly dependent spaces when
 162 mapped appropriately to higher dimensional spaces. In this paper we adopt the covariance operator
 163 based measures as our choice of dependence measure for analytical tractability.

164 Principally, z and s are independent iff $\text{Cov}(\alpha(z), \beta_s(s))$ is zero for all $\alpha(\cdot)$ and $\beta_s(\cdot)$ belong-
 165 ing to some universal RKHSs [38]. Since $z = f(x)$ and $f(\cdot) \in \mathcal{H}_x$, $\text{Cov}(\alpha(z), \beta_s(s)) =$
 166 $\text{Cov}(\alpha(f(x)), \beta_s(s))$, which necessitates application of a kernel on top of another kernel. This
 167 limits the analytical tractability of our solution. However, as we argue below, it is almost sufficient to
 168 consider transformation on s , only, in which case it reduces to $\text{Cov}(f(x), \beta_s(s))$. Let $(\mathcal{H}_s, k_s(\cdot, \cdot))$
 169 and $(\mathcal{H}_y, k_y(\cdot, \cdot))$ be separable⁴ RKHSs of functions defined on \mathbb{R}^{d_s} and \mathbb{R}^{d_y} , respectively. Consider
 170 the bi-linear functional,

$$h(\cdot, \cdot) : \mathcal{H}_x \times \mathcal{H}_s \rightarrow \mathbb{R}, h_j(f_j, \beta_s) := \text{Cov}_{x,s}(f_j(x), \beta_s(s)). \quad (5)$$

171 **Assumption 2.** We assume in the rest of this paper that the positive definite kernel functions are
 172 bounded, i.e.,

$$\mathbb{E}_x[k_x(x, x)] < \infty, \quad \mathbb{E}_s[k_s(s, s)] < \infty, \quad \text{and} \quad \mathbb{E}_y[k_y(y, y)] < \infty. \quad (6)$$

173 The assumptions in (6) guarantee that $h(\cdot, \cdot)$ in (5) is bounded [40] and therefore, invoking Riesz
 174 representation theorem [41], there exists a unique and bounded linear operator Σ_{sx} , such that

$$h(f, \beta_s) = \text{Cov}_{x,s}(f(x), \beta_s(s)) = \langle \beta_s, \Sigma_{sx} f \rangle_{\mathcal{H}_s} \quad \forall f \in \mathcal{H}_x, \forall \beta_s \in \mathcal{H}_s. \quad (7)$$

175 Based on $h(\cdot, \cdot)$, we define the linear operator $h_{f,s} : \mathcal{H}_s \rightarrow \mathbb{R}^r$ as

$$h_{f,s}(\beta_s) := \begin{bmatrix} \text{Cov}_{x,s}(f_1(x), \beta_s(s)) \\ \vdots \\ \text{Cov}_{x,s}(f_r(x), \beta_s(s)) \end{bmatrix} = \begin{bmatrix} \langle \beta_s, \Sigma_{sx} f_1 \rangle_{\mathcal{H}_s} \\ \vdots \\ \langle \beta_s, \Sigma_{sx} f_r \rangle_{\mathcal{H}_s} \end{bmatrix}.$$

176 The operator $h_{f,s}$ captures all modes of non-linear dependence, since the distribution of a low-
 177 dimensional projection of high-dimensional data is approximately normal [42], [43]. In other words,
 178 we assume that $(f(x), \beta_s(s))$ is an approximately Gaussian random vector.

⁴By separable we mean having a countable orthonormal basis set.

179 Among the different dependence measures that have been defined through the covariance operator
 180 we adopt the Hilbert-Schmidt Independence Criterion (HSIC) [37] which is defined as the Hilbert-
 181 Schmidt norm (HS-norm) of the covariance operator,

$$\text{dep}(z, \mathbf{s}) := \|\mathbf{h}_{f, \mathbf{s}}\|_{\text{HS}}^2 = \sum_{\beta_s \in \mathcal{U}_s} \|\mathbf{h}_{f, \mathbf{s}}(\beta_s)\|_2^2 = \sum_{\beta_s \in \mathcal{U}_s} \sum_{j=1}^r h^2(f_j, \beta_s) \quad (8)$$

182 where \mathcal{U}_s is a countable orthonormal basis set for \mathcal{H}_s . Note that, based on this definition, if the
 183 distribution $(\mathbf{f}(\mathbf{x}), \beta_s(\mathbf{s}))$ fails to be a normal distribution, we end up measuring mean dependency
 184 of $z = \mathbf{f}(\mathbf{x})$ from \mathbf{s} which is still much stronger than the linear dependency between z and \mathbf{s} [44].
 185 Even under this assumption, empirically (Section 4) we observe that trade-offs we obtain significantly
 186 dominate those from existing invariant representation learning algorithms.

