

600 The outline of the supplementary material is as follows. In Appendix A we discuss in detail the  
601 origin of prior knowledge of classes of prediction problems. In Appendix B we review additional  
602 related work. In Appendix C we set our notation conventions. In Appendix D we summarize a few  
603 mathematical results that are used in later proofs. In Section E we show that PCA can be cast as a  
604 degenerate setting of our formulation, and provide the proofs of the main theorems in the paper (the  
605 linear MSE setting). In Appendix F we generalize these results to an infinite dimensional Hilbert  
606 space. In Appendix G we provide two algorithms for solving the Phase 1 and Phase 2 problems  
607 in Algorithm 1. In Appendix H we provide details on the examples for the experiments with the  
608 iterative algorithm. In Appendix I we describe an experiment in which the representation, response  
609 function and predictors are modeled as a neural network (NN).

## 610 A Classes of response functions

611 As said, our approach to optimal representation is based on the assumption that a class  $\mathcal{F}$  of future  
612 prediction tasks is known. This assumption may represent prior knowledge or constraints on the  
613 response function, and can stem from various considerations. To begin, it might be hypothesized  
614 that some features are less relevant than others. As a simple intuitive example, the outer pixels in  
615 images are typically less relevant to the classification of photographed objects, regardless of their  
616 variability (which may stem from other affects, such as lighting conditions). Similarly, non-coding  
617 regions of the genotype are irrelevant for predicting phenotype. The prior knowledge may encode  
618 softer variations in relevance. Moreover, such prior assumption may be imposed on the learned  
619 function, e.g., it may be assumed that the response function respects the privacy of some features,  
620 or only weakly depends on features which provide an unfair advantage. In domain adaptation [?  
621 ], one may solve the prediction problem for feature distribution  $P_{\mathbf{x}}$  obtaining a optimal response  
622 function  $f_1$ . Then, after a change of input distribution to  $Q_{\mathbf{x}}$ , the response function learned for this  
623 feature distribution  $f_2$  may be assumed to belong to functions which are “compatible” with  $f_1$ . For  
624 example, if  $P_{\mathbf{x}}$  and  $Q_{\mathbf{x}}$  are supported on different subsets of  $\mathbb{R}^d$ , the learned response function  $f_1(x)$   
625 and  $f_2(x)$  may be assumed to satisfy some type of continuity assumptions. Similar assumptions may  
626 hold for the more general setting of transfer learning [41]. Furthermore, such assumptions may hold  
627 in a *continual learning* setting [42–45], in which a sequence of response functions is learned one  
628 task at a time. Assuming that *catastrophic forgetting* is aimed to be avoided, then starting from the  
629 second task, the choice of representation may assume that the learned response function is accurate  
630 for all previously learned tasks.

## 631 B Additional related work

632 **The information bottleneck principle** The IB principle is a prominent approach to feature rele-  
633 vance in the design of representations [16–19], and proposes to optimize the representation in order  
634 to maximize its relevance to the response  $\mathbf{y}$ . Letting  $I(\mathbf{z}; \mathbf{y})$  and  $I(\mathbf{x}; \mathbf{z})$  denote the corresponding  
635 mutual information terms [27], the IB principle aims to maximize the former while constraining the  
636 latter from above, and this is typically achieved via a Lagrangian formulation [46]. The resulting  
637 representation, however, is tailored to the joint distribution of  $(\mathbf{x}, \mathbf{y})$ , i.e., to a specific prediction  
638 task. In practice, this is achieved using a labeled dataset (Generalization bounds were derived in  
639 [47]). As in our mixed representation approach, the use of randomized representation dictated by a  
640 probability kernel  $P_{Z|X}$  is inherent to the IB principle. The IB principle was intensively utilized to  
641 hypothesize that prediction algorithms, e.g., deep neural networks (DNNs) [1] used for classifica-  
642 tion, must intrinsically include learning of efficient representations [20–24] (this spurred a debate,  
643 see, e.g., [25, 26]). However, this approach is inadequate in an unsupervised setting since the opti-  
644 mal representation depends on the response variable, and so labeled data should be provided when  
645 learning the representation. In addition, as explained in [29], while the resulting IB solution provides  
646 a fundamental limit for the problem, it also suffers from multiple theoretical and practical issues.  
647 The first main issue is that the mutual information terms are inherently difficult to estimate from  
648 finite samples [47–51], especially at high dimensions, and thus require resorting to approximations,  
649 e.g., variational bounds [52–55]. The resulting generalization bounds [47, 56] are still vacuous for  
650 modern settings [57]. The second main issue is that the IB formulation does not constrain the com-  
651 plexity of the representation and the prediction rule, which can be arbitrarily complex. These issues  
652 were addressed in [29] using the notion of *usable information*, introduced in [28]: The standard  
653 mutual information  $I(\mathbf{z}; \mathbf{y})$  can be described as the log-loss difference between a predictor for  $\mathbf{z}$

654 which does not use or does use  $\mathbf{y}$  (or vice-versa, since mutual information is symmetric). Usable  
655 information, or  $\mathcal{F}$ -information  $I_{\mathcal{F}}(\mathbf{z} \rightarrow \mathbf{y})$ , restricts the predictor to a class  $\mathcal{F}$ , which is compu-  
656 tationally constrained. Several desirable properties were established in [28] for the  $\mathcal{F}$ -information,  
657 e.g., probably approximate correct (PAC) bounds via Rademacher-complexity based bounds [58]  
658 [59, Chapter 5][60, Chapters 26-28]. In [29], the authors used the notion of  $\mathcal{F}$ -information to de-  
659 fine the *decodable IB* problem, with the goal of addressing the generalization capabilities of this IB  
660 problem. In order to explore this, the *two-player game* described in the introduction was proposed.  
661 Beyond those works, the IB framework has drawn a significant recent attention, and a remarkable  
662 number of extensions and ramifications have been proposed [61–70, 55, 71]. IB framework for  
663 self-supervised learning was recently discussed in [72].

664 **Randomization in representation learning** Randomization is classically used in data represen-  
665 tation, most notably, utilizing the seminal Johnson-Lindenstrauss Lemma [73] or more generally,  
666 *sketching* algorithms (e.g., [74–77]). Our use of randomization is different and is inspired by the  
667 classical Nash equilibrium [78]. Rather than using a single deterministic representation that was ran-  
668 domly chosen, we consider randomizing multiple representation rules. Training approaches based  
669 on mixed strategies were proposed, e.g., in the generative adversarial network (GAN) setting [79–  
670 81]. Specifically, inspired by the boosting technique [5], it was proposed in [81] to gradually add  
671 additional modes to the mix of generative models, and where the new mode added focuses on the  
672 distribution samples which are not adequately represented by the current set of modes. As men-  
673 tioned in [81], this idea dates back to the use of boosting for density estimation [82]. Our proposed  
674 iterative algorithm follows this idea, and gradually adds representation rules, so that the new rep-  
675 resentation aims to cope with response functions that are not adequately fitted by the current set of  
676 representation rules.

677 **Game theoretic formulations in statistics and machine-learning** The use of game theoretic for-  
678 mulations in statistics, between a player choosing a prediction algorithm and an adversary choosing  
679 a prediction problem (typically Nature), was established by Wald in his classical statistical decision  
680 theory [83] (see, e.g., [84, Chapter 12]). It is a common approach both in classic statistics and  
681 learning theory [85–88], as well as in modern high-dimensional statistics [59]. The effect of the  
682 representation (quantizer) on the consistency of learning algorithms when a surrogate convex loss  
683 function replaces the loss function of interest was studied in [3, 4, 86] (for binary and multiclass  
684 classification, respectively). A relation between information loss and minimal error probability was  
685 recently derived in [89].

686 Iterative algorithms for the solution of minimax games have drawn much attention in the last few  
687 years due to their importance in optimizing GANs [90, 91], adversarial training [92], and robust  
688 optimization [93]. The notion of convergence is rather delicate, even for the basic convex-concave  
689 two-player setting [94]. While the value output by the MWU algorithm [33], or improved versions  
690 [95, 96] converges to a no-regret solution, the actual strategies used by the players are, in fact,  
691 repelled away from the equilibrium point to the boundary of the probability simplex [97]. For  
692 general games, the gradient descent ascent (GDA) is a natural and practical choice, yet despite recent  
693 advances, its theory is still partial [98]. Various other algorithms have been proposed, e.g., [99–103].  
694 According to the above description, and since our algorithm is both fairly general and involves two  
695 optimization phases, deriving theoretical bounds on its convergence seems to be elusive at this point.  
696 Nevertheless, the algorithm is also modular, and its two intermediate phases (see Appendix G) can  
697 be easily upgraded to more sophisticated optimization methods. Furthermore, each of the phases  
698 can be separately analyzed.

699 **Unsupervised pretraining** From a broader perspective, our method is essentially an *unsupervised*  
700 *pretraining* method, similar to the methods which currently enable the recent success in natural  
701 language processing. Our model is much simplified compared to transformer architecture Vaswani  
702 et al. [104], but the unsupervised training aspect used for prediction tasks Devlin et al. [105] is  
703 common, and our results may shed light on these methods. For example, putting more weight  
704 on some words compared to others during training phase that uses the masked-token prediction  
705 objective.

706 **C Notation conventions**

707 For an integer  $d$ ,  $[d] := \{1, 2, \dots, d\}$ . For  $p \geq 1$ ,  $\|x\|_p := (\sum_{i=1}^d |x_i|^p)^{1/p}$  is the  $\ell_p$  norm of  
 708  $x \in \mathbb{R}^d$ . The Frobenius norm of the matrix  $A$  is denoted by  $\|A\|_F = \sqrt{\text{Tr}[A^T A]}$ . The non-  
 709 negative (resp. positive) definite cone of symmetric matrices is given by  $\mathbb{S}_+^d$  (resp.  $\mathbb{S}_{++}^d$ ). For a  
 710 given positive-definite matrix  $S \in \mathbb{S}_{++}^d$ , the Mahalanobis norm of  $x \in \mathbb{R}^d$  is given by  $\|x\|_S :=$   
 711  $\|S^{-1/2}x\|_2 = (x^T S^{-1}x)^{1/2}$ , where  $S^{1/2}$  is the symmetric square root of  $S$ . The matrix  $W :=$   
 712  $[w_1, \dots, w_r] \in \mathbb{R}^{d \times r}$  is comprised from the column vectors  $\{w_i\}_{i \in [r]} \subset \mathbb{R}^d$ . For a real symmetric  
 713 matrix  $S \in \mathbb{S}^d$ ,  $\lambda_i(S)$  is the  $i$ th largest eigenvalue, so that  $\lambda_{\max}(S) \equiv \lambda_1(S) \geq \lambda_2(S) \geq \dots \geq$   
 714  $\lambda_d(S) = \lambda_{\min}(S)$ , and in accordance,  $v_i(S)$  denote an eigenvector corresponding to  $\lambda_i(S)$  (these  
 715 are unique if there are no two equal eigenvalues, and otherwise arbitrarily chosen, while satisfying  
 716 orthogonality  $v_i^T v_j = \langle v_i, v_j \rangle = \delta_{ij}$ ). Similarly,  $\Lambda(S) := \text{diag}(\lambda_1(S), \lambda_2(S), \dots, \lambda_d(S))$  and  
 717  $V(S) := [v_1(S), v_2(S), \dots, v_d(S)]$ , so that  $S = V(S)\Lambda(S)V^T(S)$  is an eigenvalue decomposition.  
 718 For  $j \geq i$ ,  $V_{i:j} := [v_i, \dots, v_j] \in \mathbb{R}^{(j-i+1) \times d}$  is the matrix comprised of the columns indexed by  
 719  $\{i, \dots, j\}$ . The vector  $e_i \in \mathbb{R}^d$  is the  $i$ th standard basis vector, that is,  $e_i := [\underbrace{0, \dots, 0}_{i-1 \text{ terms}}, 1, \underbrace{0, \dots, 0}_{d-i \text{ terms}}]^T$ .

720 Random quantities (scalars, vectors, matrices, etc.) are denoted by boldface letters. For example,  
 721  $\mathbf{x} \in \mathbb{R}^d$  is a random vector that takes values  $x \in \mathbb{R}^d$  and  $\mathbf{R} \in \mathbb{R}^{d \times r}$  is a random matrix.  
 722 The probability law of a random element  $\mathbf{x}$  is denoted by  $L(\mathbf{x})$ . The probability of the event  $\mathcal{E}$  in some  
 723 given probability space is denoted by  $\mathbb{P}[\mathcal{E}]$  (typically understood from context). The expectation  
 724 operator is denoted by  $\mathbb{E}[\cdot]$ . The indicator function is denoted by  $\mathbb{1}\{\cdot\}$ , and the Kronecker delta is  
 725 denoted by  $\delta_{ij} := \mathbb{1}\{i = j\}$ . We do not make a distinction between minimum and infimum (or  
 726 maximum and supremum) as arbitrarily accurate approximation is sufficient for the description of  
 727 the results in this paper. The binary KL divergence between  $p_1, p_2 \in (0, 1)$  is denoted as

$$D_{\text{KL}}(p_1 \parallel p_2) := p_1 \log \frac{p_1}{p_2} + (1 - p_1) \log \frac{1 - p_1}{1 - p_2}. \quad (30)$$

728 **D Useful mathematical results**

729 In this section we provide several simplified versions of mathematical results that are used in the  
 730 proofs. The following well-known result is about the optimal low-rank approximation to a given  
 731 matrix:

732 **Theorem 9** (Eckart-Young-Mirsky [59, Example 8.1] [106, Section 4.1.4]). *For a symmetric matrix*  
 733  $S \in \mathbb{S}^d$

$$\|S_k - S\|_F \leq \min_{S' \in \mathbb{S}^d: \text{rank}(S') \leq k} \|S - S'\|_F \quad (31)$$

734 *where*

$$S_k = \sum_{i \in [k]} \lambda_i(S) \cdot v_i(S)v_i^T(S) \quad (32)$$

735 *(more generally, this is true for any unitarily invariant norm).*

736 We next review a simplified version of variational characterizations of eigenvalues of symmetric  
 737 matrices:

738 **Theorem 10** (Rayleigh quotient [37, Theorem 4.2.2]). *For a symmetric matrix  $S \in \mathbb{S}^d$*

$$\lambda_1(S) = \max_{x \neq 0} \frac{x^T S x}{\|x\|_2^2}. \quad (33)$$

739 **Theorem 11** (Courant-Fisher variational characterization [37, Theorem 4.2.6]). *For a symmetric*  
 740 *matrix  $S \in \mathbb{S}^d$ ,  $k \in [d]$ , and a subspace  $T$  of  $\mathbb{R}^d$*

$$\lambda_k(S) = \min_{T: \dim(T)=k} \max_{x \in T \setminus \{0\}} \frac{x^T S x}{\|x\|_2^2} = \max_{T: \dim(T)=d-k+1} \min_{x \in T \setminus \{0\}} \frac{x^T S x}{\|x\|_2^2}. \quad (34)$$

741 **Theorem 12** (Fan’s variational characterization [37, Corollary 4.3.39.]). For a symmetric matrix  
 742  $S \in \mathbb{S}^d$  and  $k \in [d]$

$$\lambda_1(S) + \cdots + \lambda_k(S) = \min_{U \in \mathbb{R}^{d \times k}: U^\top U = I_k} \text{Tr}[U^\top S U] \quad (35)$$

743 and

$$\lambda_{d-k+1}(S) + \cdots + \lambda_d(S) = \max_{U \in \mathbb{R}^{d \times k}: U^\top U = I_k} \text{Tr}[U^\top S U]. \quad (36)$$

744 We will use the following celebrated result from convex analysis.

745 **Theorem 13** (Carathéodory’s theorem [107, Prop. 1.3.1]). Let  $\mathcal{A} \subset \mathbb{R}^d$  be non-empty. Then, any  
 746 point  $a$  in the convex hull of  $\mathcal{A}$  can be written as a convex combination of at most  $d + 1$  points from  
 747  $\mathcal{A}$ .

