# A representation-learning game for classes of prediction tasks

Anonymous Author(s) Affiliation Address email

# Abstract

We introduce a formulation for learning dimensionality-reducing representations 1 2 of unlabeled feature vectors, when a prior knowledge on future prediction tasks 3 is available. The formulation is based on a three-player game, in which the first player chooses a representation, the second player then adversarially chooses a 4 prediction task, and the third player predicts the response based on the represented 5 features. The first and third player aim is to minimize, and the second player to 6 maximize, the *regret*: The minimal prediction loss using the representation com-7 pared to the same loss using the original features. Our first contribution is theoret-8 9 ical and addresses the mean squared error loss function, and the case in which the representation, the response to predict and the predictors are all linear functions. 10 We establish the optimal representation in pure strategies, which shows the effec-11 tiveness of the prior knowledge, and the optimal regret in mixed strategies, which 12 shows the usefulness of randomizing the representation. We prove that optimal 13 randomization requires a precisely characterized finite number of representations, 14 which is smaller than the dimension of the feature vector, and potentially much 15 smaller. Our second contribution is an efficient gradient-based iterative algorithm 16 that approximates the optimal mixed representation for a general loss function, 17 and general classes of representations, response functions and predictors. 18

# 19 1 Introduction

A common practice in modern data-science is to collect as much data as possible, even without 20 an exact knowledge of a subsequent prediction task it will be used for. The data collected is an 21 unlabeled set of feature vectors  $\{x_i\} \subset \mathbb{R}^d$  . Then, when a specific prediction task becomes of 22 interest, responses  $y_i \in \mathcal{Y}$  are collected, and a learning algorithm is trained on the pairs  $\{(x_i, y_i)\}$ . 23 Modern sources, such as high-definition images or genomic sequences, have high dimension d, and 24 this raises the question of *dimensionality-reduction*, either for a better generalization [1], for stor-25 age/communication savings [2–4], or for interpretability [5]. The goal is thus to find a representation 26  $z = R(x) \in \mathbb{R}^r$ , where  $d \gg r$ , that preserves the relevant part of the features, without a full knowl-27 edge of their utility for future prediction tasks. In this paper, we propose an unsupervised-learning 28 game-theoretic framework for this goal, whose central aspect is an assumption of prior knowledge 29 on the *class* of future prediction tasks. Our contributions are a theoretical solution in a fully linear 30 setting, under the mean squared error (MSE) loss, and an algorithm for the general setting. 31

Popular approaches to dimensionality reduction are oblivious to prior knowledge on the prediction
task. Most prominently, *principal component analysis* (PCA) [6–9], and non-linear extensions such
as kernel PCA [10] and *auto-encoders* (AE) [11–13, 1], aim that the representation *z* will maximally

preserve the *variation* in x. Nonetheless, prior knowledge may indicate that the highly varying directions in the feature space are irrelevant for future prediction tasks. From the supervised learn-

ing perspective, it is well established that efficient representations are inherent to efficient learning 37 [14, 15]. In this respect, the *information bottleneck* (IB) principle [16–19] was used to postulate 38 that efficient supervised learning learns representations which are both low-complexity and relevant 39 [20–24] (this spurred a debate, e.g., [25, 26]). The original IB formulation is based on the mutual in-40 formation functional [27], which is difficult to estimate (especially in high dimensions), and ignores 41 complexity constraints on the representation or prediction [28, 29]; see a review in Appendix B. Us-42 ing the notion of *usable information*, introduced in [28], optimal representations for the *supervised* 43 learning setting were explored in [29] via a two-player game between Alice, which selects a pre-44 diction problem of y given x, and Bob, which then selects the representation z. Alice then uses an 45 empirical risk minimizer with the standard goal of minimizing the expected risk. It was established 46 in [29] that ideal generalization is obtained for representations that optimize the *decodable IB*. 47

In this paper, we build upon [29], and propose a *three-player game* for unsupervised representation-48 learning, chosen without a specific prediction problem (Section 2). First, the representation player 49 reduce  $x \in \mathbb{R}^d$  to a representation  $z \in \mathbb{R}^r$ , where r < d. Second, the response function player 50 chooses a response (label) y rule f for x, from a given known class of (random) response functions 51  $\mathcal{F}$ . The choice of this class manifests the prior knowledge available on the type of prediction prob-52 lems that the representations will be used for. Third, the *predictor player* optimally predicts y from 53 z. The value of the game is determined by the regret: The prediction loss based on the representation 54 z compared to prediction loss based on x. The first and last player cooperate in order to minimize 55 the regret, whereas the response function player aims to maximize it. In other words, the represen-56 tation is chosen to minimize the worst-case prediction loss for any response in  $\mathcal{F}$ . The output of this 57 game is the representation chosen by the first player. In order to focus on the representation aspect 58 we side-step the generalization problem, and assume that sufficient labeled data will be provided to 59 the predictor later on in order to accurately estimate the prediction rule. 60

This formulation directly addresses the relevance of a "direction" in the feature space to the pre-61 diction tasks in F, rather than its variability, as in standard unsupervised learning (e.g., PCA and 62 AE). Compared to [29], the representation is chosen based only on the class of possible response 63 functions, rather than a specific one. Such knowledge on  $\mathcal{F}$  may stem from various considerations: 64 Domain specific, imposed by privacy or fairness constraints, or stem from transfer or continual 65 learning setting; see Appendix A for an extended discussion. Technically, the game in [29] replaces 66 the order of the first (representation) and second (response) players. From a different perspective, 67 our method is a *self-supervised learning* method, for which the prior knowledge on  $\mathcal{F}$  serves as a 68 "self-defined signal" for choosing an optimal representation, without any labeled data (see, e.g., [30] 69 and [31] for recent surveys). In addition, our game formulation naturally leads to a mixed strategies 70 solution [32], that is, allowing the representation player to randomized the representation rule, in 71 order to mix up the adversarial response player. This randomization is an inherent aspect of the IB 72 73 formulation, but its usage there is not rigorously justified. By contrast, for standard unsupervised learning, mixed representation does not improve the regret (see Proposition 14 in Appendix E.1 for 74 the PCA setting). In Appendix B we provide a thorough discussion of related work. 75

#### 76 Contributions

• Theoretical: We address the fundamental setting in which the representation, the response, and the prediction are linear functions, under the MSE loss function (Section 3). The prior knowledge on  $\mathcal{F}$  is represented by a symmetric matrix S that determines the principal directions of the function in the feature space. We establish the optimal representation and regret in pure strategies, which shows the utility of the prior information, and in mixed strategies, which shows that randomizing the representation yields *strictly lower* regret. We prove that randomizing between merely  $\ell^*$  different representation rules suffices, where  $r + 1 \leq \ell^* \leq d$  is a precisely characterized *effective dimension*.

Algorithmic: We develop an iterative gradients-based algorithm that approximates the optimal mixed representation (Section 4) for general representations/response/predictors and loss functions. The algorithm is greedy, and alternates between finding a new representation rule and an adversarial function. We empirically verify that the output mixed representation has close-to-optimal regret in the linear MSE setting. To optimize the weights of the representation, we essentially solve a minimax two-player games, and to this end, we utilize the classic multiplicative weights update (MWU) algorithm [33] (which is essentially a follow-the-regularized-leader [34, 35]).

# 91 **2** Problem formulation

We use mostly conventional notation that is detailed in Appendix C. Specifically, the eigenvalues of a positive-semidefinite matrix S are denoted as  $\lambda_{\max}(S) \equiv \lambda_1(S) \geq \cdots \geq \lambda_d(S) = \lambda_{\min}(S)$  and  $v_i(S)$  denotes an eigenvector corresponding to  $\lambda_i(S)$  such that  $V = V(S) := [v_1(S), v_2(S), \cdots, v_d(S)] \in \mathbb{R}^{d \times d}$  and  $S = V(S)\Lambda(S)V^{\top}(S)$  is an eigenvalue decomposition. For a matrix  $W \in \mathbb{R}^{d \times d}$  we let  $W_{i:j} := [w_i, \dots, w_j] \in \mathbb{R}^{(j-i+1) \times d}$  denote the matrix comprised of the columns indexed by  $\{i, \dots, j\}$ . We denote the probability law of a random variable  $\boldsymbol{x}$  as  $L(\boldsymbol{x})$ .