187 The following Lemma introduces a well-defined population expression for $\text{dep}(z, \mathbf{s})$ in (8).

Lemma 2.

$$\begin{aligned} \text{dep}(z, \mathbf{s}) = & \sum_{j=1}^r \left\{ \mathbb{E}_{\mathbf{x}, \mathbf{s}, \mathbf{x}', \mathbf{s}'} \left[f_j(\mathbf{x}) f_j(\mathbf{x}') k_s(\mathbf{s}, \mathbf{s}') \right] + \mathbb{E}_{\mathbf{x}}[f_j(\mathbf{x})] \mathbb{E}_{\mathbf{x}'}[f_j(\mathbf{x}')] \mathbb{E}_{\mathbf{s}, \mathbf{s}'}[k_s(\mathbf{s}, \mathbf{s}')] \right. \\ & \left. - 2 \mathbb{E}_{\mathbf{x}, \mathbf{s}} \left[f_j(\mathbf{x}) \mathbb{E}_{\mathbf{x}'}[f_j(\mathbf{x}')] \mathbb{E}_{\mathbf{y}'}[k_s(\mathbf{s}, \mathbf{s}')] \right] \right\} \end{aligned}$$

188 where (\mathbf{x}, \mathbf{s}) and $(\mathbf{x}', \mathbf{s}')$ are independently drawn from the joint distribution $\mathbf{p}_{\mathbf{x}\mathbf{s}}$.

189 In practice, it is necessary to empirically estimate $\text{dep}(z, \mathbf{s})$, since the population distributions are
 190 typically unknown in most real-world scenarios.

191 **Definition 3.** Let $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{s}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{s}_n, \mathbf{y}_n)\}$ be the training data, containing n i.i.d.
 192 realizations from the joint distribution $\mathbf{p}_{\mathbf{x}\mathbf{s}\mathbf{y}}$. Using, the representer theorem [45], it follows that
 193 $\mathbf{f}(\mathbf{x}) = \Theta_E [k_{\mathbf{x}}(\mathbf{x}_1, \mathbf{x}), \dots, k_{\mathbf{x}}(\mathbf{x}_n, \mathbf{x})]^T$, where $\Theta \in \mathbb{R}^{r \times n}$ is a free parameter matrix.

194 **Lemma 3.** Let an empirical estimation of covariance be

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

195 Then, the empirical estimator of $\text{dep}(z, \mathbf{s})$ is given by

$$\text{dep}^{\text{emp}}(z, \mathbf{s}) := \frac{1}{n^2} \|\Theta \mathbf{K}_{\mathbf{x}} \mathbf{H} \mathbf{L}_{\mathbf{s}}\|_F^2, \quad (9)$$

196 where $\mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{s}} \in \mathbb{R}^{n \times n}$ are Gram matrices corresponding to $\mathcal{H}_{\mathbf{x}}$ and $\mathcal{H}_{\mathbf{s}}$, respectively, $\mathbf{H} =$
 197 $\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T$, and $\mathbf{L}_{\mathbf{s}}$ is a full column-rank matrix in which $\mathbf{L}_{\mathbf{s}} \mathbf{L}_{\mathbf{s}}^T = \mathbf{K}_{\mathbf{s}}$ (Cholesky factorization).
 198 This empirical estimator in (9) has a bias of $\mathcal{O}(n^{-1})$ and a convergence rate of $\mathcal{O}(n^{-1/2})$.

199 The population and empirical dependence measures between z and \mathbf{y} i.e., $\text{dep}(z, \mathbf{y})$ and $\text{dep}^{\text{emp}}(z, \mathbf{y})$,
 200 respectively, can be defined and obtained similarly.

201 3.2 Trade-Off D

202 We now turn to the the optimization problem corresponding to the trade-off **D** in (1). Recall that
 203 $z = \mathbf{f}(\mathbf{x})$ is r -dimensional, where the dimensionality r is a free variable. A common desiderata of
 204 learned representations is that of compactness [46] in order to avoid learning representations with
 205 redundant information where different dimensions are highly correlated with each other. Therefore,
 206 going beyond the assumption that each component of $\mathbf{f}(\cdot)$ (i.e., $f_j(\cdot)$) belongs to a L_2 -universal
 207 RKHS $\mathcal{H}_{\mathbf{x}}$, we impose additional constraints on the representation. Specifically, we constrain the
 208 search space of the encoder $\mathbf{f}(\cdot)$ to learn a disentangled representation [46] as follows,

$$\mathcal{A}_r := \left\{ \left(f_1(\cdot), \dots, f_r(\cdot) \right) \mid f_i, f_j \in \mathcal{H}_{\mathbf{x}}, \text{Cov}_{\mathbf{x}}(f_i(\mathbf{x}), f_j(\mathbf{x})) + \gamma \langle f_i, f_j \rangle_{\mathcal{H}_{\mathbf{x}}} = \delta_{i,j} \right\}, \quad (10)$$

209 where the regularization term $\gamma \langle f_i, f_j \rangle_{\mathcal{H}_{\mathbf{x}}}$, encourages orthogonality and boundedness, which in turn
 210 forces the representation to be compact or non-redundant. Such disentangled representations have