## 748 E The linear MSE setting: additions and proofs

### 749 E.1 The standard principal component setting

750 In order to highlight the formulation proposed in this paper, we show, as a starting point, that the  
 751 well known PCA solution of representing  $\mathbf{x} \in \mathbb{R}^d$  with the top  $r$  eigenvectors of the covariance  
 752 matrix of  $\mathbf{x}$  can be obtained as a specific case of the regret formulation. In this setting, we take  
 753  $\mathcal{F} = \{I_d\}$ , and so  $\mathbf{y} = \mathbf{x}$  with probability 1. In addition, the predictor class  $\mathcal{Q}$  is a linear function  
 754 from the representation dimension  $r$  back to the features dimension  $d$ .

755 **Proposition 14.** Consider the linear MSE setting, with the difference that the response is  $\mathbf{y} \in \mathbb{R}^d$ ,  
 756 the loss function is the MSE  $\text{loss}(y_1, y_2) = \|y_1 - y_2\|^2$ , and the predictor is  $Q(z) = Q^\top z \in \mathbb{R}^d$  for  
 757  $Q \in \mathbb{R}^{r \times d}$ . Assume  $\mathcal{F} = \{I_d\}$  so that  $\mathbf{y} = \mathbf{x}$  with probability 1. Then,

$$\text{regret}_{\text{pure}}(\mathcal{F} \mid \Sigma_{\mathbf{x}}) = \text{regret}_{\text{mix}}(\mathcal{F} \mid \Sigma_{\mathbf{x}}) = \sum_{i=r+1}^d \lambda_i(\Sigma_{\mathbf{x}}), \quad (37)$$

758 and an optimal representation is  $R = V_{1:r}(\Sigma_{\mathbf{x}})$ .

759 The result of Proposition 14 verifies that the minimax and maximin formulations indeed generalize  
 760 the standard PCA formulation. The proof is standard and follows from the *Eckart-Young-Mirsky the-*  
 761 *orem* (e.g., [59, Example 8.1] [106, Section 4.1.4]), which determines the best rank  $r$  approximation  
 762 in the Frobenius norm.

763 *Proof of Proposition 14.* Since  $\mathcal{F} = \{I_d\}$  is a singleton, there is no distinction between pure and  
 764 mixed minimax regret. It holds that

$$\text{regret}(R, f) = \mathbb{E} [\|\mathbf{x} - Q^\top R^\top \mathbf{x}\|^2] \quad (38)$$

765 where  $A = Q^\top R^\top \in \mathbb{R}^{d \times d}$  is a rank  $r$  matrix. For any  $A \in \mathbb{R}^{d \times d}$

$$\mathbb{E} [\|\mathbf{x} - A\mathbf{x}\|^2] = \mathbb{E} [\mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top A\mathbf{x} - \mathbf{x}^\top A^\top \mathbf{x} + \mathbf{x}^\top A^\top A\mathbf{x}] \quad (39)$$

$$= \text{Tr} [\Sigma_{\mathbf{x}} - A\Sigma_{\mathbf{x}} - A^\top \Sigma_{\mathbf{x}} + A^\top A\Sigma_{\mathbf{x}}] \quad (40)$$

$$= \left\| \Sigma_{\mathbf{x}}^{1/2} - \Sigma_{\mathbf{x}}^{1/2} A \right\|_F^2 \quad (41)$$

$$= \left\| \Sigma_{\mathbf{x}}^{1/2} - B \right\|_F^2, \quad (42)$$

766 where  $B := \Sigma_{\mathbf{x}}^{1/2} A$ . By the classic *Eckart-Young-Mirsky theorem* [59, Example 8.1] [106, Section  
 767 4.1.4] (see Appendix D), the best rank  $r$  approximation in the Frobenius norm is obtained by setting

$$B^* = \sum_{i=1}^r \lambda_i(\Sigma_{\mathbf{x}}^{1/2}) \cdot v_i v_i^\top = \sum_{i=1}^r \sqrt{\lambda_i(\Sigma_{\mathbf{x}})} \cdot v_i v_i^\top \quad (43)$$

768 where  $v_i \equiv v_i(\Sigma_{\mathbf{x}}^{1/2}) = v_i(\Sigma_{\mathbf{x}})$  is the  $i$ th eigenvector of  $\Sigma_{\mathbf{x}}^{1/2}$  (or  $\Sigma_{\mathbf{x}}$ ). Then, the optimal  $A$  is

$$A^* = \sum_{i=1}^r \sqrt{\lambda_i(\Sigma_{\mathbf{x}})} \cdot \Sigma_{\mathbf{x}}^{-1/2} v_i v_i^\top = \sum_{i=1}^r \sqrt{\lambda_i(\Sigma_{\mathbf{x}})} \cdot \Sigma_{\mathbf{x}}^{-1/2} v_i v_i^\top = \sum_{i=1}^r v_i v_i^\top, \quad (44)$$

769 since  $v_i$  is also an eigenvector of  $\Sigma_{\mathbf{x}}^{-1/2}$ . Letting  $R = U(R)\Sigma(R)V^\top(R)$  and  $Q =$   
 770  $U(Q)\Sigma(Q)V^\top(Q)$  be the singular value decomposition of  $R$  and  $Q$ , respectively, it holds that

$$Q^\top R^\top = V(Q)\Sigma^\top(Q)V(Q)V(R)\Sigma^\top(R)U^\top(R). \quad (45)$$

771 Setting  $V(Q) = V(R) = I_r$ , and  $\Sigma^\top(Q) = \Sigma(R) \in \mathbb{R}^{d \times r}$  to have  $r$  ones on the diagonal (and all  
 772 other entries are zero), as well as  $U(Q) = U(R)$  to be an orthogonal matrix whose first  $r$  columns  
 773 are  $\{v_i\}_{i \in [r]}$  results that  $Q^\top R^\top = A^*$ , as required.  $\square$

## 774 E.2 Proofs of pure and mixed minimax representations

775 Before the proof of Theorem 2, we state a simple and useful lemma, which provides the pointwise  
 776 value of the regret and the optimal linear predictor for a given representation and response.

777 **Lemma 15.** *Consider the representation  $\mathbf{z} = R^\top \mathbf{x} \in \mathbb{R}^r$ . It then holds that*

$$\min_{q \in \mathbb{R}^r} \mathbb{E} \left[ (f^\top \mathbf{x} + \mathbf{n} - q^\top \mathbf{z})^2 \right] \quad (46)$$

$$= \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right] + f^\top (\Sigma_{\mathbf{x}} - \Sigma_{\mathbf{x}} R (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}}) f. \quad (47)$$

778 *Proof.* The orthogonality principle states that

$$\mathbb{E} \left[ (f^\top \mathbf{x} + \mathbf{n} - q^\top \mathbf{z}) \cdot \mathbf{z}^\top \right] = 0 \quad (48)$$

779 must hold for the optimal linear estimator. Using  $\mathbf{z} = R^\top \mathbf{x}$  and taking expectations leads to the  
 780 standard least-squares (LS) solution

$$q_{\text{LS}} = (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}} f, \quad (49)$$

781 assuming that  $R^\top \Sigma_{\mathbf{x}} R$  is invertible (which we indeed assume as if this is not the case, the represen-  
 782 tation can be reduced to a dimension lower than  $r$  in a lossless manner). The resulting regret of  $R$  is  
 783 thus given by

$$\text{regret}(R, f) = \mathbb{E} \left[ (f^\top \mathbf{x} + \mathbf{n} - q_{\text{LS}}^\top \mathbf{z})^2 \right] \quad (50)$$

$$\stackrel{(a)}{=} \mathbb{E} \left[ (f^\top \mathbf{x} + \mathbf{n})^\top (f^\top \mathbf{x} + \mathbf{n} - q_{\text{LS}}^\top \mathbf{z}) \right] \quad (51)$$

$$= \mathbb{E} \left[ (f^\top \mathbf{x} + \mathbf{n})^2 - (f^\top \mathbf{x} + \mathbf{n})^\top q_{\text{LS}}^\top \mathbf{z} \right] \quad (52)$$

$$\stackrel{(b)}{=} \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right] + f^\top \Sigma_{\mathbf{x}} f - \mathbb{E} \left[ \mathbf{x}^\top f q_{\text{LS}}^\top R^\top \mathbf{x} \right] \quad (53)$$

$$= \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right] + f^\top \Sigma_{\mathbf{x}} f - \text{Tr} \left[ f q_{\text{LS}}^\top R^\top \Sigma_{\mathbf{x}} \right] \quad (54)$$

$$= \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right] + f^\top \Sigma_{\mathbf{x}} f - q_{\text{LS}}^\top R^\top \Sigma_{\mathbf{x}} f \quad (55)$$

$$\stackrel{(c)}{=} \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right] + f^\top (\Sigma_{\mathbf{x}} - \Sigma_{\mathbf{x}} R (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}}) f, \quad (56)$$

784 where (a) follows from the orthogonality principle in (48), (b) follows from the tower property of  
 785 conditional expectation and since  $\mathbb{E}[\mathbf{x}\mathbf{n}] = \mathbb{E}[\mathbf{x} \cdot \mathbb{E}[\mathbf{n} \mid \mathbf{x}]] = 0$ , and (c) follows by substituting  $q_{\text{LS}}$   
 786 from (49).  $\square$

787 We may now prove Theorem 2.

788 *Proof of Theorem 2.* For any given  $f$ , the optimal predictor based on  $\mathbf{x} \in \mathbb{R}^d$  achieves average loss  
 789 of

$$\min_{q \in \mathbb{R}^d} \mathbb{E} \left[ (f^\top \mathbf{x} + \mathbf{n} - q^\top \mathbf{x})^2 \right] = \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right] \quad (57)$$

790 (specifically, this is obtained by setting  $R = I_d$  in Lemma 15 so that  $\mathbf{z} = \mathbf{x}$ ). Hence, the resulting  
 791 regret of  $R$  over an adversarial choice of  $f \in \mathcal{F}_S$  is

$$\max_{f \in \mathcal{F}_S} \text{regret}(R, f) = \max_{f \in \mathcal{F}} \mathbb{E} \left[ |f^\top \mathbf{x} + \mathbf{n} - q_{\text{LS}}^\top \mathbf{z}|^2 \right] - \mathbb{E} \left[ \mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}] \right]$$

$$\stackrel{(a)}{=} \max_{f \in \mathcal{F}_S} f^\top (\Sigma_{\mathbf{x}} - \Sigma_{\mathbf{x}} R (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}}) f \quad (58)$$

$$\stackrel{(b)}{=} \max_{\tilde{f}: \|\tilde{f}\|_S^2 \leq 1} \tilde{f}^\top \left( S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} - S^{1/2} \Sigma_{\mathbf{x}} R (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}} S^{1/2} \right) \tilde{f} \quad (59)$$

$$\stackrel{(c)}{=} \lambda_1 \left( S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} - S^{1/2} \Sigma_{\mathbf{x}} R (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}} S^{1/2} \right) \quad (60)$$

$$= \lambda_1 \left[ S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \Sigma_{\mathbf{x}}^{1/2} R (R^\top \Sigma_{\mathbf{x}} R)^{-1} R^\top \Sigma_{\mathbf{x}}^{1/2} \right) \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \right] \quad (61)$$

$$\stackrel{(d)}{=} \lambda_1 \left[ S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \right] \quad (62)$$

$$\stackrel{(e)}{=} \lambda_1 \left[ \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) \right], \quad (63)$$

792 where (a) follows from Lemma 15, (b) follows by letting  $\tilde{f} := S^{-1/2} f$  and recalling that any  
 793  $f \in \mathcal{F}$  must satisfy  $\|f\|_S^2 \leq 1$ , (c) follows from the *Rayleigh quotient theorem* [37, Theorem 4.2.2]  
 794 (see Appendix D), (d) follows by letting  $\tilde{R} := \Sigma_{\mathbf{x}}^{1/2} R$ , and (e) follows since  $I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top$   
 795 is an orthogonal projection (idempotent and symmetric matrix) of rank  $d - r$ .