Let  $x \in \mathcal{X}$  be a random feature vector, where  $P_x := \mathsf{L}(x)$  is known. Let  $y \in \mathcal{Y}$  be a corresponding response drawn according to a probability kernel  $y \sim f(\cdot | x = x)$ , where for brevity, we will refer to f as the *response function*. We assume  $f \in \mathcal{F}$  for some known class  $\mathcal{F}$ . Let  $z := R(x) \in \mathbb{R}^r$  be an r-dimensional representation of x where  $R: \mathcal{X} \to \mathbb{R}^r$  is chosen from a class  $\mathcal{R}$  of representation functions, and let  $Q: \mathbb{R}^r \to \mathcal{Y}$  be a prediction rule from a class  $\mathcal{Q}$ , with the loss function loss:  $\mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$ . The *regret* of the representation R for the response function f is

$$\operatorname{regret}(R, f \mid P_{\boldsymbol{x}}) := \min_{Q \in \mathcal{Q}} \mathbb{E}\left[\operatorname{loss}(\boldsymbol{y}, Q(R(\boldsymbol{x})))\right] - \min_{Q: \mathbb{R}^d \to \mathcal{Y}} \mathbb{E}\left[\operatorname{loss}(\boldsymbol{y}, Q(\boldsymbol{x}))\right].$$
(1)

The minimax regret in mixed strategies is defined via the worst case response function in  $\mathcal{F}$  as

$$\operatorname{regret}_{\min}(\mathcal{R}, \mathcal{F} \mid P_{\boldsymbol{x}}) := \min_{\mathsf{L}(\boldsymbol{R}) \in \mathcal{P}(\mathcal{R})} \max_{f \in \mathcal{F}} \mathbb{E}\left[\operatorname{regret}(\boldsymbol{R}, f \mid P_{\boldsymbol{x}})\right],$$
(2)

where  $\mathcal{P}(\mathcal{R})$  is a set of probability measures on the possible set of representations  $\mathcal{R}$ . The *minimax regret in pure strategies* restricts  $\mathcal{P}(\mathcal{R})$  to degenerated measures (deterministic), and so the expectation in (2) is removed. Our main goal is to determine the optimal representation strategy, either in pure  $\mathbb{R}^* \in \mathcal{R}$  or mixed strategies  $L(\mathbb{R}^*) \in \mathcal{P}(\mathcal{R})$ . To this end, we will also utilize the *maximin* version of (2). Specifically, let  $\mathcal{P}(\mathcal{F})$  denote a set of probability measures supported on  $\mathcal{F}$ , and assume that for any  $\mathbb{R} \in \mathcal{R}$ , there exists a measure in  $\mathcal{P}(\mathcal{R})$  that puts all its mass on  $\mathbb{R}$ . Then, the *minimax theorem* [32, Chapter 2.4] [36] implies that

$$\operatorname{regret}_{\mathsf{mix}}(\mathcal{R}, \mathcal{F} \mid P_{\boldsymbol{x}}) = \max_{\mathsf{L}(\boldsymbol{f}) \in \mathcal{P}(\mathcal{F})} \min_{R \in \mathcal{R}} \mathbb{E}\left[\operatorname{regret}(R, \boldsymbol{f} \mid P_{\boldsymbol{x}})\right].$$
(3)

The right-hand side of (3) is the maximin regret in mixed strategies, and the maximizing prob-112 ability law  $L(f^*)$  is known as the *least favorable prior*. In general,  $\operatorname{regret}_{mix}(\mathcal{R}, \mathcal{F} \mid P_x) \leq 1$ 113  $\operatorname{regret}_{\operatorname{pure}}(\mathcal{R}, \mathcal{F} \mid P_{x})$ , and the inequality can be strict. We mention that the use of expectation 114 in the definition of the mixed regret over the randomized representation, implies that the empirical 115 performance of a system based on this randomized representation achieves the mixed minimax re-116 gret value in the limit of large number of repeating representation games. The size of the dataset for 117 each of these games should be large enough to allow for accurate learning of f to be used by the 118 predictor. By contrast, the pure minimax regret guarantee is valid for a single representation, and 119 thus more conservative from this aspect. 120

#### 121 **3** The linear setting under MSE loss

In this section, we focus on linear classes and the MSE loss function. The response function class is characterized by a quadratic constraint, to wit, the class  $\mathcal{F}$  is specified by a matrix  $S \in \mathbb{S}_{++}^d$  that represents the relative importance of each direction in the feature space in determining y.

**Definition 1** (The linear MSE setting). Assume that  $\mathcal{X} = \mathbb{R}^d$ , that  $\mathcal{Y} = \mathbb{R}$  and the loss function is the MSE,  $|oss(y_1, y_2) = |y_1 - y_2|^2$ . Assume that  $\mathbb{E}[\boldsymbol{x}] = 0$  and let  $\Sigma_{\boldsymbol{x}} := \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^T] \in \mathbb{S}_{++}^d$  be its invertible covariance matrix. The classes of representations, response functions, and predictors are all linear, that is: (1) The representation is  $z = R(x) = R^\top x$  for  $R \in \mathcal{R} := \mathbb{R}^{d \times r}$  where d > r; (2) The response function is  $F \in \mathcal{F} \subset \mathbb{R}^d$ , and  $\boldsymbol{y} = f^\top \boldsymbol{x} + \boldsymbol{n} \in \mathbb{R}$ , where  $\boldsymbol{n} \in \mathbb{R}$  is a heteroscedastic noise that satisfies  $\mathbb{E}[\boldsymbol{n} \mid \boldsymbol{x}] = 0$ , and given some specified  $S \in \mathbb{S}_{++}^d$ 

$$f \in \mathcal{F}_S := \left\{ f \in \mathbb{R}^d \colon \|f\|_S^2 \le 1 \right\},\tag{4}$$

where  $||f||_{S} = ||S^{-1/2}f||_2 = (f^{\top}S^{-1}f)^{1/2}$  is the Mahalanobis norm; (3) The predictor is  $Q(z) = q^{\top}z \in \mathbb{R}$  for  $q \in \mathbb{R}^r$ . Since the regret will depend on  $P_x$  only via  $\Sigma_x$ , we will abbreviate the notation of the pure (resp. mixed) minimax regret to regret  $_{\text{pure}}(\mathcal{F} | \Sigma_x)$  (resp. regret  $_{\text{mix}}(\mathcal{F} | \Sigma_x)$ ).

In Appendix E.1 we show that standard PCA can be similarly formulated, by assuming that  $\mathcal{F}$ is a singleton containing the noiseless identity function, so that y = x surely holds, and  $\hat{x} = Q(z) \in \mathbb{R}^d$ . Proposition 14 therein shows that the pure and mixed minimax representations are both  $R = V_{1:r}(\Sigma_x)$ , and so randomization is not unnecessary. We begin with the pure minimax regret.

**Theorem 2.** For the linear MSE setting (Definition 1)

$$\operatorname{regret}_{\operatorname{pure}}(\mathcal{F}_{S} \mid \Sigma_{\boldsymbol{x}}) = \lambda_{r+1} \left( \Sigma_{\boldsymbol{x}}^{1/2} S \Sigma_{\boldsymbol{x}}^{1/2} \right).$$
(5)

139 A minimax representation matrix is

$$R^* := \Sigma_{x}^{-1/2} \cdot V_{1:r} \left( \Sigma_{x}^{1/2} S \Sigma_{x}^{1/2} \right), \tag{6}$$

140 and the worst case response function is

$$f^* := S^{1/2} \cdot v_{r+1} \left( \Sigma_{\boldsymbol{x}}^{1/2} S \Sigma_{\boldsymbol{x}}^{1/2} \right).$$
(7)

The optimal representation thus whitens the feature vector x, and then projects it on the top reigenvectors of the adjusted covariance matrix  $\Sigma_x^{1/2} S \Sigma_x^{1/2}$ , which reflects the prior knowledge that  $f \in \mathcal{F}_S$ . The proof is deferred to Appendix E.2, and its outline is as follows: Plugging the optimal predictor into the regret results a quadratic form in  $f \in \mathbb{R}^d$ , determined by a matrix which depends on the subspace spanned by the representation R. The worst-case f is the determined via the *Rayleigh quotient theorem* [37, Theorem 4.2.2], and the optimal R is found via the *Courant– Fischer variational characterization* [37, Theorem 4.2.6] (see Appendix D for a summary of useful mathematical results). We next consider the mixed minimax regret.