211 been studied in the context of independent component analysis (ICA) [39]. Now, the optimization
 212 problem in (1) reduces to,

$$\sup_{\mathbf{f} \in \mathcal{A}_r} \left\{ J(\mathbf{f}(\mathbf{x})) := (1 - \tau) \text{dep}(\mathbf{f}(\mathbf{x}), \mathbf{y}) - \tau \text{dep}(\mathbf{f}(\mathbf{x}), \mathbf{s}) \right\}, \quad 0 \leq \tau < 1, \quad (11)$$

213 where as justified earlier the target loss function $\inf_{f_Y} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\mathcal{L}_Y(f_T(\mathbf{f}(\mathbf{x})), \mathbf{y})]$ is substituted by
 214 $-\text{dep}(\mathbf{f}(\mathbf{x}), \mathbf{y})$. Fortunately, the above optimization problem lends itself to a closed-form solu-
 215 tion as given by the next theorem.

216 **Theorem 4.** A solution⁵ to the optimization problem in (11) is the eigenfunctions corresponding to r
 217 largest eigenvalues of the following generalized problem

$$\left((1 - \tau) \Sigma_{\mathbf{y}\mathbf{x}}^* \Sigma_{\mathbf{y}\mathbf{x}} - \tau \Sigma_{\mathbf{s}\mathbf{x}}^* \Sigma_{\mathbf{s}\mathbf{x}} \right) \mathbf{f} = \lambda \Sigma_{\mathbf{x}\mathbf{x}} \mathbf{f}, \quad (12)$$

218 where $\Sigma_{\mathbf{s}\mathbf{x}}$ and $\Sigma_{\mathbf{y}\mathbf{x}}$ are the covariance operators defined in (7), and $\Sigma_{\mathbf{s}\mathbf{x}}^*$ and $\Sigma_{\mathbf{y}\mathbf{x}}^*$ are the adjoint
 219 operators of $\Sigma_{\mathbf{s}\mathbf{x}}$ and $\Sigma_{\mathbf{y}\mathbf{x}}$, respectively.

220 **Remark.** If the trade-off parameter $\tau = 0$ (i.e., no semantic independence constraint is imposed), the
 221 solution in Theorem 4 resembles a supervised version of ICA in [39] which is essentially a kernelized
 222 dimensionality reduction supervised by the target attribute \mathbf{y} . On the other hand, if $\tau \rightarrow 1$ (i.e.,
 223 utility is ignored and only semantic independence is considered), the solution in Theorem 4 is the
 224 eigenfunctions corresponding to the negative eigenvalues of $\Sigma_{\mathbf{s}\mathbf{x}}^* \Sigma_{\mathbf{s}\mathbf{x}}$, which are the directions that
 225 are least explanatory of the semantic attribute \mathbf{s} .

226 An empirical version of (11) is the following optimization problem

$$\sup_{\mathbf{f} \in \mathcal{A}_r} \left\{ J^{\text{emp}}(\mathbf{f}(\mathbf{x})) := (1 - \tau) \text{dep}^{\text{emp}}(\mathbf{f}(\mathbf{x}), \mathbf{y}) - \tau \text{dep}^{\text{emp}}(\mathbf{f}(\mathbf{x}), \mathbf{s}) \right\}, \quad 0 \leq \tau < 1 \quad (13)$$

227 where $\text{dep}^{\text{emp}}(\mathbf{f}(\mathbf{x}), \mathbf{s})$ and $\text{dep}^{\text{emp}}(\mathbf{f}(\mathbf{x}), \mathbf{y})$ are given in (9).

228 **Theorem 5.** Consider the Cholesky factorization $\mathbf{K}_x = \mathbf{L}_x \mathbf{L}_x^T$, where \mathbf{L}_x is a full column-rank
 229 matrix. A solution to (13) is

$$\mathbf{f}^{\text{opt}} = \Theta^{\text{opt}} \left[k_x(\mathbf{x}_1, \cdot), \dots, k_x(\mathbf{x}_n, \cdot) \right]^T$$

230 where $\Theta^{\text{opt}} = \mathbf{U}^T (\mathbf{L}_x)^\dagger$ and the columns of \mathbf{U} are eigenvectors corresponding to r largest eigenval-
 231 ues, $\lambda_1, \dots, \lambda_r$, of the following generalized problem,

$$\left(\mathbf{L}_x^T ((1 - \tau) \tilde{\mathbf{K}}_y - \tau \tilde{\mathbf{K}}_s) \mathbf{L}_x \right) \mathbf{u} = \lambda \left(\mathbf{L}_x^T \mathbf{H} \mathbf{L}_x + n\gamma \mathbf{I} \right) \mathbf{u} \quad (14)$$

232 where γ is the regularization parameter from (10) and the supremum value of (13) is $\sum_{j=1}^r \lambda_j$.