796 Now, to find the minimizer of  $\max_{f \in \mathcal{F}_S} \text{regret}(R, f)$  over  $R$ , we note that

$$\lambda_1 \left[ \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) \right] \\ \stackrel{(a)}{=} \max_{u: \|u\|_2=1} u^\top \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top \right) u \quad (64)$$

$$\stackrel{(b)}{=} \max_{u: \|u\|_2=1, \tilde{R}^\top u=0} u^\top \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} u \quad (65)$$

$$\stackrel{(c)}{\geq} \min_{\mathcal{S}: \dim(\mathcal{S})=d-r} \max_{u: \|u\|_2=1, u \in \mathcal{S}} u^\top \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} u \quad (66)$$

$$\stackrel{(d)}{=} \lambda_{r+1} \left( \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} \right), \quad (67)$$

797 where (a) follows again from the *Rayleigh quotient theorem* [37, Theorem 4.2.2], (b) follows since  
 798  $I_d - \tilde{R} (\tilde{R}^\top \tilde{R})^{-1} \tilde{R}^\top$  is an orthogonal projection matrix, and so we may write  $u = u_\perp + u_\parallel$  so  
 799 that  $\|u_\perp\|^2 + \|u_\parallel\|^2 = 1$  and  $\tilde{R}^\top u_\perp = 0$ ; Hence replacing  $u$  with  $u_\perp$  only increases the value of  
 800 the maximum, (c) follows by setting  $\mathcal{S}$  to be a  $d - r$  dimensional subspace of  $\mathbb{R}^d$ , and (d) follows  
 801 by the *Courant–Fischer variational characterization* [37, Theorem 4.2.6] (see Appendix D). The  
 802 lower bound in (c) can be achieved by setting the  $r$  columns of  $\tilde{R} \in \mathbb{R}^{d \times r}$  to be the top eigenvectors  
 803  $\{v_i(\Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2})\}_{i \in [r]}$ . This leads to the minimax representation  $\tilde{R}^*$ . From (63), the worst case  $\tilde{f}$  is  
 804 the top eigenvector of

$$\left( I_d - \tilde{R}^* ((\tilde{R}^*)^\top \tilde{R}^*)^{-1} (\tilde{R}^*)^\top \right) \Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \tilde{R}^* ((\tilde{R}^*)^\top \tilde{R}^*)^{-1} (\tilde{R}^*)^\top \right). \quad (68)$$

805 This is a symmetric matrix, whose top eigenvector is the  $(r + 1)$ th eigenvector  $v_{r+1}(\Sigma_{\mathbf{x}}^{1/2} S \Sigma_{\mathbf{x}}^{1/2})$ .  
 806  $\square$

807 We next prove Theorem 3.

808 *Proof of Theorem 3.* We follow the proof strategy mentioned after the statement of the theorem.  
 809 We assume that  $\mathbf{n} \equiv 0$  with probability 1, since, as for the pure minimax regret, this unavoidable  
 810 additive term of  $\mathbb{E}[\mathbb{E}[\mathbf{n}^2 \mid \mathbf{x}]]$  to the loss does not affect the regret.

811 The minimax problem – a direct computation: As in the derivations leading to (63), the minimax  
 812 regret in (2) is given by

$$\text{regret}_{\text{mix}}(\mathcal{F}_S \mid \Sigma_{\mathbf{x}}) \\ = \min_{\mathbf{R} \in \mathcal{P}(\mathcal{R})} \max_{f \in \mathcal{F}_S} \mathbb{E}[\text{regret}(\mathbf{R}, f \mid \Sigma_{\mathbf{x}})] \quad (69)$$

$$= \min_{\mathbf{L}(\mathbf{R}) \in \mathcal{P}(\mathcal{R})} \max_{\tilde{\mathbf{f}}: \|\tilde{\mathbf{f}}\|_2^2 \leq 1} \mathbb{E} \left[ \tilde{\mathbf{f}}^\top \left( S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} - S^{1/2} \Sigma_{\mathbf{x}} \mathbf{R} (\mathbf{R}^\top \Sigma_{\mathbf{x}} \mathbf{R})^{-1} \mathbf{R}^\top \Sigma_{\mathbf{x}} S^{1/2} \right) \tilde{\mathbf{f}} \right] \quad (70)$$

$$= \min_{\mathbf{L}(\Sigma_{\mathbf{x}}^{-1/2} \tilde{\mathbf{R}}) \in \mathcal{P}(\mathcal{R})} \max_{\tilde{\mathbf{f}}: \|\tilde{\mathbf{f}}\|_2^2 \leq 1} \tilde{\mathbf{f}}^\top S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \mathbb{E} \left[ I_d - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top \right] \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\mathbf{f}} \quad (71)$$

$$= \min_{\mathbf{L}(\Sigma_{\mathbf{x}}^{-1/2} \tilde{\mathbf{R}}) \in \mathcal{P}(\mathcal{R})} \lambda_1 \left( S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \mathbb{E} \left[ I_d - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top \right] \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \right) \quad (72)$$

813 where  $\tilde{\mathbf{R}} = \Sigma_{\mathbf{x}}^{1/2} \mathbf{R}$ . Determining the optimal distribution of the representation directly from this  
814 expression seems to be intractable. We thus next solve the maximin problem, and then return to the  
815 maximin problem (72), set a specific random representation, and show that it achieves the maximin  
816 value. This, in turn, establishes the optimality of this choice.

817 *The maximin problem:* Let an arbitrary  $\mathbf{L}(\mathbf{f})$  be given. Then, taking the expectation of the regret  
818 over the random choice of  $\mathbf{f}$ , for any given  $R \in \mathcal{R}$ ,

$$\mathbb{E} [\text{regret}(R, \mathbf{f})] \stackrel{(a)}{=} \mathbb{E} \left[ \text{Tr} \left[ \left( S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} - S^{1/2} \Sigma_{\mathbf{x}} R (\mathbf{R}^\top \Sigma_{\mathbf{x}} R)^{-1} \mathbf{R}^\top \Sigma_{\mathbf{x}} S^{1/2} \right) \tilde{\mathbf{f}} \tilde{\mathbf{f}}^\top \right] \right] \quad (73)$$

$$\stackrel{(b)}{=} \text{Tr} \left[ \left( S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} - S^{1/2} \Sigma_{\mathbf{x}} R (\mathbf{R}^\top \Sigma_{\mathbf{x}} R)^{-1} \mathbf{R}^\top \Sigma_{\mathbf{x}} S^{1/2} \right) \tilde{\Sigma}_{\mathbf{f}} \right] \quad (74)$$

$$= \text{Tr} \left[ \tilde{\Sigma}_{\mathbf{f}}^{1/2} \left( S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} - S^{1/2} \Sigma_{\mathbf{x}} R (\mathbf{R}^\top \Sigma_{\mathbf{x}} R)^{-1} \mathbf{R}^\top \Sigma_{\mathbf{x}} S^{1/2} \right) \tilde{\Sigma}_{\mathbf{f}}^{1/2} \right] \quad (75)$$

$$\stackrel{(c)}{=} \text{Tr} \left[ \tilde{\Sigma}_{\mathbf{f}}^{1/2} S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \left( I_d - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\Sigma}_{\mathbf{f}}^{1/2} \right] \quad (76)$$

$$= \text{Tr} \left[ \left( I - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \right] \quad (77)$$

$$\stackrel{(d)}{=} \text{Tr} \left[ \left( I - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top \right) \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}}^{1/2} \left( I - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top \right) \right] \quad (78)$$

$$\stackrel{(e)}{\geq} \min_{W \in \mathbb{R}^{d \times (d-r)}: W^\top W = I_{d-r}} \text{Tr} \left[ W^\top \Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}}^{1/2} W \right] \quad (79)$$

$$\stackrel{(f)}{=} \sum_{i=r+1}^d \lambda_i (\Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}}^{1/2}) \quad (80)$$

$$= \sum_{i=r+1}^d \lambda_i \left( \tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}} S^{1/2} \right), \quad (81)$$

819 where (a) follows from Lemma 15 and setting  $\tilde{\mathbf{f}} := S^{-1/2} \mathbf{f}$ , (b) follows by setting  $\tilde{\Sigma}_{\mathbf{f}} \equiv \Sigma_{\tilde{\mathbf{f}}} =$   
820  $\mathbb{E}[\tilde{\mathbf{f}} \tilde{\mathbf{f}}^\top]$ , (c) follows by setting  $\tilde{\mathbf{R}} := \Sigma_{\mathbf{x}}^{1/2} \mathbf{R}$ , (d) follows since  $I - \tilde{\mathbf{R}} (\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1} \tilde{\mathbf{R}}^\top$  is an orthogo-  
821 nal projection (idempotent and symmetric matrix) of rank  $d - r$ , (e) follows since any orthogonal  
822 projection can be written as  $WW^\top$  where  $W \in \mathbb{R}^{d \times (d-r)}$  is an orthogonal matrix  $W^\top W = I_{d-r}$ ,  
823 (f) follows from *Fan's variational characterization* [108] [37, Corollary 4.3.39.] (see Appendix D).  
824 Equality in (e) can be achieved by letting  $\tilde{\mathbf{R}}$  be the top  $r$  eigenvectors of  $\Sigma_{\mathbf{x}}^{1/2} S^{1/2} \tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}}^{1/2}$ .

825 The next step of the derivation is to maximize the expected regret over the probability law of  $\mathbf{f}$   
826 (or  $\tilde{\mathbf{f}}$ ). Evidently,  $\mathbb{E}[\text{regret}(R, \mathbf{f})] = \sum_{i=r+1}^d \lambda_i (\tilde{\Sigma}_{\mathbf{f}} S^{1/2} \Sigma_{\mathbf{x}} S^{1/2})$  only depends on the random  
827 function  $\tilde{\mathbf{f}}$  via  $\tilde{\Sigma}_{\mathbf{f}}$ . The covariance matrix  $\tilde{\Sigma}_{\mathbf{f}}$  is constrained as follows. Recall that  $\mathbf{f}$  is supported  
828 on  $\mathcal{F}_S := \{f \in \mathbb{R}^d: \|f\|_S^2 \leq 1\}$  (see (4)), and let  $\Sigma_{\mathbf{f}} = \mathbb{E}[\mathbf{f} \mathbf{f}^\top]$  be its covariance matrix. Then, it  
829 must hold that  $\text{Tr}[S^{-1} \Sigma_{\mathbf{f}}] \leq 1$ . Then, it also holds that

$$1 \geq \text{Tr}[S^{-1} \Sigma_{\mathbf{f}}] = \text{Tr} \left[ \mathbb{E}[S^{-1} \mathbf{f} \mathbf{f}^\top] \right] \quad (82)$$

$$= \mathbb{E} \left[ \mathbf{f}^\top S^{-1} \mathbf{f} \right] \quad (83)$$

$$= \mathbb{E} \left[ \tilde{\mathbf{f}}^\top \tilde{\mathbf{f}} \right] \quad (84)$$

$$= \text{Tr} \left[ \tilde{\Sigma}_{\mathbf{f}} \right] \quad (85)$$

830 where  $\tilde{\Sigma}_{\mathbf{f}} = S^{-1/2}\Sigma_{\mathbf{f}}S^{-1/2}$ . Conversely, given any covariance matrix  $\tilde{\Sigma}_{\mathbf{f}} \in \mathbb{S}_{++}^d$  such that  
 831  $\text{Tr}[\tilde{\Sigma}_{\mathbf{f}}] \leq 1$  there exists a random vector  $\mathbf{f}$  supported on  $\mathcal{F}_S$  such that

$$\mathbb{E}[\mathbf{f}\mathbf{f}^\top] = S^{1/2}\tilde{\Sigma}_{\mathbf{f}}S^{-1/2}. \quad (86)$$

832 We show this by an explicit construction. Let  $\tilde{\Sigma}_{\mathbf{f}} = \tilde{V}_{\mathbf{f}}\tilde{\Lambda}_{\mathbf{f}}\tilde{V}_{\mathbf{f}}^\top$  be the eigenvalue decomposition of  
 833  $\tilde{\Sigma}_{\mathbf{f}}$ , and, for brevity, denote by  $\tilde{\lambda}_i \equiv \lambda_i(\tilde{\Sigma}_{\mathbf{f}})$  the diagonal elements of  $\tilde{\Lambda}_{\mathbf{f}}$ . Let  $\{\mathbf{q}_i\}_{i \in [d]}$  be a set  
 834 of independent and identically (IID) distributed random variables, so that  $\mathbf{q}_i$  is Rademacher, that is  
 835  $\mathbb{P}[\mathbf{q}_i = 1] = \mathbb{P}[\mathbf{q}_i = -1] = 1/2$ . Define the random vector

$$\mathbf{g} := \left( \mathbf{q}_1 \cdot \sqrt{\tilde{\lambda}_1}, \dots, \mathbf{q}_d \cdot \sqrt{\tilde{\lambda}_d} \right)^\top. \quad (87)$$

836 The constraint  $\text{Tr}[\tilde{\Sigma}_{\mathbf{f}}] \leq 1$  implies that  $\sum \tilde{\lambda}_i \leq 1$  and so  $\|\mathbf{g}\|^2 = \sum_{i=1}^d \tilde{\lambda}_i \leq 1$  with probability 1.  
 837 Then, letting  $\tilde{\mathbf{f}} = \tilde{V}_{\mathbf{f}}\mathbf{g}$  it also holds that  $\|\tilde{\mathbf{f}}\|_2^2 = \|\mathbf{g}\|_2^2 \leq 1$  with probability 1, and furthermore,

$$\mathbb{E}[\tilde{\mathbf{f}}\tilde{\mathbf{f}}^\top] = \tilde{V}_{\mathbf{f}}\mathbb{E}[\mathbf{g}\mathbf{g}^\top]\tilde{V}_{\mathbf{f}}^\top = \tilde{V}_{\mathbf{f}}\tilde{\Lambda}_{\mathbf{f}}\tilde{V}_{\mathbf{f}}^\top = \tilde{\Sigma}_{\mathbf{f}}. \quad (88)$$