149 **Theorem 3.** For the linear MSE setting (Definition 1)

$$\operatorname{regret}_{\mathsf{mix}}(\mathcal{F}_S \mid \Sigma_{\boldsymbol{x}}) = \frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \lambda_i^{-1}},\tag{8}$$

150 where  $\lambda_i \equiv \lambda_i (S^{1/2} \Sigma_x S^{1/2})$  and  $\ell^*$  is any member of

$$\left\{\ell \in [d] \setminus [r]: (\ell - r) \cdot \lambda_{\ell}^{-1} \le \sum_{i=1}^{\ell} \lambda_i^{-1} \le (\ell - r) \cdot \lambda_{\ell+1}^{-1}\right\}$$
(9)

- 151 (with  $\lambda_{d+1} \equiv 0$ ). Furthermore:
- The covariance matrix of the least favorable prior of  $\boldsymbol{f}$ : Let  $\Lambda_{\ell} := \operatorname{diag}(\lambda_1, \dots, \lambda_{\ell^*}, 0, \dots, 0)$ , and let  $V \equiv V(S^{1/2}\Sigma_{\boldsymbol{x}}S^{1/2})$ . Then, the covariance matrix of the least favorable prior of  $\boldsymbol{f}$  is

$$\Sigma_{f}^{*} := \frac{V^{\top} \Lambda_{\ell^{*}}^{-1} V}{\sum_{i=1}^{\ell^{*}} \lambda_{i}^{-1}}.$$
(10)

• The probability law of the minimax representation: Let  $\overline{A} \in \{0,1\}^{\ell^* \times \binom{\ell^*}{r}}$  be a matrix whose columns are the members of the set

$$\overline{\mathcal{A}} := \{ \overline{a} \in \{0, 1\}^{\ell^*} : \|\overline{a}\|_1 = \ell^* - r \}$$
(11)

156 (in an arbitrary order). Let  $\overline{b} = (b_1, \dots, b_{\ell^*})^\top$  be such that

$$b_i = (\ell^* - r) \cdot \frac{\lambda_i^{-1}}{\sum_{j=1}^{\ell^*} \lambda_j^{-1}}.$$
(12)

Then, there exists a solution  $p \in [0,1]^{\binom{\ell^*}{r}}$  with support size at most  $\ell^* + 1$  to  $\overline{A}p = \overline{b}$ . For  $j \in [\binom{\ell^*}{r}]$ , let  $\mathcal{I}_j := \{i \in [\ell^*]: \overline{A}_{ij} = 0\}$  be the zero indices on the *j*th column of  $\overline{A}$ , and let  $V_{\mathcal{I}_j}$  denote the *r* columns of *V* whose index is in  $\mathcal{I}_j$ . A minimax representation is

$$\boldsymbol{R}^* = \Sigma_{\boldsymbol{x}}^{-1/2} V_{\mathcal{I}_j} \tag{13}$$

160 with probability  $p_j$ , for  $j \in [\binom{\ell^*}{r}]$ .

Interestingly, while the eigenvalues  $\lambda_i(\Sigma_x^{1/2}S\Sigma_x^{1/2}) = \lambda_i(S^{1/2}\Sigma_xS^{1/2})$  are equal, the pure minimax regret utilizes the eigenvectors of  $\Sigma_x^{1/2}S\Sigma_x^{1/2}$  whereas the mixed minimax regret utilizes those 161 162 of  $S^{1/2}\Sigma_x S^{1/2}$ , which are possibly different. The proof of Theorem 3 is also in Appendix E.2, and 163 is substantially more complicated and longer than for the pure regret. We use a two-step indirect 164 approach, since it seems challenging to directly maximize over  $L(\mathbf{R})$ . First, we solve the maximin 165 problem (3), and find the least favorable prior  $L(f^*)$ . Second, we propose a probability law for the 166 representation  $L(\mathbf{R})$ , and show that its regret equals the maximin value, and thus also the minimax. 167 With more detail, in the first step, we show that the regret only depends on L(f) via  $\Sigma_f = \mathbb{E}[ff^+]$ , 168 and we explicitly construct a probability law that is both fully supported on  $\mathcal{F}_S$  and has this co-169 variance matrix. This reduces the problem from optimizing L(f) to optimizing  $\Sigma_f$ , whose solution 170 (Lemma 16) leads to the least favorable  $\Sigma_f^*$ , and then to the maximin value. In the second step, 171 we explicitly construct a representation that achieves the maximin regret. Concretely, we construct 172 representation matrices that use r of the  $\ell^*$  principal components of  $\Sigma_x^{1/2} S \Sigma_x^{1/2}$ , where  $\ell^* > r$ . 173 The defining property of  $\ell^*$  (9) established in the maximin solution is utilized to find weights on the 174  $\binom{\ell^*}{r}$  possible representations, that achieves the maximin solution, and thus also the minimax. The 175 proof uses Carathéodory's theorem (see Appendix D) which also establishes that the optimal  $\{p_i\}$ 176 is supported on at most  $\ell^* + 1$  matrices, much less than  $\binom{\ell^*}{r}$ . We next make a few comments: 177

1. Computing the mixed minimax probability: This requires solving  $\overline{A}^{\dagger} p = \overline{b}$  for a probability 178 vector p, which is a linear-program feasibility problem that is routinely solved [38]. For illustration, 179 if r = 1 then  $\overline{A} \in \{0,1\}^{\ell^* \times \ell^*}$  is a square all ones matrix, except for a zero diagonal, and  $p_j =$ 180  $1 - (\ell^* - 1)\lambda_j^{-1}/(\sum_{i=1}^{\ell^*} \lambda_i^{-1})$  for  $j \in [\ell^*]$ . Similarly, the case  $\ell^* = r + 1$  is solved by setting 181  $p_j = (\lambda_i^{-1})/(\sum_{i'=1}^{\ell^*} \lambda_{i'}^{-1})$  on the  $\ell^*$  standard basis vectors. Nonetheless, the dimension of p is  $\binom{\ell^*}{r}$ 182 and thus increases fast as  $\Theta((\ell^*)^r)$ , and this approach may be intractable. However, in this case the 183 algorithm we present in Section 4 can be used. As we empirically show, it approximately achieves 184 the optimal regret, and the number of atoms is not much larger than  $\ell^* + 1$ . 185

186 2. <u>Required randomness</u>: The regret formulation (2) assumes that the actual realization of the rep-187 resentation rule is known to the predictor. Formally, this can be conveyed to the predictor using an 188 small header of less than  $\log_2(\ell^* + 1) \leq \log(d + 1)$  bits. Practically, this is unnecessary and an 189 efficient predictor can be learned from a labeled data set (z, y).