233 **Corollary 5.1. Embedding Dimensionality:** A useful corollary of Theorem 5 is optimal embedding
 234 dimensionality:

$$\arg \sup_r \left\{ \sup_{\mathbf{f} \in \mathcal{A}_r} \left\{ J^{\text{emp}}(\mathbf{f}(\mathbf{x})) := (1 - \tau) \text{dep}^{\text{emp}}(\mathbf{f}(\mathbf{x}), \mathbf{y}) - \tau \text{dep}^{\text{emp}}(\mathbf{f}(\mathbf{x}), \mathbf{s}) \right\} \right\},$$

235 which is the number of positive eigenvalues of the generalized eigenvalue problem in (14). To
 236 intuitively examine this result, consider two extreme cases: i) If there is no semantic independence
 237 constraint (i.e., $\tau = 0$), adding more dimensions to the optimum r will not harm the representation
 238 power of \mathbf{z} . ii) If we only care about semantic independence and ignore the target task (i.e., $\tau \rightarrow 1$),
 239 the optimal r would be equal to zero, indicating that a null representation is the best for discarding all
 240 semantic information. In this case, adding more dimension to \mathbf{z} will necessarily violate the semantic
 241 independence constraint. More discussion can be found in the supplementary material.

242 In the following Theorem, we prove that the empirical solution converges to its population counterpart.

243 **Theorem 6.** Assume that $k_s(\cdot, \cdot)$ and $k_y(\cdot, \cdot)$ are bounded by one and $f_k^2(\mathbf{x}_i)$ is bounded by M for
 244 any $k = 1, \dots, r$ and $i = 1, \dots, n$ for which $\mathbf{f} = (f_1, \dots, f_r) \in \mathcal{A}_r$. For any $n > 1$ and $0 < \delta < 1$,
 245 with probability at least $1 - \delta$, we have

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

246 where $0.22 \leq a \leq 1$ is a constant.

⁵The term 'solution' in any optimization problem in this paper refers to a global optima.

247 **3.3 Trade-Off L**

248 We recall that label space trade-off arises when the representation z is ideal and is free to be designed
 249 optimally i.e., it does not necessarily depend on the input data x or the encoder’s hypothesis class.
 250 However, we assume that the representation z is a direct effect of the target and sensitive variables (y
 251 and s). Following [47], we use an additive noise model as

$$z = f_L(y, s) + e, \quad e \perp\!\!\!\perp y, e \perp\!\!\!\perp s \quad (15)$$

252 where $f_L(\cdot, \cdot) : \mathbb{R}^{d_y} \times \mathbb{R}^{d_s} \rightarrow \mathbb{R}^r$ is a Borel-measurable function. Following Section 3.1,
 253 we deploy $-\text{dep}(z, y)$, defined similar to $\text{dep}(z, s)$ in (8), as a proxy for the loss function
 254 $\inf_{g_Y \in \mathcal{H}_y} \mathbb{E}_{x, y} [\mathcal{L}_T(g_T(z), y)]$. Recall that, the desired optimization problem is given in (2). Instead
 255 of directly optimizing over $z \in L^2$, we optimize over all Borel-measurable functions $f_L(\cdot, \cdot)$ by
 256 ignoring e since it is independent of both y and s :

$$\sup_{f_L \in \mathcal{A}_r(y, s)} \left\{ (1 - \tau) \text{dep}(f_L(y, s), y) - \tau \text{dep}(f_L(y, s), s) \right\}, \quad (16)$$

257 where $\mathcal{A}_r(y, s)$ is defined similar to \mathcal{A}_r in (10) by using (y, s) instead of x in the definition. Recall
 258 that $\mathcal{A}_r(y, s)$ ensures that z will not contain highly correlated (entangled) dimensions, and thus be
 259 minimally redundant or maximally compact.

260 **Remark.** The optimization problem in (16) and its empirical counterpart can be solved similar to
 261 that of trade-off **D** in Theorems 5 and 6 where x is replaced with (y, s) .

262 **3.4 Trade-Off F**

263 Here we define and discuss the trade-off achievable by practical realizations of representation learning
 264 algorithms with either fairness, invariance or semantic independence constraints.

265 **Definition 4.** *Feasible Space Trade-Off* arises from the statistical dependence between the target
 266 feature y and the sensitive attribute s conditioned on the given input data x , the choice of hypothesis
 267 class for the learners involved, and the choice of dependence measure adopted. This setting can be
 268 formalized as,

$$\inf_{f \in \mathcal{H}_x} \left\{ (1 - \tau) \inf_{g_Y \in \mathcal{H}_y} \mathbb{E}_{x, y} \left[\mathcal{L}_Y(g_Y(f(x)), y) \right] + \tau \widetilde{\text{dep}}(f(x), s) \right\}, \quad 0 \leq \tau < 1, \quad (17)$$

269 where \mathcal{H}_x and \mathcal{H}_y are the hypothesis class for the encoder network and target predictor, respectively,
 270 $\mathcal{L}_Y(\cdot, \cdot)$ denotes the loss function of target task, and $\widetilde{\text{dep}}(f(x), s)$ is a parametric or non-parametric
 271 surrogate measure of dependency quantifying the dependency between representation vector $z =$
 272 $f(x)$ and the sensitive attribute s .