838 Consequently, letting  $\mathbf{f} = S^{1/2}\tilde{\mathbf{f}}$  assures that  $\|\mathbf{f}\|_S = \|\tilde{\mathbf{f}}\|_2 \leq 1$  and  $\mathbb{E}[\mathbf{f}\mathbf{f}^\top] = S^{1/2}\tilde{\Sigma}_{\mathbf{f}}S^{-1/2}$ , as  
 839 was required to obtain. Therefore, instead of maximizing over probability laws on  $\mathcal{P}(\mathcal{F}_S)$ , we may  
 840 equivalently maximize over  $\tilde{\Sigma}_{\mathbf{f}} \in \mathbb{S}_{++}^d$  such that  $\text{Tr}[\tilde{\Sigma}_{\mathbf{f}}] \leq 1$ , i.e., to solve

$$\text{regret}_{\text{mix}}(\mathcal{F}_S \mid \Sigma_{\mathbf{x}}) = \max_{\tilde{\Sigma}_{\mathbf{f}}: \text{Tr}[\tilde{\Sigma}_{\mathbf{f}}] \leq 1} \sum_{i=r+1}^d \lambda_i(\tilde{\Sigma}_{\mathbf{f}}S^{1/2}\Sigma_{\mathbf{x}}S^{1/2}). \quad (89)$$

841 The optimization problem in (89) is solved in Lemma 16, and is provided after this proof. Setting  
 842  $\Sigma = S^{1/2}\Sigma_{\mathbf{x}}S^{1/2}$  in Lemma 16, and letting  $\lambda_i \equiv \lambda_i(S^{1/2}\Sigma_{\mathbf{x}}S^{1/2})$ , the solution is given by

$$\frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \frac{1}{\lambda_i}} \quad (90)$$

843 where  $\ell^* \in [d] \setminus [r]$  satisfies

$$\frac{\ell^* - r}{\lambda_{\ell^*}} \leq \sum_{i=1}^{\ell^*} \frac{1}{\lambda_i} \leq \frac{\ell^* - r}{\lambda_{\ell^*+1}}. \quad (91)$$

844 Lemma 16 also directly implies that an optimal  $\tilde{\Sigma}_{\mathbf{f}}$  is given as in (10). The value in (90) is exactly  
 845  $\text{regret}_{\text{mix}}(\mathcal{F}_S \mid \Sigma_{\mathbf{x}})$  claimed by the theorem, and we next show it is indeed achievable by a properly  
 846 constructed random representation.

847 *The minimax problem – a solution via the maximin certificate:* Given the value of the regret game in  
 848 mixed strategies found in (90), we may also find a minimax representation in mixed strategies. To  
 849 this end, we return to the minimax expression in (72), and propose a random representation which  
 850 achieves the maximin value in (90). Note that for any given  $\tilde{R}$ , the matrix  $I_d - \tilde{R}(\tilde{R}^\top \tilde{R})^{-1}\tilde{R}^\top$  is an  
 851 orthogonal projection, that is, a symmetric matrix whose eigenvalues are all either 0 or 1, and it has at  
 852 most  $r$  eigenvalues equal to zero. We denote its eigenvalue decomposition by  $I_d - \tilde{R}(\tilde{R}^\top \tilde{R})^{-1}\tilde{R}^\top =$   
 853  $U\Omega U^\top$ . Then, any probability law on  $\tilde{R}$  induces a probability law on  $U$  and  $\Omega$  (and vice-versa).  
 854 To find the mixed minimax representation, we propose setting  $U = V(\Sigma_{\mathbf{x}}^{1/2}S\Sigma_{\mathbf{x}}^{1/2}) \equiv V$  with  
 855 probability 1, that is, to be deterministic, and thus only randomize  $\Omega$ . With this choice, and by  
 856 denoting, for brevity,  $\Lambda \equiv \Lambda(\Sigma_{\mathbf{x}}^{1/2}S\Sigma_{\mathbf{x}}^{1/2}) = \Lambda(S^{1/2}\Sigma_{\mathbf{x}}S^{1/2})$ , the value of the objective function in  
 857 (72) is given by

$$\begin{aligned} & \lambda_1 \left( S^{1/2}\Sigma_{\mathbf{x}}^{1/2}V \cdot \mathbb{E}[\Omega] \cdot V^\top \Sigma_{\mathbf{x}}^{1/2}S^{1/2} \right) \\ &= \lambda_1 \left( \mathbb{E}[\Omega] \cdot V^\top \Sigma_{\mathbf{x}}^{1/2}S\Sigma_{\mathbf{x}}^{1/2}V \right) \\ &= \lambda_1 (\mathbb{E}[\Omega] \cdot \Lambda). \end{aligned} \quad (92)$$

$$= \lambda_1 (\mathbb{E}[\Omega] \cdot \Lambda). \quad (93)$$

858 Now, the distribution of  $\Omega$  is equivalent to a distribution on its diagonal, which is supported on  
 859 the finite set  $\mathcal{A} := \{a \in \{0, 1\}^d : \|a\|_1 \geq d - r\}$ . Our goal is thus to find a probability law on  $\mathbf{a}$ ,  
 860 supported on  $\mathcal{A}$ , which solves

$$\min_{L(\Omega)} \max_{i \in [d]} \lambda_i (\mathbb{E}[\Omega] \cdot \Lambda) = \min_{L(\mathbf{a})} \max_{i \in [d]} \mathbb{E}[\mathbf{a}_i] \lambda_i \quad (94)$$

861 where  $\lambda_i \equiv \lambda_i(S^{1/2} \Sigma_{\mathbf{x}} S^{1/2})$  are the diagonal elements of  $\Lambda$ . Consider  $\ell^*$ , the optimal dimension  
 862 of the maximin problem, which satisfies (91). We then set  $\mathbf{a}_{\ell^*+1} = \dots = \mathbf{a}_d = 1$  to hold with  
 863 probability 1, and so it remains to determine the probability law of  $\bar{\mathbf{a}} := (\mathbf{a}_1, \dots, \mathbf{a}_{\ell^*})$ , supported  
 864 on  $\bar{\mathcal{A}} := \{a \in \{0, 1\}^{\ell^*} : \|a\|_1 \geq \ell^* - r\}$ . Clearly, reducing  $\|a\|_1$  only reduces the objective function  
 865  $\max_{i \in [d]} \mathbb{E}[\mathbf{a}_i] \lambda_i$ , and so we may in fact assume that  $\bar{\mathbf{a}}$  is supported on  $\bar{\mathcal{A}} := \{\bar{a} \in \{0, 1\}^{\ell^*} : \|\bar{a}\|_1 =$   
 866  $\ell^* - r\}$ , a finite subset of cardinality  $\binom{\ell^*}{r}$ . Suppose that we find a probability law  $L(\bar{\mathbf{a}})$  supported  
 867 on  $\bar{\mathcal{A}}$  such that

$$\mathbb{E}[\mathbf{a}_i] = (\ell^* - r) \cdot \frac{1/\lambda_i}{\sum_{i=1}^{\ell^*} 1/\lambda_i} := b_i, \quad (95)$$

868 for all  $i \in [\ell^*]$ . Then, since  $\mathbb{E}[\mathbf{a}_i] = 1$  for  $i \in [d] \setminus [\ell^*]$

$$\max_{i \in [d]} \mathbb{E}[\mathbf{a}_i] \lambda_i = \max \left\{ \frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \frac{1}{\lambda_i}}, \lambda_{\ell^*+1}, \dots, \lambda_d \right\} \quad (96)$$

$$= \max \left\{ \frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \frac{1}{\lambda_i}}, \lambda_{\ell^*+1} \right\} \quad (97)$$

$$\stackrel{(*)}{=} \frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \frac{1}{\lambda_i}}, \quad (98)$$

869 where  $(*)$  follows from the condition on  $\ell^*$  in the right inequality of (91). This proves that such  
 870 probability law achieves the minimax regret in mixed strategies. This last term is  $\text{regret}_{\text{mix}}(\mathcal{F}_S \mid \Sigma_{\mathbf{x}})$   
 871 claimed by the theorem. It remains to construct  $L(\bar{\mathbf{a}})$  which satisfies (95). To this end, note that the  
 872 set

$$\mathcal{C} := \left\{ c \in [0, 1]^{\ell^*} : \|c\|_1 = \ell^* - r \right\} \quad (99)$$

873 is convex and compact, and  $\bar{\mathcal{A}}$  is the set of its *extreme points* ( $\mathcal{C}$  is the convex hull of  $\bar{\mathcal{A}}$ ). Letting  
 874  $\bar{\mathbf{b}} = (b_1, \dots, b_{\ell^*})^\top$  as denoted in (95), it holds that  $\bar{b}_i \geq 0$  and  $\{\bar{b}_i\}_{i=1}^{\ell^*}$  is a non-decreasing sequence.  
 875 Using the condition on  $\ell^*$  in the left inequality of (91), it then holds that

$$\bar{b}_1 \leq \dots \leq \bar{b}_{\ell^*} = (\ell^* - r) \cdot \frac{1/\lambda_{\ell^*}}{\sum_{i=1}^{\ell^*} 1/\lambda_i} \leq 1. \quad (100)$$

876 Hence,  $\bar{\mathbf{b}} \in \mathcal{C}$ . By Carathéodory's theorem [107, Prop. 1.3.1] (see Appendix D), any point inside a  
 877 convex compact set in  $\mathbb{R}^{\ell^*}$  can be written as a convex combination of at most  $\ell^* + 1$  extreme points.  
 878 Thus, there exists  $\{p_{\bar{a}}\}_{\bar{a} \in \bar{\mathcal{A}}}$  such that  $p_{\bar{a}} \in [0, 1]$  and  $\sum_{\bar{a} \in \bar{\mathcal{A}}} p_{\bar{a}} = 1$  so that  $\bar{\mathbf{b}} = \sum_{\bar{a} \in \bar{\mathcal{A}}} p_{\bar{a}} \cdot \bar{\mathbf{a}}$ , and  
 879 moreover the support of  $p_{\bar{a}}$  has cardinality at most  $\ell^* + 1$ . Let  $\bar{\mathcal{A}} \in \{0, 1\}^{\ell^* \times |\bar{\mathcal{A}}|}$  be such that its  $j$ th  
 880 column is given by the  $j$ th member of  $\bar{\mathcal{A}}$  (in an arbitrary order). Let  $p \in [0, 1]^{|\bar{\mathcal{A}}|}$  be a vector whose  
 881  $j$ th element corresponds to the  $j$ th member of  $\bar{\mathcal{A}}$ . Then,  $p$  is the solution to  $\bar{\mathcal{A}} p = \bar{\mathbf{b}}$ , and as claimed  
 882 above, such a solution with at most  $\ell^* + 1$  nonzero entries always exists. Setting  $\mathbf{a} = (\bar{\mathbf{a}}, \underbrace{1, \dots, 1}_{d - \ell^* \text{ terms}})$

883 with probability  $p_{\bar{a}}$  then assures that (95) holds, as was required to be proved.

884 Given the above, we observe that setting  $\tilde{R}$  as in the theorem induces a distribution on  $\Omega$  for which  
 885 the random entries of its diagonal  $\mathbf{a}$  satisfy (95), and thus achieve  $\text{regret}_{\text{mix}}(\mathcal{F}_S \mid \Sigma_{\mathbf{x}})$ .  $\square$

886 We next turn to complete the proof of Theorem 3 by solving the optimization problem in (89).  
 887 Assume that  $\Sigma \in \mathbb{S}_{++}^d$  is a strictly positive covariance matrix  $\Sigma \succ 0$ , and consider the optimization  
 888 problem

$$v_r^* = \max_{\tilde{\Sigma}_{\mathcal{F}} \in \mathbb{S}_{++}^d} \sum_{i=r+1}^d \lambda_i(\tilde{\Sigma}_{\mathcal{F}} \Sigma)$$

$$\text{subject to } \text{Tr}[\tilde{\Sigma}_{\mathbf{f}}] \leq 1 \quad (101)$$

889 for some  $r \in [d-1]$ . Note that the objective function refers to the maximization of the  $d-r$  minimal  
890 eigenvalues of  $\Sigma^{1/2}\tilde{\Sigma}_{\mathbf{f}}\Sigma^{1/2}$ .