<sup>190</sup> 3. The rank of  $\sum_{f=1}^{k}$  The rank of the covariance matrix of the least favorable prior is an *effective* <sup>191</sup> *dimension*, satisfying (see (8))

$$\ell^* = \underset{\ell \in [d] \setminus [r]}{\operatorname{arg\,max}} \frac{1 - (r/\ell)}{\frac{1}{\ell} \sum_{i=1}^{\ell} \lambda_i^{-1}}.$$
(14)

By convention,  $\{\lambda_i^{-1}\}_{i \in [d]}$  is a monotonic non-decreasing sequence, and so is the partial Cesàro 192 mean  $\psi(\ell) := \frac{1}{\ell} \sum_{i=1}^{\ell} \lambda_i^{-1}$ . For example, if  $\lambda_i = i^{-\alpha}$  with  $\alpha > 0$  then  $\psi(\ell) = \Theta(\ell^{\alpha})$ . If, e.g.,  $\psi(\ell) = \ell^{\alpha}$ , then it is easily derived that  $\ell^* \approx \min\{\frac{\alpha+1}{\alpha}r, d\}$ . So, if  $\alpha \ge \frac{r}{d-r}$  is large enough and the decay rate of  $\{\lambda_i\}$  is fast enough then  $\ell^* < d$ , and otherwise  $\ell^* = d$ . As the decay rate of  $\{\lambda_i\}$  becomes faster, the rank of  $\Sigma_f^*$  decreases to r. Importantly,  $\ell^* \ge r+1$  always holds, and so 193 194 195 196 the optimal mixed representation is not deterministic even if  $S^{1/2}\Sigma_{x}S^{1/2}$  has less than r significant 197 eigenvalues (which can be represented by a single matrix  $R \in \mathbb{R}^{d \times r}$ ). Hence, the mixed minimax 198 regret is always strictly lower than the pure minimax regret. Thus, even when  $S = I_d$ , and no 199 valuable prior knowledge is known on the response function, the mixed minimax representation is 200 different from the standard PCA solution of top r eigenvectors of  $\Sigma_x$ . 201

4. <u>Uniqueness of the optimal representation</u>: Since one can always post-multiply  $R^{\top}x$  by some invertible matrix, and then pre-multiply  $z = R^{\top}x$  by its inverse, the following simple observation holds: When  $\mathcal{R}$  and  $\mathcal{Q}$  are not further restricted, then if  $\mathbf{R}$  is a minimax representation, and  $W(\mathbf{R}) \in \mathbb{R}^{r \times r}$  is an invertible matrix, then  $\mathbf{R} \cdot W(\mathbf{R})$  is also a minimax representation.

<sup>206</sup> 5. Infinite-dimensional features: Theorems 2 and 3 assume a finite dimensional feature space, but <sup>207</sup> as we show in Appendix F, the results can be easily generalized to an infinite dimensional Hilbert <sup>208</sup> space  $\mathcal{X}$ , in the more restrictive setting that the noise n is statistically independent of x.

**Example 4.** Assume  $S = I_d$ , and denote, for brevity,  $V \equiv V(\Sigma_x) := [v_1, \dots, v_d]$  and  $\Lambda \equiv \Lambda(\Sigma_x) := \text{diag}(\lambda_1, \dots, \lambda_d)$ . The optimal minimax representation in pure strategies (Theorem 2) is



Figure 1: Left: Pure and mixed minimax regret and  $\ell_*$  for Example 4, for d = 50, r = 25, with  $\lambda_i = \sigma_i^2 \propto i^{-\alpha}$ . Right: Pure and mixed minimax regret and  $\ell_*$  for Example 5, for d = 50, r = 25, with  $\sigma_i^2 \propto i^{-\alpha}$  and  $s_i \propto i^2$ . The trend of  $\ell_*$  is reversed for  $\alpha > 2$ .

211 then

$$R^* = \Sigma_{\boldsymbol{x}}^{-1/2} \cdot V_{1:r} = V \Lambda_{\boldsymbol{x}}^{-1/2} V^{\top} V_{1:r} = V \Lambda_{\boldsymbol{x}}^{-1/2} \cdot [e_1, \dots, e_r] = \left[ \lambda_1^{-1/2} \cdot v_1, \dots, \lambda_r^{-1/2} \cdot v_r \right],$$
(15)

which is comprised of the top r eigenvectors of  $\Sigma_x$ , scaled so that  $v_i^{\top} x$  has unit variance. By Comment 4 above,  $V_{1:r}$  is also an optimal minimax representation. The worst case response is  $f = v_{r+1}(\Sigma_x)$  and, as expected, since R uses the first r principal directions

$$\operatorname{regret}_{\operatorname{pure}}(\mathcal{F} \mid \Sigma_{\boldsymbol{x}}) = \lambda_{r+1}. \tag{16}$$

<sup>215</sup> The minimax regret in mixed strategies (Theorem 3) is different, and given by

$$\operatorname{regret}_{\mathsf{mix}}(\mathcal{F} \mid \Sigma_{\boldsymbol{x}}) = \frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \lambda_i^{-1}},\tag{17}$$

where  $\ell^*$  is determined by the decay rate of the eigenvalues of  $\Sigma_x$  (see (9)). The least favorable covariance matrix is given by (Theorem 3)

$$\Sigma_{\boldsymbol{f}}^* = \left[\sum_{i=1}^{\ell^*} \lambda_i^{-1}\right]^{-1} \cdot V \operatorname{diag}\left(\lambda_1^{-1}, \dots, \lambda_{\ell^*}^{-1}, 0, \dots, 0\right) \cdot V^{\top}.$$
(18)

Intuitively, the least favorable  $\Sigma_{f}^{*}$  equalizes the first  $\ell^{*}$  eigenvalues of  $\Sigma_{x}\Sigma_{f}^{*}$  (and nulls the other d -  $\ell^{*}$ ) so that the representation is indifferent to these  $\ell^{*}$  directions. As evident from the regret, the "equalization" of the *i*th eigenvalue adds a term of  $\lambda_{i}^{-1}$  to the denominator, and if  $\lambda_{i}$  is too small then  $v_{i}$  is not chosen for the representation, as agrees with Comment 3 above (a fast decay of  $\{\lambda_{i}\}$ reduces  $\ell_{*}$  away from d). The mixed minimax representation sets

$$\boldsymbol{R}^{*} = \Sigma_{\boldsymbol{x}}^{-1/2} \cdot V_{\mathcal{I}_{j}} = \left[ \lambda_{i_{j,1}}^{-1/2} \cdot v_{i_{j,1}}, \dots, \lambda_{i_{j,r}}^{-1/2} \cdot v_{i_{j,r}} \right]$$
(19)

with probability  $p_j$ , where  $\mathcal{I}_j \equiv \{i_{j,1}, \ldots, i_{j,r}\}$  (the derivation is similar to (15)). Thus, the optimal representation chooses a random subset of r vectors from  $\{v_1, \ldots, v_{\ell^*}\}$ . See the left panel of Figure 1 for a numerical example.

Example 5. To demonstrate the effect of prior knowledge on the response function, we assume  $\Sigma_{\boldsymbol{x}} = \operatorname{diag}(\sigma_1^2, \ldots, \sigma_d^2)$  and  $S = \operatorname{diag}(s_1, \ldots, s_d)$ , where  $\sigma_1^2 \ge \sigma_2^2 \ge \cdots \ge \sigma_d^2$  (but  $\{s_i\}_{i \in [d]}$ are not necessarily ordered). Letting  $f = (f_1, \ldots, f_d)$ , the class of response functions is  $\mathcal{F}_S :=$   $\{f \in \mathbb{R}^d: \sum_{i=1}^d (f_i^2/s_i) \le 1\}$ , and so coordinates  $i \in [d]$  with a large  $s_i$  have large influence on the response. Let  $(i_{(1)}, \ldots, i_{(d)})$  be a permutation of [d] so that  $\sigma_{i(j)}^2 s_{i(j)}$  it the *j*th largest value of  $(\sigma_i^2 s_i)_{i \in [d]}$ . The pure minimax regret is (Theorem 2)

$$\operatorname{regret}_{\operatorname{pure}}(\mathcal{F} \mid \Sigma_{\boldsymbol{x}}) = \sigma_{i_{r+1}}^2 s_{i_{r+1}}.$$
(20)