273 This setting corresponds to the trade-off **F** in Figure 1(b), and is necessarily dominated by the
 274 Data Space Trade-Off **D**. Multiple factors may lead to such sub-optimal trade-offs. These include,
 275 hypothesis classes that are not universal RKHSs (e.g., [4] considered the case where \mathcal{H}_x is universal,
 276 but \mathcal{H}_s and \mathcal{H}_y are linear RKHSs), the surrogate dependence measure $\widetilde{\text{dep}}(f(x), s)$ does not account
 277 for all non-linear dependencies (e.g., [3, 2, 21, 4] which consider adversarially learned dependence
 278 measures), sub-optimal optimization of (17) in terms of achieving only local optima but not the
 279 global optima (e.g., when the hypothesis class is deep neural networks that are optimized through
 280 stochastic gradient descent, or through stochastic gradient descent-ascent in the case of adversarial
 281 representation learning[3, 21, 2]), and combinations thereof.

282 **4 Numerical Estimation of Trade-Offs**

283 In this section, we demonstrate the practical utility of the analytical results developed in the paper
 284 and validate our theoretical insights. For this purpose, we design an illustrative toy example that
 285 conforms to the setting studied in the paper and numerically quantify the trade-offs that we introduced.
 286 Experimental validation on more tasks can be found in the supplementary material.

287 Consider the following Gaussian mixture model from which we generate 4000, 2000, and 2000

$$v = [v_1, v_2] \sim \frac{1}{2} \left(\mathcal{N}(m, \Sigma) + \mathcal{N}(m', \Sigma) \right), \quad m = [0, 1], \quad m' = [1, 1], \quad \Sigma = \text{diag}(0.1^2, 0.1^2)$$

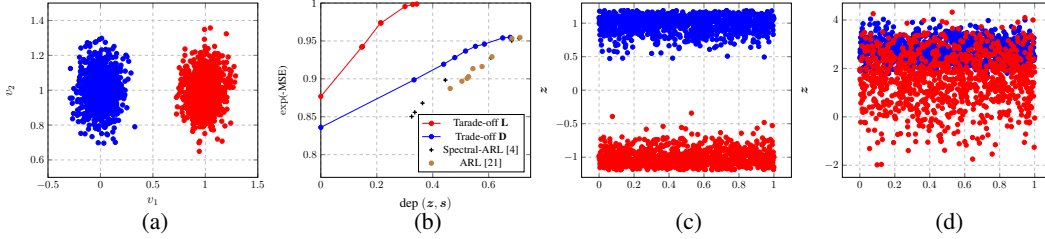


Figure 3: (a): A mixture of two Gaussians which generates the input data as $\mathbf{x} = v_1$, the sensitive attribute as $\mathbf{s} = v_1^3$, and the target attribute as $\mathbf{y} = [v_1, v_2^3]$. (b): Two fundamental trade-offs, \mathbf{L} and \mathbf{D} , together with two baseline feasible trade-offs \mathbf{F} , ARL optimized with SGDA [21] and global optima of ARL with a linear RKHS [4]. (c), (d): The learned embedding for $\tau = 0$ and $\tau = 0.5$, respectively. An invariant representation should collapse v_1 i.e., the two colors should fully overlap with each other in the embedding. The overlap is partial for $\tau = 0.5$ and as $\tau \rightarrow 1$, the optimal representation is zero.

288 independent samples for training, validation and testing, respectively. Figure 3(a) shows the test
 289 samples where the samples generated with m and m' are in blue and red, respectively. The input data
 290 \mathbf{x} is set to v_1 (the first entry of \mathbf{v}), the sensitive attribute \mathbf{s} is v_1^3 , and the target attribute \mathbf{y} is $[v_1, v_2^3]$.
 291 In this problem both input data and target attribute are dependent on the sensitive attribute. We choose
 292 all three RKHS \mathcal{H}_x , \mathcal{H}_y , and \mathcal{H}_s to be Gaussian, which is a universal RKHS. The optimal \mathbf{z} is learned
 293 for the trade-off \mathbf{D} through the closed-form solution in Theorem 5 for different invariance parameter
 294 values τ in $[0, 1)$. Then, this optimal embedding is fed to a target task predictor which is a multi-layer
 295 perceptron (MLP) with two hidden layers, and 4, 8 neurons and optimize the mean-squared error
 296 (MSE). The x-axis is a normalized version of the dependence measure used in our optimization, while
 297 the y-axis quantifies utility normalized to $[0, 1]$ as $\exp(-\text{MSE})$. The same procedure is implemented
 298 for trade-off \mathbf{L} , except that the input data is \mathbf{v} , instead of \mathbf{x} . These trade-offs are shown in Figure 3(b).
 299 We choose the input data to be \mathbf{v} instead of (\mathbf{y}, \mathbf{s}) for trade-off \mathbf{L} since (\mathbf{y}, \mathbf{s}) is fully generated from
 300 \mathbf{v} and therefore, \mathbf{v} perfectly explains (\mathbf{y}, \mathbf{s}) . For $\tau = 0$ and $\tau = 0.5$, the optimal embeddings are
 301 illustrated in Figure 3, (c) and (d), respectively. Since the sensitive attribute is only related to v_1 ,
 302 an invariant embedding should collapse the corresponding dimension and cause the two colors to
 303 overlap with each other.