891 **Lemma 16.** *Let*

$$a_{\ell} := \frac{\ell - r}{\sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)}}. \quad (102)$$

892 *The optimal value of (101) is  $v^* = \max_{[d]\setminus[r]} a_{\ell}$  and  $\ell^* \in \arg \max_{[d]\setminus[r]} a_{\ell}$  iff*

$$\frac{\ell^* - r}{\lambda_{\ell^*}(\Sigma)} \leq \sum_{i=1}^{\ell^*} \frac{1}{\lambda_i(\Sigma)} \leq \frac{\ell^* - r}{\lambda_{\ell^*+1}(\Sigma)}. \quad (103)$$

893 *An optimal solution is*

$$\tilde{\Sigma}_{\mathbf{f}}^* = \left[ \sum_{i=1}^{\ell^*} \frac{1}{\lambda_i(\Sigma)} \right]^{-1} \cdot V(\Sigma) \text{diag} \left( \frac{1}{\lambda_1(\Sigma)}, \dots, \frac{1}{\lambda_{\ell^*}(\Sigma)}, 0, \dots, 0 \right) \cdot V(\Sigma)^{\top}. \quad (104)$$

894 *Proof.* Let  $\bar{\Sigma}_{\mathbf{f}} = \Sigma^{1/2}\tilde{\Sigma}_{\mathbf{f}}\Sigma^{1/2}$ , let  $\bar{\Sigma}_{\mathbf{f}} = \bar{U}_{\mathbf{f}}\bar{\Lambda}_{\mathbf{f}}\bar{U}_{\mathbf{f}}^{\top}$  be its eigenvalue decomposition, and, for  
895 brevity, denote  $\bar{\lambda}_i \equiv \lambda_i(\bar{\Sigma}_{\mathbf{f}})$ . Then, the trace operation appearing in the constraint of (101) can be  
896 written as

$$\text{Tr}[\tilde{\Sigma}_{\mathbf{f}}] = \text{Tr} \left[ \Sigma^{-1/2}\bar{\Sigma}_{\mathbf{f}}\Sigma^{-1/2} \right] \quad (105)$$

$$= \text{Tr} \left[ \Sigma^{-1/2}\bar{U}_{\mathbf{f}}\bar{\Lambda}_{\mathbf{f}}\bar{U}_{\mathbf{f}}^{\top}\Sigma^{-1/2} \right] \quad (106)$$

$$= \text{Tr} \left[ \Sigma^{-1/2} \left( \sum_{i=1}^d \lambda_i \bar{u}_i \bar{u}_i^{\top} \right) \Sigma^{-1/2} \right] \quad (107)$$

$$= \sum_{i=1}^d \bar{\lambda}_i \cdot (\bar{u}_i^{\top} \Sigma^{-1} \bar{u}_i) \quad (108)$$

$$= \sum_{i=1}^d c_i \bar{\lambda}_i, \quad (109)$$

897 where  $\bar{u}_i = v_i(\bar{U}_{\mathbf{f}})$  (that is, the  $i$ th column of  $\bar{U}_{\mathbf{f}}$ ), and  $c_i := \bar{u}_i^{\top} \Sigma^{-1} \bar{u}_i$  (which satisfies  $c_i > 0$ ).  
898 Thus, the optimization problem in (101) over  $\tilde{\Sigma}_{\mathbf{f}}$  is equivalent to an optimization problem over  
899  $\{\bar{\lambda}_i, \bar{u}_i\}_{i \in [d]}$ , given by

$$\begin{aligned} v_r^* &= \max_{\{\bar{u}_i, \bar{\lambda}_i\}_{i \in [d]}} \sum_{i=r+1}^d \bar{\lambda}_i \\ &\text{subject to } \sum_{i=1}^d c_i \bar{\lambda}_i \leq 1, \\ &\quad c_i = \bar{u}_i^{\top} \Sigma^{-1} \bar{u}_i, \\ &\quad \bar{u}_i^{\top} \bar{u}_j = \delta_{ij}, \\ &\quad \bar{\lambda}_1 \geq \bar{\lambda}_2 \geq \dots \geq \bar{\lambda}_d \geq 0. \end{aligned} \quad (110)$$

900 To solve the optimization problem (110), let us fix feasible  $\{\bar{u}_i\}_{i \in [d]}$ , so that  $\{c_i\}_{i \in [d]}$  are fixed too.  
901 This results the problem

$$\begin{aligned} v_r^*(\{\bar{u}_i\}) &\equiv v_r^*(\{c_i\}) = \max_{\{\bar{\lambda}_i\}_{i \in [d]}} \sum_{i=r+1}^d \bar{\lambda}_i \\ &\text{subject to } \sum_{i=1}^d c_i \bar{\lambda}_i \leq 1, \end{aligned}$$

$$\bar{\lambda}_1 \geq \bar{\lambda}_2 \geq \dots \geq \bar{\lambda}_d \geq 0. \quad (111)$$

902 The objective function of (111) is linear in  $\{\bar{\lambda}_i\}_{i \in [d]}$  and its constraint set is a convex bounded  
 903 polytope. So the solution to (111) must be obtained on the boundary of the constraint set. Clearly,  
 904 the optimal value satisfies  $v_r^*({c_i}) \geq 0$ , and thus the solution  $\{\bar{\lambda}_i^*\}_{i \in [d]}$  must be obtained when the  
 905 constraint  $\sum_{i=1}^d c_i \bar{\lambda}_i \leq 1$  is satisfied with equality. Indeed, if this is not the case then one may scale  
 906 all  $\bar{\lambda}_i^*$  by a constant larger than 1, and obtain larger value of the objective, while still satisfying the  
 907 constraint.

908 To find the optimal solution to (111), we consider feasible points for which  $\ell := \max\{i \in [d]: \bar{\lambda}_i >$   
 909  $0\}$  is fixed. Let  $\{\bar{\lambda}_i^*\}_{i \in [d]}$  be the optimal solution of (111), under the additional constraint that  
 910  $\bar{\lambda}_{\ell+1} = \dots = \bar{\lambda}_d = 0$ . We next prove that  $\bar{\lambda}_1^* = \dots = \bar{\lambda}_\ell^*$  must hold. To this end, assume by  
 911 contradiction that there exists  $j \in [\ell]$  so that  $\bar{\lambda}_{j-1}^* > \bar{\lambda}_j^* > 0$ . There are two cases to consider,  
 912 to wit, whether  $j - 1 < r + 1$  and so only  $\bar{\lambda}_j$  appears in the objective of (111), or, otherwise,  
 913  $j - 1 \geq r + 1$  and then  $\bar{\lambda}_{j-1} + \bar{\lambda}_j$  appears in the objective of (111). Assuming the first case, let  
 914  $\alpha = \bar{\lambda}_{j-1}^* c_{j-1} + \bar{\lambda}_j^* c_j$  and consider the optimization problem

$$\begin{aligned} & \max_{\hat{\lambda}_{j-1}, \hat{\lambda}_j} \hat{\lambda}_j \\ & \text{subject to } \hat{\lambda}_{j-1} c_{j-1} + \hat{\lambda}_j c_j = \alpha, \\ & \hat{\lambda}_{j-1} \geq \hat{\lambda}_j > 0. \end{aligned} \quad (112)$$

915 It is easy to verify that the optimum of this problem is  $\hat{\lambda}_{j-1}^* = \hat{\lambda}_j^* = \frac{\alpha}{c_{j-1} + c_j}$ . Thus, if  $\bar{\lambda}_{j-1}^* > \bar{\lambda}_j^*$   
 916 then one can replace this pair with  $\bar{\lambda}_{j-1}^* = \bar{\lambda}_j^* = \hat{\lambda}_{j-1}^* = \hat{\lambda}_j^*$  so that the value of the constraint  
 917  $\sum_{i=1}^d \bar{\lambda}_i c_i$  remains the same, and thus  $(\bar{\lambda}_1^*, \dots, \hat{\lambda}_{j-1}^*, \hat{\lambda}_j^*, \bar{\lambda}_{j+1}^*, \dots, \bar{\lambda}_d^*)$  is a feasible point, while the  
 918 objective function value of (111) is smaller; a contradiction. Therefore, it must hold for the first  
 919 case that  $\bar{\lambda}_{j-1}^* = \bar{\lambda}_j^*$ . For the second case, in a similar fashion, let now  $\alpha = \bar{\lambda}_{j-1}^* c_{j-1} + \bar{\lambda}_j^* c_j$ , and  
 920 consider the optimization problem

$$\begin{aligned} & \max_{\hat{\lambda}_{j-1}, \hat{\lambda}_j} \hat{\lambda}_j + \hat{\lambda}_{j-1} \\ & \text{subject to } \hat{\lambda}_{j-1} c_{j-1} + \hat{\lambda}_j c_j = \alpha, \\ & \hat{\lambda}_{j-1} \geq \hat{\lambda}_j > 0. \end{aligned} \quad (113)$$

921 The solution for this optimization problem is at one of the two extreme points of the feasible interval  
 922 for  $\hat{\lambda}_j$ . Since  $\lambda_j^* > 0$  was assumed it therefore must hold that  $\hat{\lambda}_{j-1}^* = \hat{\lambda}_j^*$ , and hence also  $\bar{\lambda}_{j-1}^* = \bar{\lambda}_j^*$ .  
 923 Thus,  $\lambda_{j-1}^* < \lambda_j^*$  leads to a contradiction. From the above, we deduce that the optimal solution of  
 924 (111) under the additional constraint that  $\bar{\lambda}_{\ell+1} = \dots = \bar{\lambda}_d = 0$  is

$$\bar{\lambda}_1^* = \dots = \bar{\lambda}_\ell^* = \frac{1}{\sum_{i=1}^\ell c_i} \quad (114)$$

$$\bar{\lambda}_{\ell+1}^* = \dots = \bar{\lambda}_d^* = 0, \quad (115)$$

925 and that the optimal value is  $\frac{\ell-r}{\sum_{i=1}^{\ell-r} c_i}$ . Since  $\ell \in [d] \setminus [r]$  can be arbitrarily chosen, we deduce that the  
 926 value of (111) is

$$v^*({c_i}) = \max_{\ell \in [d] \setminus [r]} \frac{\ell - r}{\sum_{i=1}^\ell c_i}. \quad (116)$$

927 For any given  $\ell \in [d] \setminus [r]$ , we may now optimize over  $\{\bar{u}_i\}$ , which from (116) is equivalent to  
 928 minimizing  $\sum_{i=1}^\ell c_i$ . It holds that

$$\min_{\{\bar{u}_i\}} \sum_{i=1}^\ell c_i = \min_{\{\bar{u}_i: \bar{u}_i^\top \bar{u}_j = \delta_{ij}\}} \sum_{i=1}^\ell \bar{u}_i^\top \Sigma^{-1} \bar{u}_i \quad (117)$$

$$= \min_{\{\bar{u}_i: \bar{u}_i^\top \bar{u}_j = \delta_{ij}\}} \text{Tr} \left[ \Sigma^{-1} \sum_{i=1}^\ell \bar{u}_i \bar{u}_i^\top \right] \quad (118)$$

$$\stackrel{(a)}{=} \min_{\dot{U} \in \mathbb{R}^{d \times \ell}; \dot{U}^\top \dot{U} = I_\ell} \text{Tr} \left[ \Sigma^{-1} \dot{U} \dot{U}^\top \right] \quad (119)$$

$$= \min_{\dot{U} \in \mathbb{R}^{d \times \ell}; \dot{U}^\top \dot{U} = I_\ell} \text{Tr} \left[ \dot{U}^\top \Sigma^{-1} \dot{U} \right] \quad (120)$$

$$\stackrel{(b)}{=} \sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)}, \quad (121)$$

929 where in (a)  $\dot{U} \in \mathbb{R}^{d \times \ell}$  whose  $\ell$  columns are  $\{\bar{u}_i\}_{i \in [\ell]}$  and  $\dot{U}^\top \dot{U} = I_\ell$ , and in (b) we have used  
 930 *Fan's variational characterization* [108] [37, Corollary 4.3.39.] (see Appendix D). Substituting  
 931 back to (116) results that

$$v_r^* = \max_{\ell \in [d] \setminus [r]} \frac{\ell - r}{\sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)}} = \max_{\ell \in [d] \setminus [r]} a_\ell. \quad (122)$$

932 Let us denote that maximizer index by  $\ell^*$ . Then, Fan's characterization is achieved by setting  
 933  $\bar{U}_f = V$  (so that the  $\ell^*$  columns of  $\dot{U}$  are the  $\ell^*$  eigenvectors  $v_i(\Sigma)$ , corresponding to the  $\ell^*$  largest  
 934 eigenvalues of  $\Sigma$ ), so that

$$\bar{\Sigma}_f^* = \left[ \sum_{i=1}^{\ell^*} \frac{1}{\lambda_i(\Sigma)} \right]^{-1} \cdot V \cdot \text{diag} \left( \underbrace{1, \dots, 1}_{\ell^* \text{ terms}}, 0, \dots, 0 \right) \cdot V^\top, \quad (123)$$

935 and then

$$\hat{\Sigma}_f^* = \Sigma^{-1/2} \bar{\Sigma}_f^* \Sigma^{-1/2} \quad (124)$$

$$= \left[ \sum_{i=1}^{\ell^*} \frac{1}{\lambda_i(\Sigma)} \right]^{-1} \cdot V \Lambda^{-1/2} V^\top V \cdot \text{diag} (1, \dots, 1, 0, \dots, 0) V^\top V \Lambda^{-1/2} V^\top \quad (125)$$

$$= \left[ \sum_{i=1}^{\ell^*} \frac{1}{\lambda_i(\Sigma)} \right]^{-1} \cdot V \cdot \text{diag} \left( \frac{1}{\lambda_1(\Sigma)}, \dots, \frac{1}{\lambda_{\ell^*}(\Sigma)}, 0, \dots, 0 \right) \cdot V^\top \quad (126)$$

936 as claimed in (104).

937 To complete the proof, it remains to characterize  $\ell^*$ , which belongs to the set possible indices max-  
 938 imizing  $\{a_\ell\}_{\ell \in [d] \setminus [r]}$ . Since  $\ell^*$  maximizes  $a_\ell$  it must be a local maximizer, that is, it must hold that  
 939  $a_{\ell^*-1} \leq a_{\ell^*} \geq a_{\ell^*+1}$ . By simple algebra, these conditions are equivalent to those in (103). It  
 940 remains to show that any  $\ell \in [d] \setminus [r]$  which satisfies (103) has the same value, and thus any local  
 941 maxima is a global maxima. We will show this by proving that the sequence  $\{a_\ell\}_{\ell=r}^d$  is *unimodal*,  
 942 as follows. Let  $\Delta_\ell := a_{\ell+1} - a_\ell$  be the discrete derivative of  $\{a_\ell\}_{\ell \in [d]}$ , and consider the sequence  
 943  $\{\Delta_\ell\}_{\ell \in [d] \setminus [r]}$ . We show that as  $\ell$  increases from  $r$  to  $d$ ,  $\{\Delta_\ell\}_{\ell \in [d] \setminus [r]}$  is only changing its sign at  
 944 most once. To this end, we first note that

$$\Delta_\ell = \frac{\ell + 1 - r}{\sum_{i=1}^{\ell+1} \frac{1}{\lambda_i(\Sigma)}} - \frac{\ell - r}{\sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)}} = \frac{\sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)} - (\ell - r) \frac{1}{\lambda_{\ell+1}(\Sigma)}}{\left[ \sum_{i=1}^{\ell+1} \frac{1}{\lambda_i(\Sigma)} \right] \left[ \sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)} \right]}. \quad (127)$$

945 Since the denominator of (127) is strictly positive, it suffices to prove that the sequence comprised  
 946 of the numerator of (127), to wit  $\{\zeta_\ell\}_{\ell \in [d] \setminus [r]}$  with

$$\zeta_\ell := \sum_{i=1}^{\ell} \frac{1}{\lambda_i(\Sigma)} - (\ell - r) \frac{1}{\lambda_{\ell+1}(\Sigma)}, \quad (128)$$