The optimal representation is  $R = [e_{i_{(1)}}, e_{i_{(2)}}, \dots, e_{i_{(r)}}]$ , that is, uses the most influential coordinates, according to  $\{s_i\}$ , which may be different from the r principal directions of  $\Sigma_x$ . For the minimax regret in mixed strategies, Theorem 3 results

$$\operatorname{regret}_{\operatorname{mix}}(\mathcal{F} \mid \Sigma_{\boldsymbol{x}}) = \frac{\ell^* - r}{\sum_{j=1}^{\ell^*} (s_{i_j} \sigma_{i_j}^2)^{-1}}$$
(21)

for  $\ell^* \in [d] \setminus [r]$  satisfying (9), and the covariance matrix of the least favorable prior is given by

$$\Sigma_{\boldsymbol{f}}^{*} = \frac{\sum_{j=1}^{\ell^{*}} \sigma_{i_{j}}^{-2} \cdot e_{i_{j}} e_{i_{j}}^{\top}}{\sum_{j=1}^{\ell^{*}} (s_{i_{j}} \sigma_{i_{j}}^{2})^{-1}}.$$
(22)

That is, up to a scale factor  $(\sum_{i=1}^{\ell^*} s_i^{-1} \sigma_i^{-2})^{-1}$ , the matrix is diagonal so that the *k*th term on the diagonal is  $\sum_{f}^{*}(k,k) = \sigma_{k}^{-2}$  if  $k = i_j$  for some  $j \in [\ell^*]$  and  $\sum_{f}^{*}(k,k) = 0$  otherwise. As in Example 4,  $\sum_{f}^{*}$  equalizes the first  $\ell^*$  eigenvalues of  $\sum_{x} \sum_{f}$  (and nulls the other  $d - \ell^*$ ). However, it does so in a manner that chooses them according to their influence on  $f^{\top}x$ . The random minimax representation in mixed strategies is

$$\boldsymbol{R}^{*} = \left[\sigma_{i_{j,1}}^{-1} \cdot e_{i_{j,1}}, \dots, \sigma_{i_{j,r}}^{-1} \cdot e_{i_{j,r}}\right]$$
(23)

with probability  $p_j$ . Again, all the first  $\ell^*$  coordinates are used, and not just the top r. See the right panel of Figure 1 for a numerical example. We finally remark that, naturally, in the non-diagonal case, the minimax regret will also depend on the relative alignment between S and  $\Sigma_x$ .

# <sup>244</sup> 4 An iterative algorithm for general classes and loss functions

In this section, we develop an iterative algorithm for finding the optimal representation in mixed 245 strategies, i.e., solving (2) for general classes and loss functions. Since optimizing general probabil-246 ity measures over  $\mathcal{R}$  is formidable, we restrict the optimization to finite mixed representations, i.e., 247 assume that  $\mathbf{R} = R^{(j)} \in \mathcal{R}$  with probability  $p^{(j)}$ , where  $j \in [m]$  (which suffices for the linear MSE 248 setting of Section 3, but possibly sub-optimal in general). Furthermore, the algorithm's operation 249 will require randomization also for the response player, and so we set  $f = f^{(i)} \in \mathcal{F}$  with probability 250  $o^{(i)}$  where  $i \in [\overline{m}]$ , and  $\overline{m} = m_0 + m$  for some  $m_0 \ge 0$ . The resulting optimization problem then 251 becomes 252

$$\min_{\{p^{(j)},R^{(j)}\in\mathcal{R}\}} \max_{\{o^{(i)},f^{(i)}\in\mathcal{F}\}} \min_{\{Q^{(j,i)}\in\mathcal{Q}\}} \sum_{j\in[m]} \sum_{i\in[\overline{m}]} p^{(j)} \cdot o^{(i)} \cdot \mathbb{E}\left[\mathsf{loss}(f^{(i)}(\boldsymbol{x}),Q^{(j,i)}(R^{(j)}(\boldsymbol{x})))\right],$$
(24)

under the constraints  $p^{(j)} \ge 0$  and  $\sum_j p^{(j)} = 1$ , and  $o^{(i)} \ge 0$  and  $\sum_i o^{(i)} = 1$ . Note that the prediction rule  $Q^{(j,i)}$  is determined based on both  $R^{(j)}$  and  $f^{(i)}$ , and that the ultimate goal of solving (24) is just to extract the optimal **R**.

A high level description of the algorithm is to gradually add more representations to the support size 256 of  $\mathbf{R}$  up to m, where next k will denote the current number of representations,  $k \in [m]$ . Initialization 257 requires an representation  $R^{(1)}$ , as well as a *set* of functions  $\{f^{(i)}\}_{i\in m_0}$ , so that the final support 258 size of f will be  $\overline{m} = m_0 + m$ . Finding this initial representation and the set of functions is based 259 on the specific loss function and a possible set of representation/predictors. At iteration  $k \in [m]$ , the 260 main loop of the algorithm has two phases. In the first phase, a new adversarial function is added 261 to the set of functions, as the worse function for the current random representation. In the second 262 263 phase, a new representation atom is added to the set of possible representations. This representation is determined based on the given set of functions. Concretely, the two phases operate as follows: 264

• Phase 1 – Given k representations  $\{R^{(j)}\}_{j \in (k)}$  with weights  $\{p^{(j)}\}_{j \in [k]}$ , the algorithm determines the function  $f^{(m_0+k)}$  as the worst function for this random representation (optimal adversarial action of the response function player). Specifically,

$$\operatorname{reg}_{k} := \operatorname{regret}_{\operatorname{mix}}(\{R^{(j)}, p^{(j)}\}_{j \in [k]}, \mathcal{F} \mid P_{\boldsymbol{x}})$$

$$(25)$$

$$:= \max_{f \in \mathcal{F}} \min_{\{Q^{(j)} \in \mathcal{Q}\}_{j \in [k]}} \sum_{j \in [k]} p^{(j)} \cdot \mathbb{E}\left[ \mathsf{loss}(f(\boldsymbol{x}), Q^{(j)}(R^{(j)}(\boldsymbol{x}))) \right]$$
(26)

is solved, and  $f^{(m_0+k)}$  is set to be the maximizer. This simplifies (24) in the sense that m is replaced by k, the random representation  $\mathbf{R}$  is kept fixed, and  $f \in \mathcal{F}$  is optimized as a pure strategy (the previous functions  $\{f^{(i)}\}_{i \in [m_0+k-1]}$  are ignored).

• Phase 2 – Adding a representation atom: Given fixed  $\{f^{(j)}\}_{j \in [m_0+k]}$  and  $\{R^{(j)}\}_{j \in [k]}$ , a new 271 representation  $R^{(k+1)}$  is found as the most incrementally valuable representation atom. Specifically,

$$\min_{R^{(k+1)} \in \mathcal{R}} \operatorname{regret}_{\operatorname{mix}} \left\{ \{R^{(j_1)}\}_{j_1 \in [k+1]}, \{f^{(j_2)}\}_{j_2 \in [m_0+k]} \mid P_{\boldsymbol{x}} \right\}$$

$$:= \min_{R^{(k+1)} \in \mathcal{R}} \min_{\{p^{(j_1)}\}_{j_1 \in [k+1]}} \max_{\{o^{(j_2)}\}_{j_2 \in [m_0+k]}} \min_{\{Q^{(j_1,j_2)} \in \mathcal{Q}\}_{j_1 \in [k+1], j_2 \in [m_0+k]}} \min_{\{q^{(j_1)} \in \mathcal{Q}\}_{j_1 \in [k+1], j_2 \in [m_0+k]}} \sum_{j_1 \in [k+1]} \sum_{j_2 \in [m_0+k]} p^{(j_1)} \cdot o^{(j_2)} \cdot \mathbb{E} \left[ \operatorname{loss}(f^{(j_1)}(\boldsymbol{x}), Q^{(j_1,j_2)}(R^{(j_1)}(\boldsymbol{x}))) \right] \quad (27)$$

is solved, the solution  $R^{(k+1)}$  is added to the set of representations, and the weights are updated to 273 the optimal  $\{p^{(j_1)}\}_{j_1 \in [k+1]}$ . Compared to (24), here the response functions and current k represen-274 tations are kept fixed, and only their weights  $\{p^{(j_1)}\}$  are optimized, along with  $R^{(k+1)}$ . 275