304 We make the following observations, (a) Trade-off \mathbf{L} dominates trade-off \mathbf{D} as expected. (b) The
 305 trade-offs \mathbf{F} obtained by the baselines are dominated by trade-off \mathbf{D} . Adversarial representation
 306 learning [3, 21, 2] uses sub-optimal optimization (SGDA), while Spectral-ARL [4] uses a global
 307 optimum solution but restricts the hypothesis class in (3) to linear RKHS. As such, the baselines are
 308 unable to match the global optimal solution of (13), and (c) At $\tau = 0.5$ the embedding does indeed
 309 collapse v_1 to an extent leading to partial overlap between the two mixtures.

310 5 Conclusions and Societal Impact

311 This paper developed the theoretical underpinnings for identifying and determining the fundamental
 312 trade-offs and limits of representation learning under competing objectives. These trade-offs included
 313 i) label space trade-off which is solely induced by the statistical relation between target task and
 314 semantic attribute; ii) data space trade-off which is due to the statistical dependence between the
 315 input data and both target and semantic attributes. Further, we found closed-form solutions for the
 316 global optima, both the population and empirical versions, for the underlying optimization problems,
 317 and thus quantify the trade-offs *exactly*. Our results shed light on the regions of the trade-off that are
 318 feasible or impossible to achieve by learning algorithms. Numerical results suggest that commonly
 319 used adversarial representation learning based techniques are unable to reach the optimal trade-offs.

320 The theoretical results in this paper are useful for algorithmic fairness, privacy-preservation, and
 321 domain generalization applications of representation learning. Such systems are being widely
 322 deployed in a variety of practical applications: search engines, social media, law enforcement,
 323 healthcare, consumer devices, financial and judicial risk assessments, face analysis, and many more.
 324 Therefore, providing theoretical limits of performance is critically important for informed framing
 325 of regulatory policies, deployment of such solutions, and gaining societal trust. As such, we do not
 326 anticipate any adverse societal impacts from this work.

327 **References**

- 328 [1] P. Roy and V. N. Boddeti, “Mitigating information leakage in image representations: A maximum entropy approach,” in *IEEE Conference on Computer Vision and Pattern Recognition*,
 329 2019.
 330
- 331 [2] D. Madras, E. Creager, T. Pitassi, and R. Zemel, “Learning adversarially fair and transferable
 332 representations,” *arXiv preprint arXiv:1802.06309*, 2018.
- 333 [3] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand,
 334 and V. Lempitsky, “Domain-adversarial training of neural networks,” *The Journal of Machine*
 335 *Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- 336 [4] B. Sadeghi, R. Yu, and V. Boddeti, “On the global optima of kernelized adversarial representation
 337 learning,” in *IEEE International Conference on Computer Vision*, pp. 7971–7979, 2019.
- 338 [5] H. Zhao and G. J. Gordon, “Inherent tradeoffs in learning fair representations,” *arXiv preprint*
 339 *arXiv:1906.08386*, 2019.
- 340 [6] D. McNamara, C. S. Ong, and R. C. Williamson, “Costs and benefits of fair representation
 341 learning,” in *AAAI/ACM Conference on AI, Ethics, and Society*, pp. 263–270, 2019.
- 342 [7] S. Dutta, D. Wei, H. Yueksel, P.-Y. Chen, S. Liu, and K. Varshney, “Is there a trade-off between
 343 fairness and accuracy? a perspective using mismatched hypothesis testing,” in *International*
 344 *Conference on Machine Learning*, pp. 2803–2813, PMLR, 2020.
- 345 [8] H. Zhao, J. Chi, Y. Tian, and G. J. Gordon, “Trade-offs and guarantees of adversarial representa-
 346 tion learning for information obfuscation,” *arXiv preprint arXiv:1906.07902*, 2019.
- 347 [9] N. Martinez, M. Bertran, and G. Sapiro, “Minimax pareto fairness: A multi objective perspec-
 348 tive,” in *International Conference on Machine Learning*, pp. 6755–6764, PMLR, 2020.
- 349 [10] H. Zhao, C. Dan, B. Aragam, T. S. Jaakkola, G. J. Gordon, and P. Ravikumar, “Fundamental
 350 limits and tradeoffs in invariant representation learning,” *arXiv preprint arXiv:2012.10713*,
 351 2020.
- 352 [11] Y. Ganin and V. Lempitsky, “Unsupervised domain adaptation by backpropagation,” in *Internat-
 353 ional Conference on Machine Learning*, 2015.
- 354 [12] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, “Adversarial discriminative domain adaptation,”
 355 in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- 356 [13] H. Zhao, S. Zhang, G. Wu, J. M. Moura, J. P. Costeira, and G. J. Gordon, “Adversarial multiple
 357 source domain adaptation,” *Advances in Neural Information Processing Systems*, vol. 31,
 358 pp. 8559–8570, 2018.
- 359 [14] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel, “Fairness through awareness,” in
 360 *Innovations in Theoretical Computer Science Conference*, pp. 214–226, 2012.
- 361 [15] S. Ruggieri, “Using t-closeness anonymity to control for non-discrimination.,” *Trans. Data*
 362 *Priv.*, vol. 7, no. 2, pp. 99–129, 2014.
- 363 [16] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, “Certifying
 364 and removing disparate impact,” in *ACM SIGKDD International Conference on Knowledge*
 365 *Discovery and Data Mining*, pp. 259–268, 2015.
- 366 [17] F. Calmon, D. Wei, B. Vinzamuri, K. N. Ramamurthy, and K. R. Varshney, “Optimized pre-
 367 processing for discrimination prevention,” in *Advances in Neural Information Processing*
 368 *Systems*, pp. 3992–4001, 2017.
- 369 [18] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, “Learning fair representations,” in
 370 *International Conference on Machine Learning*, pp. 325–333, 2013.
- 371 [19] H. Edwards and A. Storkey, “Censoring representations with an adversary,” *arXiv preprint*
 372 *arXiv:1511.05897*, 2015.
- 373 [20] A. Beutel, J. Chen, Z. Zhao, and E. H. Chi, “Data decisions and theoretical implications when
 374 adversarially learning fair representations,” *arXiv preprint arXiv:1707.00075*, 2017.
- 375 [21] Q. Xie, Z. Dai, Y. Du, E. Hovy, and G. Neubig, “Controllable invariance through adversarial
 376 feature learning,” in *Advances in Neural Information Processing Systems*, pp. 585–596, 2017.