947 is only changing its sign at most once. Indeed, this claim is true because  $\zeta_r = \sum_{i=1}^r \frac{1}{\lambda_i(\Sigma)} > 0$  and  
 948 because  $\{\zeta_\ell\}_{\ell \in [d] \setminus [r]}$  is a monotonic non-increasing sequence,

$$\zeta_\ell - \zeta_{\ell+1} = (\ell - r + 1) \left[ \frac{1}{\lambda_{\ell+2}(\Sigma)} - \frac{1}{\lambda_{\ell+1}(\Sigma)} \right] \geq 0. \quad (129)$$

949 Therefore,  $\{\zeta_\ell\}_{\ell \in [d] \setminus [r]}$  has at most a single sign change (its has a positive value at  $\ell = r$  and is  
 950 monotonically non-increasing with  $\ell$  up to  $\ell = d$ ), and so is  $\{\Delta_\ell\}_{\ell=r}^d$ . The single sign change  
 951 property of the finite difference  $\{\Delta_\ell\}_{\ell=r}^d$  is equivalent to the fact that  $\{a_\ell\}_{\ell=r}^d$  is *unimodal*. Thus,  
 952 any local maximizer of  $a_\ell$  is also a global maximizer.  $\square$

953 **F The Hilbert space MSE setting**

954 In this section, we show that the regret expressions in Section 3 can be easily generalized to an  
 955 infinite dimensional Hilbert space, for responses with noise that is statistically independent of the  
 956 features. We still assume the MSE loss function ( $\mathcal{Y} = \mathbb{R}$ , and  $\text{loss}(y_1, y_2) = (y_1 - y_2)^2$ ), and that  
 957 the predictor is a linear function. However, we allow the the representation and response function to  
 958 be functions in a Hilbert space. As will be evident, the resulting regret is not very different from the  
 959 finite-dimensional case. Formally, this is defined as follows:

960 **Definition 17** (The Hilbert space MSE setting). Assume that  $\mathbf{x} \sim P_{\mathbf{x}}$  is supported on a compact  
 961 subset  $\mathcal{X} \subset \mathbb{R}^d$ , and let  $L_2(P_{\mathbf{x}})$  be the Hilbert space of functions from  $\mathcal{X} \rightarrow \mathbb{R}$  such that  $\mathbb{E}[f^2(\mathbf{x})] =$   
 962  $\int_{\mathcal{X}} f^2(\mathbf{x}) \cdot dP_{\mathbf{x}} < \infty$ , with the inner product,

$$\langle f, g \rangle := \int_{\mathcal{X}} f(\mathbf{x})g(\mathbf{x}) \cdot dP_{\mathbf{x}} \quad (130)$$

963 for  $f, g \in L_2(P_{\mathbf{x}})$ . Let  $\{\phi_j(x)\}_{j=1}^{\infty}$  be an orthonormal basis for  $L_2(P_{\mathbf{x}})$ .

964 A representation is comprised of a set of functions  $\{\psi_i\}_{i \in [r]} \subset L_2(P_{\mathbf{x}})$ ,  $\psi_i: \mathcal{X} \rightarrow \mathbb{R}$ , so that

$$\mathcal{R} := \{R(x) = (\psi_1(x), \dots, \psi_r(x))^{\top} \in \mathbb{R}^r\}. \quad (131)$$

965 Let  $\{\lambda_j\}_{j \in \mathbb{N}}$  be a positive monotonic non-increasing sequence for which  $\lambda_j \downarrow 0$  as  $j \rightarrow \infty$ , and let  
 966  $\mathcal{F}$  be the set of functions from  $\mathcal{X} \rightarrow \mathbb{R}$  such that given  $f \in \mathcal{F}$ , the response is given by

$$\mathbf{y} = f(\mathbf{x}) + \mathbf{n} \in \mathbb{R} \quad (132)$$

967 where

$$f \in \mathcal{F}_{\{\lambda_j\}} := \left\{ f(x) = \sum_{j=1}^{\infty} f_j \phi_j(x) : \{f_j\}_{j \in \mathbb{N}} \in \ell_2(\mathbb{N}), \sum_{j=1}^{\infty} \frac{f_j^2}{\lambda_j} \leq 1 \right\}, \quad (133)$$

968 where  $\mathbf{n} \in \mathbb{R}$  is a homoscedastic noise that is statistically independent of  $\mathbf{x}$  and satisfies  $\mathbb{E}[\mathbf{n}] = 0$ .  
 969 Infinite-dimensional ellipsoids such as  $\mathcal{F}_{\{\lambda_j\}}$  naturally arise in reproducing kernel Hilbert spaces  
 970 (RKHS) [59, Chapter 12] [60, Chapter 16], in which  $\{\lambda_j\}$  is the eigenvalues of the kernel. In this  
 971 case, the set  $\mathcal{F}_{\{\lambda_i\}} = \{f: \|f\|_{\mathcal{H}} \leq 1\}$  where  $\|\cdot\|_{\mathcal{H}}$  is the norm of the RKHS  $\mathcal{H}$ . For example,  $\mathcal{H}$   
 972 could be the first-order Sobolev space of functions with finite first derivative energy.

973 Let the set of predictor functions be the set of linear functions from  $\mathbb{R}^d \rightarrow \mathbb{R}$ , that is

$$\mathcal{Q} := \{Q(z) = q^{\top} z = \sum_{i=1}^r q_i \cdot \psi_i(x), q \in \mathbb{R}^r\}. \quad (134)$$

974 We denote the pure (resp. mixed) minimax regret as  $\text{regret}_{\text{pure}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}})$  (resp.  
 975  $\text{regret}_{\text{mix}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}})$ ). We begin with pure strategies.

976 **Theorem 18.** For the Hilbert space MSE setting (Definition 17)

$$\text{regret}_{\text{pure}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}}) = \lambda_{r+1}. \quad (135)$$

977 A minimax representation is

$$R^*(x) = (\phi_1(x), \dots, \phi_r(x))^{\top}, \quad (136)$$

978 and the worst case response function is  $f^* = \sqrt{\lambda_{r+1}} \cdot \phi_{r+1}$ .

979 We now turn to the minimax representation in mixed strategies.

980 **Theorem 19.** For the Hilbert space MSE setting (Definition 17)

$$\text{regret}_{\text{mix}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}}) = \frac{\ell^* - r}{\sum_{i=1}^{\ell^*} \frac{1}{\lambda_i}}, \quad (137)$$

981 where  $\ell^*$  is defined as (9) of Theorem 3 (with the replacement  $d \rightarrow \mathbb{N}_+$ ). Let  $\{\mathbf{b}_j\}_{j=1}^{\infty}$  be an IID  
 982 sequence of Rademacher random variables,  $\mathbb{P}[\mathbf{b}_i = 1] = \mathbb{P}[\mathbf{b}_i = -1] = 1/2$ . Then, a least  
 983 favorable prior  $\mathbf{f}^*$  is

$$\mathbf{f}_i^* = \begin{cases} \mathbf{b}_i \cdot \frac{1}{\sqrt{\sum_{i=1}^{\ell^*} \frac{1}{\lambda_i}}}, & 1 \leq i \leq \ell^* \\ 0, & i \geq \ell^* + 1 \end{cases}, \quad (138)$$

984 and a law of minimax representation is to choose

$$\mathbf{R}^*(x) = \{\phi_{\mathcal{I}_j}(x)\}_{j=1}^r \quad (139)$$

985 with probability  $p_j$ ,  $j \in [\binom{\ell^*}{r}]$ , defined as in Theorem 3.

986 **Discussion** Despite having countably infinite possible number of representations, the optimal rep-  
 987 resentation only utilizes a *finite* set of orthogonal functions, as determined by the radius of  $\mathcal{F}_{\{s_i\}}$ .  
 988 The proof of Theorems 18 and 19 is obtained by reducing the infinite dimensional problem to a  $d$ -  
 989 dimensional problem via an approximation argument, then showing the the finite dimensional case  
 990 is similar to the problem of Section 3, and then taking limit  $d \uparrow \infty$ .

## 991 F.1 Proofs

992 Let us denote the  $d$ -dimensional slice of  $\mathcal{F}_{\{\lambda_j\}}$  by

$$\mathcal{F}_{\{\lambda_j\}}^{(d)} := \{f(x) \in \mathcal{F}_{\{\lambda_j\}} : f_j = 0 \text{ for all } j \geq d+1\}. \quad (140)$$

993 Further, let us consider the restricted representation class, in which the representation functions  
 994  $\psi_i(t)$  belong to the span of the first  $d$  basis functions, that is

$$\mathcal{R}^{(d)} := \{R(x) \in \mathcal{R} : \psi_i(x) \in \text{span}(\{\phi_i\}_{i \in [d]}) \text{ for all } i \in [r]\}. \quad (141)$$

995 The following proposition implies that the regret in the infinite-dimensional Hilbert space is obtained  
 996 as the limit of finite-dimensional regrets, as the one characterized in Section 3:

997 **Proposition 20.** *It holds that*

$$\text{regret}_{\text{pure}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}}) = \lim_{d \uparrow \infty} \text{regret}_{\text{pure}}(\mathcal{R}^{(d)}, \mathcal{F}_{\{\lambda_j\}}^{(d)} \mid P_{\mathbf{x}}) \quad (142)$$

998 and

$$\text{regret}_{\text{mix}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}}) = \lim_{d \uparrow \infty} \text{regret}_{\text{mix}}(\mathcal{R}^{(d)}, \mathcal{F}_{\{\lambda_j\}}^{(d)} \mid P_{\mathbf{x}}). \quad (143)$$

999 *Proof.* Let  $\{c_{ij}\}_{j \in \mathbb{N}}$  be the coefficients of the orthogonal expansion of  $\psi_i$ ,  $i \in [r]$ , that is,  $\psi_i =$   
 1000  $\sum_{j=1}^{\infty} c_{ij} \phi_j$ . With a slight abuse of notation, we also let  $c_i := (c_{i1}, c_{i2} \dots) \in \ell_2(\mathbb{N})$ . We use a  
 1001 sandwich argument. On one hand,

$$\text{regret}_{\text{pure}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}}) = \min_{R \in \mathcal{R}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}} \text{regret}(R, f) \quad (144)$$

$$\geq \min_{R \in \mathcal{R}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}^{(d)}} \text{regret}(R, f) \quad (145)$$

$$\stackrel{(*)}{=} \min_{R \in \mathcal{R}^{(d)}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}^{(d)}} \text{regret}(R, f) \quad (146)$$

$$= \text{regret}_{\text{pure}}(\mathcal{R}^{(d)}, \mathcal{F}_{\{\lambda_j\}}^{(d)} \mid P_{\mathbf{x}}), \quad (147)$$

1002 where  $(*)$  follows from the following reasoning: For any  $(R \in \mathcal{R}, f \in \mathcal{F}_{\{\lambda_j\}}^{(d)})$ ,

$$\text{regret}(R, f) = \min_{q \in \mathbb{R}^r} \mathbb{E} \left[ \left( \sum_{j=1}^d f_j \phi_j(\mathbf{x}) + \mathbf{n} - \sum_{j=1}^{\infty} \sum_{i=1}^r q_i c_{ij} \phi_j(\mathbf{x}) \right)^2 \right] - \mathbb{E} [\mathbf{n}^2] \quad (148)$$

$$\stackrel{(a)}{=} \min_{q \in \mathbb{R}^r} \mathbb{E} \left[ \left( \sum_{j=1}^d f_j \phi_j(\mathbf{x}) - \sum_{j=1}^{\infty} \sum_{i=1}^r q_i c_{ij} \phi_j(\mathbf{x}) \right)^2 \right] \quad (149)$$

$$\stackrel{(b)}{=} \min_{q \in \mathbb{R}^r} \sum_{j=1}^d \left( f_j - \sum_{i=1}^r q_i c_{ij} \right)^2 + \sum_{j=d+1}^{\infty} \left( \sum_{i=1}^r q_i c_{ij} \right)^2, \quad (150)$$

1003 where here (a) follows since the noise  $\mathbf{n}$  is independent of  $\mathbf{x}$ , and since, similarly to the finite-  
 1004 dimensional case (Section 3), the prediction loss based on the features  $x \in \mathcal{X}$  is  $\mathbb{E}[\mathbf{n}^2]$ , for any  
 1005 given  $f \in \mathcal{F}$ , (b) follows from Parseval's identity and the orthonormality of  $\{\phi_j\}_{j \in \mathbb{N}}$ . So,

$$\begin{aligned} & \min_{R \in \mathcal{R}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}^{(d)}} \text{regret}(R, f) \\ &= \min_{\{c_{ij}\}_{i \in [r], j \in \mathbb{N}}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}^{(d)}} \min_{q \in \mathbb{R}^r} \sum_{j=1}^d \left( f_j - \sum_{i=1}^r q_i c_{ij} \right)^2 + \sum_{j=d+1}^{\infty} \left( \sum_{i=1}^r q_i c_{ij} \right)^2. \end{aligned} \quad (151)$$

1006 Evidently, since  $\sum_{j=d+1}^{\infty} (\sum_{i=1}^r q_i c_{ij})^2 \geq 0$ , an optimal representation may satisfy that  $c_{ij} = 0$  for  
 1007 all  $j \geq d+1$ . Thus, the optimal representation belongs to  $\mathcal{R}^{(d)}$ .