The procedure is described in Algorithm 1, where, following the main loop,  $m^*$ 276  $\arg\min_{k\in[m]} \operatorname{reg}_k$  representation atoms are chosen and the output is  $\{R^{(j)}, p^{(j)}\}_{j\in[m^*]}$ . Algorithm 277 1 relies on solvers for the Phase 1 (26) and Phase 2 (27) problems. In Appendix G we propose two 278 algorithms for these problems, which are based on gradient steps for updating the adversarial re-279 sponse and the new representation, and on the MWU algorithm [33] (follow-the-regularized-leader 280 [35]) for updating the weights. In short, the Phase 1 algorithm updates the response function f via 281 a projected gradient step of the expected loss, and then adjusts the predictors  $\{Q^{(j)}\}$  to the updated 282 response function f and the current representations  $\{R^{(j)}\}_{j \in [k]}$ . The Phase 2 algorithm only up-283 dates the new representation  $R^{(k+1)}$  via projected gradient steps, while keeping  $\{R^{(j)}\}_{j \in [k]}$  fixed. 284 Given the representations  $\{R^{(j)}\}_{j \in [k+1]}$  and the functions  $\{f^{(i)}\}_{i \in [m_0+k]}$ , a predictor  $Q^{(j,i)}$  is then fitted to each representation-function pair, which also determines the loss for this pair. The weights 285 286  $\{p^{(j)}\}_{j \in [k+1]}$  and  $\{o^{(i)}\}_{i \in [m_0+k]}$  are updated towards the equilibrium of the two-player game deter-287 mined by the loss of the predictors  $\{Q^{(j,i)}\}_{i \in [k+1], i \in [m_0+k]}$  via the MWU algorithm.

Algorithm 1 Solver of (24): An iterative algorithm for learning mixed representations.

1: input  $P_{\boldsymbol{x}}, \mathcal{R}, \mathcal{F}, \mathcal{Q}, d, r, m, m_0$ 2: input  $R^{(1)}, \{f^{(j)}\}_{j \in [m_0]}$ > Feature distribution, classes, dimensions and parameters ▷ Initial representation and initial function (set) 3: begin

4: for k = 1 to m do

**phase 1:**  $f^{(m_0+k)}$  is set by a solver of (26) and 5:

$$\operatorname{reg}_{k} \leftarrow \operatorname{regret}_{\mathsf{mix}}(\{R^{(j)}, p^{(j)}\}_{j \in [k]}, \mathcal{F} \mid P_{\boldsymbol{x}})$$

$$(28)$$

▷ Solved using Algorithm 2

- **phase 2:**  $R^{(k+1)}, \{p_k^{(j)}\}_{j \in [k+1]}$  is set by a solver of (27)  $\triangleright$  Solved using Algorithm 2; step can be removed if k = m6:
- 7: end for
- 8: set  $m^* = \arg\min_{k \in [m]} \operatorname{reg}_k$
- 9: return  $\{R^{(j)}\}_{j \in [m^*]}$  and  $p_{m_*} = \{p_k^{(j)}\}_{j \in [m^*]}$

We next outline two examples, where full details can be found in Appendix H. 289

**Example 6.** We validate that efficiency of Algorithm 1 in the linear MSE setting (Section 3), for 290 which a closed-form solution exists. We ran Algorithm 1 on randomly drawn diagonal  $\Sigma_x$ , and 291 computed the ratio between the regret obtained by the algorithm to the theoretical value. The left 292 panel of Figure 2 shows that the ratio is between 1.15 - 1.2 in a wide range of d values. We mention 293 again that Algorithm 1 is useful even for this setting since finding an  $(\ell^* + 1)$ -sparse solution to 294  $\overline{A}p = \overline{b}$  is computationally difficult when  $\binom{\ell^*}{r}$  is very large. For example, in the largest dimension 295 of the experiment, the potential number of representation matrices is  $\binom{d}{r} = \binom{19}{5} = 11,628.$ 296

Our next example pertains to a logistic regression setting, under the cross-entropy loss function. 297

**Definition 7** (The linear cross-entropy setting). Assume that  $\mathcal{X} = \mathbb{R}^d$ , that  $\mathcal{Y} = \{\pm 1\}$  and that 298  $\mathbb{E}[\boldsymbol{x}] = 0$ . Assume that the class of representation is linear  $z = R(x) = R^{\top}x$  for some  $R \in \mathcal{R} :=$ 299

<sup>288</sup> 



Figure 2: Results of Algorithm 1. Left: r = 5, varying d. The ratio between the regret achieved by Algorithm 1 and the theoretical regret in the linear MSE setting. Right: r = 3, varying d. The regret achieved by Algorithm 1 in the linear cross entropy setting, various m.

<sup>300</sup>  $\mathbb{R}^{d \times r}$  where d > r. Assume that a response function and a prediction rule determine the probability <sup>301</sup> that y = 1 via logistic regression modeling, as  $f(\boldsymbol{y} = \pm 1 \mid x) = 1/[1 + \exp(\mp f^{\top} x)]$ . Assume <sup>302</sup> the cross-entropy loss function, where given that the prediction that  $\boldsymbol{y} = 1$  with probability q results <sup>303</sup> the loss loss $(y, q) := -\frac{1}{2}(1+y)\log q - \frac{1}{2}(1-y)\log(1-q)$ . The set of predictor functions is <sup>304</sup>  $\mathcal{Q} := \{Q(z) = 1/[1 + \exp(-q^{\top} \boldsymbol{z})], q \in \mathbb{R}^r\}$ . As for the linear case, we assume that  $f \in \mathcal{F}_S$ <sup>305</sup> for some  $S \in \mathbb{S}^d_{++}$ . It is not difficult to show that the regret is then given by the expected binary <sup>306</sup> Kullback-Leibler (KL) divergence

$$\operatorname{regret}(R, f \mid P_{\boldsymbol{x}}) = \min_{q \in \mathbb{R}^r} \mathbb{E}\left[ D_{\operatorname{KL}}\left( [1 + \exp(-f^{\top} \boldsymbol{x})]^{-1} \mid | [1 + \exp(-q^{\top} R^{\top} \boldsymbol{x})]^{-1} \right) \right].$$
(29)

**Example 8.** We ran Algorithm 1 on empirical distributions of features drawn from an isotropic normal distribution, in the linear cross-entropy setting. Algorithm 1 is suitable in this setting since gradients of the regret have closed-form (see Appendix H). The right panel of Figure 2 shows the reduced regret obtained by increasing the support size m of the random representation, and thus the effectiveness of mixed representations.

<sup>312</sup> We refer the reader to Appendix I for additional experiments with Algorithm 1.

# 313 5 Conclusion

We proposed a game-theoretic formulation for learning representations of unlabeled features when prior knowledge (or assumptions) on the class of future prediction tasks is available. We focused on the fundamental of linear MSE setting, and derived the optimal solution. Beyond the lower regret that is directly obtained from utilizing the prior knowledge, our results also revealed the importance of using randomized representations. We have then proposed an iterative algorithm suitable for general classes of functions and losses, and exemplified its effectiveness.

We next discuss *limitations* and potential future research: (1) We have focused on the elementary and 320 simplified class  $\mathcal{F}_S = \{f : \|f\|_S \le 1\}$ , mainly for theoretical investigations. A natural refinement to 321 non-linear functions is the general class  $\mathcal{F}_{S_x} := \mathbb{E}\left[ \|\nabla_x f(\boldsymbol{x})\|_{S_{\boldsymbol{x}}}^2 \right] \leq 1$ , where  $\{S_x\}_{x \in \mathbb{R}^d}$  is now 322 locally specified (somewhat similarly to the regularization term used in contractive AE [39], though 323 for different reasons). (2) Since the proposed iterative algorithm includes optimization over three 324 325 players, it is of interest to develop version of the algorithm with lower computational optimization cost. (3) We have assumed that  $\mathcal{F}_S$  is given in advance, and a natural follow-up goal is to efficiently 326 learn S from previous experience, e.g., improving S from one episode to another in a meta-learning 327 setup [40]. (4) It is interesting to evaluate the effectiveness of the learned representation in our 328 formulation, as an initialization for further optimization when labeled data is collected. One may 329 330 postulate that since our learned representation is *uniformly* good for all response functions in the 331 class, it may serve as a universal initialization for such training.