- 377 [22] B. H. Zhang, B. Lemoine, and M. Mitchell, “Mitigating unwanted biases with adversarial
378 learning,” in *AAAI/ACM Conference on AI, Ethics, and Society*, 2018.
- 379 [23] J. Song, P. Kalluri, A. Grover, S. Zhao, and S. Ermon, “Learning controllable fair representa-
380 tions,” in *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- 381 [24] M. Bertran, N. Martinez, A. Papadaki, Q. Qiu, M. Rodrigues, G. Reeves, and G. Sapiro, “Ad-
382 versarially learned representations for information obfuscation and inference,” in *International
383 Conference on Machine Learning*, 2019.
- 384 [25] E. Creager, D. Madras, J.-H. Jacobsen, M. A. Weis, K. Swersky, T. Pitassi, and R. Zemel,
385 “Flexibly fair representation learning by disentanglement,” in *International Conference on
386 Machine Learning (ICML)*, 2019.
- 387 [26] F. Locatello, G. Abbati, T. Rainforth, S. Bauer, B. Schölkopf, and O. Bachem, “On the fair-
388 ness of disentangled representations,” in *Advances in Neural Information Processing Systems*,
389 pp. 14611–14624, 2019.
- 390 [27] J. Mary, C. Calauzenes, and N. El Karoui, “Fairness-aware learning for continuous attributes
391 and treatments,” in *International Conference on Machine Learning*, pp. 4382–4391, PMLR,
392 2019.
- 393 [28] J. Hamm, “Minimax filter: Learning to preserve privacy from inference attacks,” *The Journal of
394 Machine Learning Research*, vol. 18, no. 1, pp. 4704–4734, 2017.
- 395 [29] M. Coavoux, S. Narayan, and S. B. Cohen, “Privacy-preserving neural representations of text,”
396 *arXiv preprint arXiv:1808.09408*, 2018.
- 397 [30] T. Xiao, Y.-H. Tsai, K. Sohn, M. Chandraker, and M.-H. Yang, “Adversarial learning of
398 privacy-preserving and task-oriented representations,” in *Proceedings of the AAAI Conference
399 on Artificial Intelligence*, vol. 34, pp. 12434–12441, 2020.
- 400 [31] M. Dusmanu, J. L. Schönberger, S. N. Sinha, and M. Pollefeys, “Privacy-preserving visual
401 feature descriptors through adversarial affine subspace embedding,” in *IEEE Conference on
402 Computer Vision and Pattern Recognition*, 2021.
- 403 [32] E. Adeli, Q. Zhao, A. Pfefferbaum, E. V. Sullivan, L. Fei-Fei, J. C. Niebles, and K. M. Pohl,
404 “Representation learning with statistical independence to mitigate bias,” in *IEEE/CVF Winter
405 Conference on Applications of Computer Vision*, pp. 2513–2523, 2021.
- 406 [33] V. Grari, O. E. Hajouji, S. Lamprier, and M. Detyniecki, “Learning unbiased representations via
407 re’nyi minimization,” *arXiv preprint arXiv:2009.03183*, 2020.
- 408 [34] B. K. Sriperumbudur, K. Fukumizu, and G. R. Lanckriet, “Universality, characteristic kernels
409 and rkhs embedding of measures.,” *Journal of Machine Learning Research*, vol. 12, no. 7, 2011.
- 410 [35] D. Greenfeld and U. Shalit, “Robust learning with the hilbert-schmidt independence criterion,”
411 in *International Conference on Machine Learning*, pp. 3759–3768, PMLR, 2020.
- 412 [36] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm,
413 “Mutual information neural estimation,” in *International Conference on Machine Learning*,
414 pp. 531–540, PMLR, 2018.
- 415 [37] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf, “Measuring statistical dependence
416 with hilbert-schmidt norms,” in *International Conference on Algorithmic Learning Theory*,
417 pp. 63–77, Springer, 2005.
- 418 [38] A. Gretton, R. Herbrich, A. Smola, O. Bousquet, and B. Schölkopf, “Kernel methods for
419 measuring independence,” *Journal of Machine Learning Research*, vol. 6, no. Dec, pp. 2075–
420 2129, 2005.
- 421 [39] F. R. Bach and M. I. Jordan, “Kernel independent component analysis,” *Journal of Machine
422 Learning Research*, vol. 3, no. Jul, pp. 1–48, 2002.
- 423 [40] K. Fukumizu, F. R. Bach, and A. Gretton, “Statistical consistency of kernel canonical correlation
424 analysis.,” *Journal of Machine Learning Research*, vol. 8, no. 2, 2007.
- 425 [41] E. Kreyszig, *Introductory functional analysis with applications*, vol. 1. wiley New York, 1978.
- 426 [42] P. Diaconis and D. Freedman, “Asymptotics of graphical projection pursuit,” *The Annals of
427 Statistics*, pp. 793–815, 1984.