1008 On the other hand,

$$\text{regret}_{\text{pure}}(\mathcal{R}, \mathcal{F}_{\{\lambda_j\}} \mid P_{\mathbf{x}}) = \min_{R \in \mathcal{R}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}} \text{regret}(R, f) \quad (152)$$

$$\leq \min_{R \in \mathcal{R}^{(d)}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}} \text{regret}(R, f) \quad (153)$$

$$\stackrel{(*)}{\leq} \min_{R \in \mathcal{R}^{(d)}} \max_{f \in \mathcal{F}_{\{\lambda_j\}}^{(d)}} \text{regret}(R, f) + \lambda_{d+1} \quad (154)$$

$$= \text{regret}_{\text{pure}}(\mathcal{R}^{(d)}, \mathcal{F}_{\{\lambda_j\}}^{(d)} \mid P_{\mathbf{x}}) + \lambda_{d+1}, \quad (155)$$

1009 where (\*) follows from the following reasoning: For any  $(R \in \mathcal{R}^{(d)}, f \in \mathcal{F}_{\{\lambda_j\}})$ ,

$$\text{regret}(R, f) = \min_{q \in \mathbb{R}^r} \mathbb{E} \left[ \left( \sum_{j=1}^{\infty} f_j \phi_j(\mathbf{x}) + \mathbf{n} - \sum_{j=1}^{\infty} \sum_{i=1}^r q_i c_{ij} \phi_j(\mathbf{x}) \right)^2 \right] - \mathbb{E}[\mathbf{n}^2] \quad (156)$$

$$\stackrel{(a)}{=} \min_{q \in \mathbb{R}^r} \sum_{j=1}^d \left( f_j - \sum_{i=1}^r q_i c_{ij} \right)^2 + \sum_{j=d+1}^{\infty} f_j^2 \quad (157)$$

$$\stackrel{(b)}{\leq} \min_{q \in \mathbb{R}^r} \sum_{j=1}^d \left( f_j - \sum_{i=1}^r q_i c_{ij} \right)^2 + \lambda_{d+1}, \quad (158)$$

1010 where (a) follows similarly to the analysis made in the previous step, and (b) follows since for any  
 1011  $f \in \mathcal{F}_{\{\lambda_j\}}$  it holds that

$$\sum_{j=d+1}^{\infty} f_j^2 \leq \lambda_{d+1} \sum_{j=d+1}^{\infty} \frac{f_j^2}{\lambda_j} \leq \lambda_{d+1} \sum_{j=1}^{\infty} \frac{f_j^2}{\lambda_j} \leq \lambda_{d+1}. \quad (159)$$

1012 Combining (147) and (155) and using  $\lambda_{d+1} \downarrow 0$  completes the proof for the pure minimax regret.  
 1013 The proof for the mixed minimax is analogous and thus is omitted.  $\square$

1014 We also use the following simple and technical lemma.

1015 **Lemma 21.** For  $R \in \mathcal{R}^{(d)}$  and  $f \in \mathcal{F}^{(d)}$ <sup>1</sup>

$$\text{regret}(R, f) = f^\top (I_d - R^\top (R R^\top)^{-1} R) f, \quad (160)$$

1016 where  $R \in \mathbb{R}^{r \times d}$  is the matrix of coefficients of the orthogonal expansion of  $\psi_i = \sum_{j=1}^d c_{ij} \phi_j$  for  
 1017  $i \in [r]$ , so that  $R(i, j) = c_{ij}$ .

<sup>1</sup>Note that any  $f \in \mathcal{F}^{(d)}$  may be uniquely identified with a  $d$ -dimensional vector  $f \in \mathbb{R}^d$ . With a slight abuse of notation we do not distinguish between the two.

1018 *Proof.* It holds that

$$\text{regret}(R, f) = \min_{q \in \mathbb{R}^r} \mathbb{E} \left[ \left( \sum_{j=1}^d f_j \phi_j(\mathbf{x}) + \mathbf{n} - \sum_{i=1}^r q_i \sum_{j=1}^d c_{ij} \phi_j(\mathbf{x}) \right)^2 \right] - \mathbb{E} [\mathbf{n}^2] \quad (161)$$

$$= \min_{q \in \mathbb{R}^r} \mathbb{E} \left[ \sum_{j=1}^d \left( f_j - \sum_{i=1}^r q_i c_{ij} \right) \phi_j(\mathbf{x}) \right] \quad (162)$$

$$= \min_{q \in \mathbb{R}^r} \sum_{j=1}^d \left( f_j - \sum_{i=1}^r q_i c_{ij} \right)^2 \quad (163)$$

$$= \min_{q \in \mathbb{R}^r} \sum_{j=1}^d \left[ f_j^2 - 2f_j \sum_{i=1}^r q_i c_{ij} + \sum_{i_1=1}^r \sum_{i_2=1}^r q_{i_1} c_{i_1 j} q_{i_2} c_{i_2 j} \right] \quad (164)$$

$$= \min_{q \in \mathbb{R}^r} f^\top f - 2q^\top Rf + q^\top RR^\top q \quad (165)$$

$$= f^\top (I_d - R^\top (RR^\top)^{-1} R) f, \quad (166)$$

1019 where the last equality is obtained by the minimizer  $q^* = (RR^\top)^{-1} Rf$ .  $\square$

1020 *Proof of Theorems 18 and 19.* By Proposition 20, we may first consider the finite dimensional case,  
 1021 and then take the limit  $d \uparrow \infty$ . By Lemma 21, in the  $d$ -dimensional case (for both the representation  
 1022 and the response function), the regret is formally as in the linear setting under the MSE of Theorem  
 1023 2, by setting therein  $\Sigma_{\mathbf{x}} = I_d$ , and  $S = \text{diag}(\lambda_1, \dots, \lambda_d)$  (c.f. Lemma 15). The claim of the  
 1024 Theorem 18 then follows by taking  $d \uparrow \infty$  and noting that  $\lambda_{d+1} \downarrow 0$ . The proof of Theorem 19 is  
 1025 analogous and thus omitted.  $\square$

## 1026 G Iterative algorithms for the Phase 1 and Phase 2 problems

1027 In this section we describe our proposed algorithms for the solution Phase 1 and Phase 2 problems  
 1028 of Algorithm 1. Those algorithms are general, and only require providing gradients of the regret  
 1029 function (1) and an initial representation and a set of adversarial functions. These are individually  
 1030 determined for each setting. See Section H for the way these are determined in Examples 6 and 8.

### 1031 G.1 Phase 1: finding a new adversarial function

1032 We propose an algorithm to solve the Phase 1 problem (26), which is again based on an iterative  
 1033 algorithm. We denote the function's value at the  $t$ th iteration by  $f_{(t)}$ . The proposed Algorithm 2  
 1034 operates as follows. At initialization, the function  $f_{(1)} \in \mathcal{F}$  is arbitrarily initialized (say at random),  
 1035 and then the optimal predictor  $Q^{(j)}$  is found for each of the  $k$  possible representations  $R^{(j)}$ ,  $j \in [k]$ .  
 1036 Then, the algorithm iteratively repeats the following steps, starting with  $t = 2$ : (1) Updating the  
 1037 function from  $f_{(t-1)}$  to  $f_{(t)}$  based on a gradient step of

$$\sum_{j \in [k]} p^{(j)} \cdot \mathbb{E} \left[ \text{loss}(f_{(t-1)}(\mathbf{x}), Q^{(j)}(R^{(j)}(\mathbf{x}))) \right], \quad (167)$$

1038 that is, the weighted loss function of the previous iteration function, which is then followed by a  
 1039 projection to the feasible class of functions  $\mathcal{F}$ , denoted as  $\Pi_{\mathcal{F}}(\cdot)$  (2) Finding the optimal predictor  
 1040  $Q^{(j)}$  for the current function  $f_{(t)}$  and the given representations  $\{R^{(j)}\}_{j \in [k]}$ , and computing the  
 1041 respective loss for each representation,

$$L^{(j)} := \mathbb{E} \left[ \text{loss}(f_{(t)}(\mathbf{x}), Q^{(j)}(R^{(j)}(\mathbf{x}))) \right]. \quad (168)$$

1042 This loop iterates for  $T_f$  iterations, or until convergence.

---

**Algorithm 2** A procedure for finding a new function via the solution of (26)

---

```

1: procedure PHASE 1 SOLVER( $\{R^{(j)}, p^{(j)}\}_{j \in [k]}, \mathcal{F}, \mathcal{Q}, d, r, P_{\mathbf{x}}$ )
2:   begin
3:     Initialize  $T_f$  ▷ Number of iterations parameters
4:     Initialize  $\eta_f$  ▷ Step size parameter
5:     Initialize  $f_{(1)} \in \mathcal{F}$  ▷ Function initialization, e.g., at random
6:     for  $j = 1$  to  $k$  do
7:       set  $Q^{(j)} \leftarrow \arg \min_{Q \in \mathcal{Q}} \mathbb{E} [\text{loss}(f_{(1)}(\mathbf{x}), Q(R^{(j)}(\mathbf{x})))]$ 
8:     end for
9:     for  $t = 2$  to  $T_f$  do
10:      update  $f_{(t-1/2)} = f_{(t-1)} + \eta_f \cdot \sum_{j \in [k]} p^{(j)}_{(t-1)} \cdot \nabla_f \mathbb{E} [\text{loss}(f_{(t-1)}(\mathbf{x}), Q^{(j)}(R^{(j)}(\mathbf{x})))]$ 
      ▷ A gradient update of the function
11:      project  $f_{(t)} = \Pi_{\mathcal{F}}(f_{(t-1/2)})$  ▷ Projection on the class  $\mathcal{F}$ 
12:      for  $j = 1$  to  $k$  do
13:        set  $Q^{(j)} \leftarrow \arg \min_{Q \in \mathcal{Q}} \mathbb{E} [\text{loss}(f_{(t)}(\mathbf{x}), Q(R^{(j)}(\mathbf{x})))]$  ▷ Update of predictors
14:        set  $L^{(j)} \leftarrow \mathbb{E} [\text{loss}(f_{(t)}(\mathbf{x}), Q^{(j)}(R^{(j)}(\mathbf{x})))]$  ▷ Compute loss of each representation
15:      end for
16:    end for
17:    return  $f_{(T)}$ , and the regret  $\sum_{j \in [k]} p^{(j)} \cdot L^{(j)}$ 
18:  end procedure

```

---

1043 **Design choices and possible variants of the basic algorithm** At initialization, we have chosen a  
1044 simple random initialization for  $f_{(1)}$ , but it may also be initialized based on some prior knowledge  
1045 of the adversarial function. For the update of the predictors, we have specified a full computation of  
1046 the optimal predictor, which can be achieved in practice by running another iterative algorithm such  
1047 as stochastic gradient descent (SGD) until convergence. If this is too computationally expensive,  
1048 the number of gradient steps may be limited. The update of the function is done via projected SGD  
1049 with a constant step size  $\eta_f$ , yet it is also possible to modify the step size with the iteration, e.g., the  
1050 common choice  $\eta_f/\sqrt{t}$  at step  $t$  Hazan [35]. Accelerated algorithms, e.g., moment-based may also  
1051 be deployed.

1052 **Convergence analysis** A theoretical analysis of the convergence properties of the algorithm ap-  
1053 pears to be challenging. Evidently, this is a minimax game between the response player and the  
1054 predictor player, but not a concave-convex game. As described in Appendix B, even concave-convex  
1055 games are not well understood at this point. We thus opt to validate this algorithm numerically.

## 1056 G.2 Phase 2: finding a new representation

1057 We propose an iterative algorithm to solve the Phase 2 problem (27), and thus finding a new repre-  
1058 sentation  $R^{(k+1)}$ . To this end, we first note that the objective function in (27) can be separated into  
1059 a part that depends on existing representations and a part that depends on the new one, specifically,  
1060 as

$$\begin{aligned}
& \sum_{j_1 \in [k]} \sum_{j_2 \in [m_0+k]} p^{(j_1)} \cdot o^{(j_2)} \cdot \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(j_1, j_2)}(R^{(j_1)}(\mathbf{x}))) \right] \\
& \quad + \sum_{j_2 \in [m_0+k]} p^{(k+1)} \cdot o^{(j_2)} \cdot \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(k+1, j_2)}(R^{(k+1)}(\mathbf{x}))) \right] \\
& = \sum_{j_1 \in [k]} \sum_{j_2 \in [m_0+k]} p^{(j_1)} \cdot o^{(j_2)} \cdot L^{(j_1, j_2)} \\
& \quad + \sum_{j_2 \in [m_0+k]} p^{(k+1)} \cdot o^{(j_2)} \cdot \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(k+1, j_2)}(R^{(k+1)}(\mathbf{x}))) \right], \tag{169}
\end{aligned}$$

1061 where

$$L^{(j_1, j_2)} := \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(j_1, j_2)}(R^{(j_1)}(\mathbf{x}))) \right], \tag{170}$$

1062 and the predictors  $\{Q^{(j_1, j_2)}\}_{j_1 \in [k], j_2 \in [m_0 + k]}$  can be optimized independently of the new representa-  
 1063 tion  $R^{(k+1)}$ . We propose an iterative algorithm for this problem, and denote the new representation  
 1064 at the  $t$ th iteration of the algorithm by  $R_{(t)}^{(k+1)}$ . The algorithm's input is a set of  $m_0 + k$  adver-  
 1065 sarial functions  $\{f^{(i)}\}_{i \in [m_0 + k]}$ , and the current set of representations  $\{R^{(j)}\}_{j \in [k]}$ . Based on these,  
 1066 the algorithm may find the optimal predictor for  $f^{(j_2)}$  based on the representation  $R^{(j_1)}$ , and thus  
 1067 compute the loss

$$L_*^{(j_1, j_2)} := \min_{Q \in \mathcal{Q}} \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q(R^{(j_1)}(\mathbf{x}))) \right] \quad (171)$$

1068 for  $j_1 \in [k]$  and  $j_2 \in [m_0 + k]$ . In addition, the new representation is arbitrarily initialized (say, at  
 1069 random) as  $R_{(1)}^{(k+1)}$ , and the predictors  $\{Q_{(1)}^{(k+1, j_2)}\}_{j_2 \in [m_0 + k]}$  are initialized as the optimal predictors  
 1070 for  $f^{(j_2)}$  given the representation  $R_{(1)}^{(k+1)}$ . The algorithm keeps track of weights for the represen-  
 1071 tations (including the new one), which are initialized uniformly, i.e.,  $p_{(1)}^{(j_1)} = \frac{1}{k+1}$  for  $j_1 \in [k+1]$   
 1072 (including a weight for the new representation). The algorithm also keeps track of weights for the  
 1073 functions, which are also initialized uniformly as  $o_{(1)}^{(j_2)} = \frac{1}{m_0 + k}$  for  $j_2 \in [m_0 + k]$ . Then, the  
 1074 algorithm iteratively repeats the following steps, starting with  $t = 2$ : (1) Updating the new represen-  
 1075 tation from  $R_{(t-1)}^{(k+1)}$  to  $R_{(t)}^{(k+1)}$  based on a gradient step of the objective function (27) as a function of  
 1076  $R^{(k+1)}$ . Based on the decomposition in (169) the term of the objective which depends on  $R^{(k+1)}$  is

$$p_{(t-1)}^{(k+1)} \sum_{j_2 \in [m_0 + k]} o_{(t-1)}^{(j_2)} \cdot \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(k+1, j_2)}(R^{(k+1)}(\mathbf{x}))) \right], \quad (172)$$