#### **Broader impact**

The research described in this paper is foundational, and does not aim for any specific application. 333 Nonetheless, the learned representation is based on a prior assumption on the class of response 334 functions, and the choice of this prior may have positive or negative impacts: For example, a risk 335 of this choice of prior is that the represented features completely ignore a viable feature for making 336 future predictions. A benefit that can stem from choosing a proper prior is that the representation 337 will null the effect of features that lead to unfair advantages for some particular group, in future 338 predictions. Anyhow, the results presented in the paper are indifferent to such future utilization, and 339 any usage of these results should take into account the aforementioned possible implications. 340

#### 341 **References**

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- <sup>343</sup> [2] John N. Tsitsiklis. Decentralized detection. 1989.
- [3] XuanLong Nguyen, Martin J. Wainwright, and Michael I. Jordan. On surrogate loss functions and *f*-divergences. *The Annals of Statistics*, 37(2):876 904, 2009. doi: 10.1214/08-AOS595. URL https://doi.org/10.1214/08-AOS595.
- [4] John Duchi, Khashayar Khosravi, and Feng Ruan. Multiclass classification, information,
   divergence and surrogate risk. *Annals of Statistics*, 46(6B):3246–3275, 2018. ISSN 0090 5364. doi: 10.1214/17-AOS1657.
- [5] Robert E. Schapire and Yoav Freund. *Boosting: Foundations and Algorithms*. The MIT Press, 2012. ISBN 0262017180.
- [6] Karl Pearson. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin philosophical magazine and journal of science*, 2(11):559–572, 1901.
- [7] Ian Jolliffe. Principal component analysis. *Encyclopedia of statistics in behavioral science*,
   2005.
- [8] John P. Cunningham and Zoubin Ghahramani. Linear dimensionality reduction: Survey, insights, and generalizations. *The Journal of Machine Learning Research*, 16(1):2859–2900, 2015.
- [9] Iain M. Johnstone and Debashis Paul. PCA in high dimensions: An orientation. *Proceedings* of the IEEE, 106(8):1277–1292, 2018.
- [10] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component anal ysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.
- [11] Mark A. Kramer. Nonlinear principal component analysis using autoassociative neural net works. *AIChE journal*, 37(2):233–243, 1991.
- [12] Geoffrey E. Hinton and Ruslan R. Salakhutdinov. Reducing the dimensionality of data with
   neural networks. *science*, 313(5786):504–507, 2006.
- [13] Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Y. Ng. Unsupervised learning
   of hierarchical representations with convolutional deep belief networks. *Communications of the ACM*, 54(10):95–103, 2011.
- [14] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol,
   and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep
   network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.
- 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.
- [16] Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method.
   *arXiv preprint physics/0004057*, 2000.

- [17] Gal Chechik, Amir Globerson, Naftali Tishby, and Yair Weiss. Information bottleneck for
   Gaussian variables. *Advances in Neural Information Processing Systems*, 16, 2003.
- [18] Noam Slonim, Nir Friedman, and Naftali Tishby. Multivariate information bottleneck. *Neural computation*, 18(8):1739–1789, 2006.
- [19] Peter Harremoës and Naftali Tishby. The information bottleneck revisited or how to choose
   a good distortion measure. In 2007 IEEE International Symposium on Information Theory,
   pages 566–570. IEEE, 2007.
- [20] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle.
   In 2015 ieee information theory workshop, pages 1–5. IEEE, 2015.
- [21] Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via
   information. *arXiv preprint arXiv:1703.00810*, 2017.
- [22] Ravid Shwartz-Ziv. Information flow in deep neural networks. arXiv preprint
   arXiv:2202.06749, 2022.
- [23] Alessandro Achille and Stefano Soatto. Emergence of invariance and disentanglement in deep representations. *The Journal of Machine Learning Research*, 19(1):1947–1980, 2018.
- [24] Alessandro Achille and Stefano Soatto. Information dropout: Learning optimal representa tions through noisy computation. *IEEE transactions on pattern analysis and machine intelli- gence*, 40(12):2897–2905, 2018.
- [25] Andrew M. Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan D. Tracey, and David D. Cox. On the information bottleneck theory of deep learning.
   *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124020, 2019.
- [26] Bernhard C. Geiger. On information plane analyses of neural network classifiers– A review.
   *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [27] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2006.
   ISBN 0471241954.
- [28] Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. A theory of
   usable information under computational constraints. *arXiv preprint arXiv:2002.10689*, 2020.
- [29] Yann Dubois, Douwe Kiela, David J. Schwab, and Ramakrishna Vedantam. Learning optimal
   representations with the decodable information bottleneck. *Advances in Neural Information Processing Systems*, 33:18674–18690, 2020.
- [30] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive
   predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [31] Ravid Shwartz-Ziv and Yann LeCun. To compress or not to compress–self-supervised learn ing and information theory: A review. *arXiv preprint arXiv:2304.09355*, 2023.
- [32] Guillermo Owen. *Game theory*. Emerald Group Publishing, 2013.
- [33] Yoav Freund and Robert E. Schapire. Adaptive game playing using multiplicative weights.
   *Games and Economic Behavior*, 29(1-2):79–103, 1999.
- [34] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends R in Machine Learning*, 4(2):107–194, 2012.
- [35] Elad Hazan. Introduction to online convex optimization. *Foundations and Trends* (R) *in Optimization*, 2(3-4):157–325, 2016.
- [36] Maurice Sion. On general minimax theorems. 1958.
- [37] Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge university press, 2012.
- [38] Dimitris Bertsimas and John N. Tsitsiklis. *Introduction to linear optimization*, volume 6.
   Athena scientific Belmont, MA, 1997.

- [39] Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio. Contrac tive auto-encoders: Explicit invariance during feature extraction. In *Proceedings of the 28th international conference on international conference on machine learning*, pages 833–840,
   2011.
- [40] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5149–5169, 2021.
- [41] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jen nifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79:
   151–175, 2010.
- [42] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic
   intelligence. In *International conference on machine learning*, pages 3987–3995. PMLR,
   2017.
- [43] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual
   learning. *arXiv preprint arXiv:1710.10628*, 2017.
- 438 [44] Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv* 439 *preprint arXiv:1904.07734*, 2019.
- [45] Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning.
   In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition,
   pages 11254–11263, 2019.
- [46] Stephen Boyd, Stephen P. Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge
   university press, 2004.
- [47] Ohad Shamir, Sivan Sabato, and Naftali Tishby. Learning and generalization with the information bottleneck. *Theoretical Computer Science*, 411(29-30):2696–2711, 2010.
- [48] XuanLong Nguyen, Martin J. Wainwright, and Michael I. Jordan. Estimating divergence
   functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on In- formation Theory*, 56(11):5847–5861, 2010.
- [49] Ben Poole, Sherjil Ozair, Aäron van den Oord, Alexander A. Alemi, and George Tucker. On
   variational lower bounds of mutual information. In *NeurIPS Workshop on Bayesian Deep Learning*, 2018.
- [50] Tailin Wu, Ian Fischer, Isaac L. Chuang, and Max Tegmark. Learnability for the information
   bottleneck. In *Uncertainty in Artificial Intelligence*, pages 1050–1060. PMLR, 2020.
- [51] David McAllester and Karl Stratos. Formal limitations on the measurement of mutual infor mation. In *International Conference on Artificial Intelligence and Statistics*, pages 875–884.
   PMLR, 2020.
- [52] Matthew Chalk, Olivier Marre, and Gasper Tkacik. Relevant sparse codes with variational
   information bottleneck. *Advances in Neural Information Processing Systems*, 29, 2016.
- [53] Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational
   information bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- [54] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio,
   Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference on machine learning*, pages 531–540. PMLR, 2018.
- [55] Behrooz Razeghi, Flavio P. Calmon, Deniz Gunduz, and Slava Voloshynovskiy. Bottlenecks
   CLUB: Unifying information-theoretic trade-offs among complexity, leakage, and utility.
   *arXiv preprint arXiv:2207.04895*, 2022.
- [56] Matías Vera, Pablo Piantanida, and Leonardo Rey Vega. The role of information complexity
   and randomization in representation learning. *arXiv preprint arXiv:1802.05355*, 2018.