- 428 [43] P. Hall and K.-C. Li, “On almost linearity of low dimensional projections from high dimensional
429 data,” *The Annals of Statistics*, pp. 867–889, 1993.
- 430 [44] A. Rényi, “On measures of dependence,” *Acta Mathematica Academiae Scientiarum Hungarica*,
431 vol. 10, no. 3-4, pp. 441–451, 1959.
- 432 [45] J. Shawe-Taylor, N. Cristianini, *et al.*, *Kernel methods for pattern analysis*. Cambridge university
433 press, 2004.
- 434 [46] Y. Bengio, A. Courville, and P. Vincent, “Representation learning: A review and new per-
435 spectives,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8,
436 pp. 1798–1828, 2013.
- 437 [47] J. M. Mooij, J. Peters, D. Janzing, J. Zscheischler, and B. Schölkopf, “Distinguishing cause from
438 effect using observational data: methods and benchmarks,” *The Journal of Machine Learning
439 Research*, vol. 17, no. 1, pp. 1103–1204, 2016.

440 Checklist

- 441 1. For all authors...
- 442 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
443 contributions and scope? [Yes] See Section 3.2 and Section 3.3. Particularly, see
444 Theorem 4 and Theorem 5.
- 445 (b) Did you describe the limitations of your work? [Yes] See the discussion above equation
446 (5) and below equation (8).
- 447 (c) Did you discuss any potential negative societal impacts of your work? [N/A] See
448 Section 5.
- 449 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
450 them? [Yes]
- 451 2. If you are including theoretical results...
- 452 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 3.1.
453 Particularly, see Assumption 1 and Assumption 2 and discussion above equation (5)
454 and below equation (8).
- 455 (b) Did you include complete proofs of all theoretical results? [Yes] See supplementary
456 material for the proofs of all Lemmas and Theorems.
- 457 3. If you ran experiments...
- 458 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
459 mental results (either in the supplemental material or as a URL)? [Yes] See supplemen-
460 tary material.
- 461 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
462 were chosen)? [Yes] See Section 4 and supplementary material.
- 463 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
464 ments multiple times)? [Yes] Only one of the baseline methods, ARL, requires running
465 multiple times with different random seeds. The error bar of ARL results is given in
466 supplementary material.
- 467 (d) Did you include the total amount of compute and the type of resources used (e.g.,
468 type of GPUs, internal cluster, or cloud provider)? [No] This paper is not about
469 computational complexity and/or execution time.
- 470 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 471 (a) If your work uses existing assets, did you cite the creators? [Yes] See supplementary
472 material for the citation to the publicly available repository that we used.
- 473 (b) Did you mention the license of the assets? [N/A]
- 474 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
475 We are not using any new assets.
- 476 (d) Did you discuss whether and how consent was obtained from people whose data you’re
477 using/curating? [N/A]

- 478 (e) Did you discuss whether the data you are using/curating contains personally identifiable
479 information or offensive content? [N/A]
- 480 5. If you used crowdsourcing or conducted research with human subjects...
- 481 (a) Did you include the full text of instructions given to participants and screenshots, if
482 applicable? [N/A]
- 483 (b) Did you describe any potential participant risks, with links to Institutional Review
484 Board (IRB) approvals, if applicable? [N/A]
- 485 (c) Did you include the estimated hourly wage paid to participants and the total amount
486 spent on participant compensation? [N/A]