- [57] Borja Rodriguez Galvez. The information bottleneck: Connections to other problems, learn ing and exploration of the ib curve, 2019.
- [58] Peter L. Bartlett and Shahar Mendelson. Rademacher and Gaussian complexities: Risk
   bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [59] Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48.
   Cambridge University Press, 2019.
- [60] Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to
   *algorithms*. Cambridge university press, 2014.
- [61] Karthik Sridharan and Sham M. Kakade. An information theoretic framework for multi-view
   learning. 2008.
- [62] Rana Ali Amjad and Bernhard C. Geiger. Learning representations for neural network-based
   classification using the information bottleneck principle. *IEEE transactions on pattern anal- ysis and machine intelligence*, 42(9):2225–2239, 2019.
- [63] Artemy Kolchinsky, Brendan D. Tracey, and David H. Wolpert. Nonlinear information bot tleneck. *Entropy*, 21(12):1181, 2019.
- [64] D. J. Strouse and David J. Schwab. The information bottleneck and geometric clustering.
   *Neural computation*, 31(3):596–612, 2019.
- [65] Ankit Pensia, Varun Jog, and Po-Ling Loh. Extracting robust and accurate features via a
   robust information bottleneck. *IEEE Journal on Selected Areas in Information Theory*, 1(1):
   131–144, 2020.
- [66] Shahab Asoodeh and Flavio P Calmon. Bottleneck problems: An information and estimation theoretic view. *Entropy*, 22(11):1325, 2020.
- [67] Vudtiwat Ngampruetikorn and David J. Schwab. Perturbation theory for the information
   bottleneck. *Advances in Neural Information Processing Systems*, 34:21008–21018, 2021.
- [68] Xi Yu, Shujian Yu, and José C Príncipe. Deep deterministic information bottleneck with
   matrix-based entropy functional. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3160–3164. IEEE, 2021.
- [69] Vudtiwat Ngampruetikorn and David J. Schwab. Information bottleneck theory of
   high-dimensional regression: Relevancy, efficiency and optimality. *arXiv preprint arXiv:2208.03848*, 2022.
- [70] Deniz Gündüz, Zhijin Qin, Inaki Estella Aguerri, Harpreet S Dhillon, Zhaohui Yang, Aylin
   Yener, Kai Kit Wong, and Chan-Byoung Chae. Beyond transmitting bits: Context, semantics,
   and task-oriented communications. *IEEE Journal on Selected Areas in Communications*, 41 (1):5–41, 2022.
- [71] Vudtiwat Ngampruetikorn and David J. Schwab. Generalized information bottleneck for Gaussian variables. *arXiv preprint arXiv:2303.17762*, 2023.
- [72] Vudtiwat Ngampruetikorn, William Bialek, and David Schwab. Information-bottleneck
   renormalization group for self-supervised representation learning. *Bulletin of the American Physical Society*, 65, 2020.
- [73] William B. Johnson. Extensions of Lipschitz mappings into a Hilbert space. *Contemp. Math.*,
   26:189–206, 1984.
- [74] Santosh S. Vempala. *The random projection method*, volume 65. American Mathematical
   Soc., 2005.
- <sup>513</sup> [75] Michael W. Mahoney et al. Randomized algorithms for matrices and data. *Foundations and* <sup>514</sup> *Trends*® *in Machine Learning*, 3(2):123–224, 2011.

- [76] David P. Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends* (*R*) *in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- [77] Fan Yang, Sifan Liu, Edgar Dobriban, and David P. Woodruff. How to reduce dimension with
   PCA and random projections? *IEEE Transactions on Information Theory*, 67(12):8154–8189,
   2021.
- [78] John F. Nash Jr. Equilibrium points in n-person games. *Proceedings of the national academy* of sciences, 36(1):48–49, 1950.
- [79] Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and
   equilibrium in generative adversarial nets (gans). In *International Conference on Machine Learning*, pages 224–232. PMLR, 2017.
- [80] Paulina Grnarova, Kfir Y. Levy, Aurelien Lucchi, Thomas Hofmann, and Andreas
   Krause. An online learning approach to generative adversarial networks. *arXiv preprint arXiv:1706.03269*, 2017.
- [81] Ilya O. Tolstikhin, Sylvain Gelly, Olivier Bousquet, Carl-Johann Simon-Gabriel, and Bernhard Schölkopf. Adagan: Boosting generative models. *Advances in neural information processing systems*, 30, 2017.
- [82] Max Welling, Richard Zemel, and Geoffrey E. Hinton. Self supervised boosting. *Advances in neural information processing systems*, 15, 2002.
- [83] Abraham Wald. Contributions to the theory of statistical estimation and testing hypotheses.
   *The Annals of Mathematical Statistics*, 10(4):299–326, 1939.
- [84] Larry Wasserman. All of statistics: A concise course in statistical inference, volume 26.
   Springer, 2004.
- [85] Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *Annals of Statistics*, pages 1564–1599, 1999.
- [86] Peter D. Grünwald and A. Philip Dawid. Game theory, maximum entropy, minimum discrepancy and robust bayesian decision theory. *The Annals of Statistics*, 32(4):1367–1433, 2004.
- [87] David Haussler and Manfred Opper. Mutual information, metric entropy and cumulative
   relative entropy risk. *The Annals of Statistics*, 25(6):2451–2492, 1997.
- [88] Farzan Farnia and David Tse. A minimax approach to supervised learning. Advances in Neural Information Processing Systems, 29, 2016.
- [89] Jorge Silva and Felipe Tobar. On the interplay between information loss and operation loss in
   representations for classification. In *International Conference on Artificial Intelligence and Statistics*, pages 4853–4871. PMLR, 2022.
- [90] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
   Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [91] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and
   Anil A. Bharath. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1):53–65, 2018.
- [92] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian
   Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [93] Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust optimization*, vol ume 28. Princeton university press, 2009.

- [94] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and
   Xi Chen. Improved techniques for training GANs. *Advances in neural information pro- cessing systems*, 29, 2016.
- [95] Constantinos Daskalakis, Alan Deckelbaum, and Anthony Kim. Near-optimal no-regret al gorithms for zero-sum games. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 235–254. SIAM, 2011.
- [96] Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable
   sequences. Advances in Neural Information Processing Systems, 26, 2013.
- [97] James P. Bailey and Georgios Piliouras. Multiplicative weights update in zero-sum games. In
   *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 321–338,
   2018.
- [98] Guodong Zhang, Yuanhao Wang, Laurent Lessard, and Roger B. Grosse. Near-optimal local
   convergence of alternating gradient descent-ascent for minimax optimization. In *Interna- tional Conference on Artificial Intelligence and Statistics*, pages 7659–7679. PMLR, 2022.
- [99] Florian Schäfer and Anima Anandkumar. Competitive gradient descent. Advances in Neural Information Processing Systems, 32, 2019.
- [100] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. The numerics of GANs. Advances
   *in neural information processing systems*, 30, 2017.
- [101] Alistair Letcher, David Balduzzi, Sébastien Racaniere, James Martens, Jakob Foerster, Karl
   Tuyls, and Thore Graepel. Differentiable game mechanics. *The Journal of Machine Learning Research*, 20(1):3032–3071, 2019.
- [102] Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezeshki, Rémi Le Priol, Gabriel
   Huang, Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved
   game dynamics. In *The 22nd International Conference on Artificial Intelligence and Statis- tics*, pages 1802–1811. PMLR, 2019.
- [103] Guodong Zhang and Yuanhao Wang. On the suboptimality of negative momentum for minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 2098–2106. PMLR, 2021.
- [104] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
   Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [105] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [106] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [107] Dimitri Bertsekas, Angelia Nedic, and Asuman Ozdaglar. *Convex analysis and optimization*,
   volume 1. Athena Scientific, 2003.
- [108] Ky Fan. On a theorem of Weyl concerning eigenvalues of linear transformations i. *Proceed-ings of the National Academy of Sciences*, 35(11):652–655, 1949.