1077 that is, the loss function of the previous iteration new representation, weighted according to the  
 1078 current function weights  $o_{(t-1)}^{(j_2)}$ . Since the multiplicative factor  $p_{(t-1)}^{(k+1)}$  is common to all terms, it is  
 1079 removed from the gradient computation (this aids in the choice of the gradient step). This gradient  
 1080 step is then possibly followed by normalization or projection, which we denote by the operator  
 1081  $\Pi_{\mathcal{R}}(\cdot)$ . For example, in the linear case, it make sense to normalize  $R^{(k+1)}$  to have unity norm (in  
 1082 some matrix norm of choice). After updating the new representation to  $R_{(t)}^{(k+1)}$ , optimal predictors  
 1083 are found for each function, the loss is computed

$$L_{(t)}^{(k+1, j_2)} := \min_{Q \in \mathcal{Q}} \mathbb{E} \left[ \text{loss}(f^{(j_2)}(\mathbf{x}), Q(R_{(t)}^{(k+1)}(\mathbf{x}))) \right] \quad (173)$$

1084 for all  $j_2 \in [m_0 + k]$ , and the optimal predictor is updated to  $\{Q_{(t)}^{(k+1, j_2)}\}_{j_2 \in [m_0 + k]}$  based on this  
 1085 solution. (2) Given the current new representation  $R_{(t)}^{(k+1)}$ , the loss matrix

$$\{L_{(t)}^{(j_1, j_2)}\}_{j_1 \in [k], j_2 \in [m_0 + k]} \quad (174)$$

1086 is constructed where for  $j_1 \in [k]$  it holds that  $L_{(t)}^{(j_1, j_2)} = L^{(j_1, j_2)}$  for all  $t$  (i.e., the loss of previous  
 1087 representations and functions is kept fixed). This is considered to be the loss matrix of a two-player  
 1088 zero-sum game between the representation player and the function player, where the representation  
 1089 player has  $k + 1$  possible strategies and the function player has  $m_0 + k$  strategies. The weights  
 1090  $\{p_{(t)}^{(j_1)}\}_{j_1 \in [k+1]}$  and  $\{o_{(t)}^{(j_2)}\}_{j_2 \in [m_0 + k]}$  are then updated according to the MWU rule. Specifically, for  
 1091 an *inverse temperature parameter*  $\beta$  (or a *regularization parameter*), the update is given by

$$p_{(t)}^{(j)} = \frac{p_{(t-1)}^{(j)} \cdot \beta^{L^{(j)}}}{\sum_{\tilde{j} \in [k]} p_{(t-1)}^{(\tilde{j})} \cdot \beta^{L^{(\tilde{j})}}} \quad (175)$$

1092 for the representation weights and, analogously, by

$$o_{(t)}^{(j)} = \frac{o_{(t-1)}^{(j)} \cdot \beta^{-L^{(j)}}}{\sum_{\tilde{j} \in [k]} o_{(t-1)}^{(\tilde{j})} \cdot \beta^{-L^{(\tilde{j})}}} \quad (176)$$

1093 for the function weights (as the function player aims to maximize the loss). This can be considered as  
 1094 a regularized gradient step on the probability simplex, or more accurately, a *follow-the-regularized-*  
 1095 *leader* [35]. The main reasoning of this algorithm is that at each iteration the weights  $\{p^{(j)}\}_{j \in [k+1]}$

1096 and  $\{o^{(j)}\}_{j \in [m_0+k]}$  are updated towards the solution of the two-player zero-sum game with payoff  
 1097 matrix  $\{-L_{(t)}^{(j_1, j_2)}\}_{j_1 \in [k+1], j_2 \in [m_0+k]}$ . In turn, based only on the function weights  $\{o^{(j)}\}_{j \in [m_0+k]}$ ,  
 1098 the new representation is updated to  $R_{(t)}^{(k+1)}$ , which then changes the pay-off matrix at the next  
 1099 iteration. It is well known that the MWU solved two-player zero-sum game [33], in which the  
 1100 representation player can choose the weights and the function player can choose the function.  
 1101 This loop iterates for  $T_{\text{stop}}$  iterations, and then the optimal weights are given by the average over the  
 1102 last  $T_{\text{avg}}$  iterations [33], i.e.,

$$p_*^{(j)} = \frac{1}{T_{\text{avg}}} \sum_{t=T_{\text{stop}}-T_{\text{avg}}+1}^{T_{\text{stop}}} p_{(t)}^{(j)}, \quad (177)$$

1103 and

$$o_*^{(j)} = \frac{1}{T_{\text{avg}}} \sum_{t=T_{\text{stop}}-T_{\text{avg}}+1}^{T_{\text{stop}}} o_{(t)}^{(j)}. \quad (178)$$

1104 In the last  $T_R - T_{\text{stop}}$  iterations, only the representation  $R_{(t)}^{(k+1)}$  and the predictors are updated. The  
 1105 algorithm then outputs  $R_{(T)}^{(k+1)}$  as the new representation and the weights  $\{p_*^{(j)}\}_{j \in [k+1]}$ .

1106 **Design choices and possible variants of the basic algorithm** At initialization, we have cho-  
 1107 sen a simple random initialization for  $R_{(1)}^{(k+1)}$ , but it may also be initialized based on some prior  
 1108 knowledge of the desired new representation. The initial predictors  $\{Q_{(1)}^{(k+1, j_2)}\}_{j_2 \in [m_0+k]}$  will then  
 1109 be initialized as the optimal predictors for  $R_{(1)}^{(k+1)}$  and  $\{f^{(j_2)}\}_{j_2 \in [m_0+k]}$ . We have initialized the  
 1110 representation and function weights uniformly. A possibly improved initialization for the function  
 1111 weights is to put more mass on the more recent functions, that is, for large values of  $j_2$ , or to use  
 1112 the minimax strategy of the function player in the two-player zero-sum game with payoff matrix  
 1113  $\{-L_{(t)}^{(j_1, j_2)}\}_{j_1 \in [k], j_2 \in [m_0+k]}$  (that is, a game which does not include the new representation). As  
 1114 in the Phase 1 algorithm, the gradient update of the new representation can be replaced by a more  
 1115 sophisticated algorithm, the computation of the optimal predictors can be replaced with (multiple)  
 1116 update steps, and the step size may also be adjusted. For the MWU update, we use the proposed  
 1117 scaling in [33]

$$\beta = \frac{1}{1 + \sqrt{\frac{c \ln m}{T}}} \quad (179)$$

1118 for some constant  $c$ . It is well known that using the last iteration of a MWU algorithm may fail [97],  
 1119 while averaging the weights value of all iterations provides the optimal value of a two-player zero-  
 1120 sum games [33]. For improved accuracy, we compute the average weights over the last  $T_{\text{avg}}$  iterations  
 1121 (thus disregarding the initial iterations). We then halt the weights update and let the function and  
 1122 predictor update to run for  $T - T_{\text{stop}}$  iterations in order to improve the convergence of  $R^{(k+1)}$ .  
 1123 Finally, the scheduling of the steps may be more complex, e.g., it is possible that running multiple  
 1124 gradient steps follows by multiple MWU steps may improve the result.

## 1125 H Details for the examples of Algorithm 1

1126 As mentioned, the solvers of the Phase 1 and Phase 2 problems of Algorithm 1 require the gradients  
 1127 of the regret (1) as inputs, as well as initial representation and set of adversarial functions. We next  
 1128 provide these details for the examples in Section 4. The code for the experiments was written in  
 1129 Python 3.6 the code is available at this link. The optimization of hyperparameters was done using  
 1130 the Optuna library. The hardware used is standard and detailed appear in Table 1.

### 1131 H.1 Details for Example 6: the linear MSE setting

1132 In this setting, the expectation over the feature distribution can be carried out analytically, and the  
 1133 regret is given by

$$\text{regret}(R, f \mid \Sigma_{\mathbf{x}}) = \mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x})^2 \right] \quad (180)$$

---

**Algorithm 3** A procedure for finding a new representation  $R^{(k+1)}$  via the solution of (27)

---

```

1: procedure PHASE 2 SOLVER( $\{R^{(j_1)}\}_{j \in [k]}, \{f^{(j_2)}\}_{j_2 \in [m_0+k]}, \mathcal{R}, \mathcal{F}, \mathcal{Q}, d, r, P_{\mathbf{x}}$ )
2:   begin
3:     Initialize  $T_R, T_{\text{stop}}, T_{\text{avg}}$  ▷ Number of iterations parameters
4:     Initialize  $\eta_R$  ▷ Step size parameter
5:     Initialize  $\beta \in (0, 1)$  ▷ Inverse temperature parameter
6:     Initialize  $f_{(1)} \in \mathcal{F}$  ▷ Function initialization, e.g., at random
7:     Initialize  $p_{(1)}^{(j)} \leftarrow 0$  for  $j \in [k]$  and  $p_{(1)}^{(k+1)} \leftarrow 0$  ▷ A uniform weight initialization for the
representations
8:     Initialize  $o_{(1)}^{(j_2)} \leftarrow \frac{1}{m_0+k}$  for  $j_2 \in [k]$  ▷ A uniform weight initialization for the functions
9:     for  $j_1 = 1$  to  $k$  do
10:      for  $j_2 = 1$  to  $m_0 + k$  do
11:        Set  $Q^{(j_1, j_2)} \leftarrow \arg \min_{Q \in \mathcal{Q}} \mathbb{E} [\text{loss}(f^{(j_2)}(\mathbf{x}), Q(R^{(j_1)}(\mathbf{x})))]$ 
▷ Optimal predictors for existing representations and input functions
12:        Set  $L^{(j_1, j_2)} \leftarrow \min_{Q \in \mathcal{Q}} \mathbb{E} [\text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(j_1, j_2)}(R^{(j_1)}(\mathbf{x})))]$  ▷ The minimal loss
13:      end for
14:    end for
15:    for  $j_2 = 1$  to  $m_0 + k$  do
16:      Initialize  $R_{(1)}^{(k+1)}$  ▷ Arbitrarily, e.g., at random
17:      Set  $Q_{(1)}^{(k+1, j_2)} \leftarrow \arg \min_{Q \in \mathcal{Q}} \mathbb{E} [\text{loss}(f^{(j_2)}(\mathbf{x}), Q(R_{(1)}^{(k+1)}(\mathbf{x})))]$  for  $j_2 \in [m_0 + k]$ 
▷ Optimal predictors for new representation and input functions
18:    end for
19:    for  $t = 2$  to  $T_R$  do
20:      update ▷ A gradient update of the new representation

$$R_{(t-1/2)}^{(k+1)} = R_{(t-1)}^{(k+1)} + \eta_R \cdot \sum_{j_2 \in [m_0+k]} o_{(t-1)}^{(j_2)} \cdot \nabla_{R^{(k+1)}} \mathbb{E} [\text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(k+1, j_2)}(R_{(t-1)}^{(k+1)}(\mathbf{x})))]$$

21:      projection  $R_{(t)}^{(k+1)} = \Pi_{\mathcal{R}}(R_{(t-1/2)}^{(k+1)})$  ▷ Standardization based on the class  $\mathcal{R}$ 
22:      for  $j = 1$  to  $k$  do
23:        Set  $Q^{(k+1, j_2)} \leftarrow \arg \min_{Q \in \mathcal{Q}} \mathbb{E} [\text{loss}(f^{(j_2)}(\mathbf{x}), Q(R_{(t)}^{(k+1)}(\mathbf{x})))]$ 
▷ Update of predictors for the new representation
24:         $L_{(t)}^{(k+1, j_2)} \leftarrow \mathbb{E} [\text{loss}(f^{(j_2)}(\mathbf{x}), Q^{(k+1, j_2)}(R_{(t)}^{(k+1)}(\mathbf{x})))]$  ▷ Compute loss
25:      end for
26:      Set  $L_{(t)}^{(j_1, j_2)} \leftarrow L^{(j_1, j_2)}$  for  $j_1 \in [k]$  and  $j_2 \in [m_0 + k]$ 
27:      if  $t < T_{\text{stop}}$  then
28:        update  $p_{(t)}^{(j)} \leftarrow \frac{p_{(t-1)}^{(j)} \cdot \beta^{L^{(j)}}}{\sum_{\tilde{j} \in [k]} p_{(t-1)}^{(\tilde{j})} \cdot \beta^{L^{(\tilde{j})}}}$  for  $j \in [k]$  ▷ A MWU
29:        update  $o_{(t)}^{(j)} \leftarrow \frac{o_{(t-1)}^{(j)} \cdot \beta^{-L^{(j)}}}{\sum_{\tilde{j} \in [m_0+k]} o_{(t-1)}^{(\tilde{j})} \cdot \beta^{-L^{(\tilde{j})}}}$  for  $j \in [m_0 + k]$  ▷ A MWU
30:      else if  $t = T_{\text{stop}}$  then
31:        update  $p_{(t)}^{(j)} = p_{(t)}^{(j)} \leftarrow \frac{1}{T_{\text{avg}}} \sum_{t=T_{\text{stop}}-T_{\text{avg}}+1}^{T_{\text{stop}}} p_{(t)}^{(j)}$  for  $j \in [k]$ 
▷ Optimal weights by averaging last  $T_{\text{avg}}$  iterations
32:        update  $o_{(t)}^{(j)} \leftarrow \frac{1}{T_{\text{avg}}} \sum_{t=T_{\text{stop}}-T_{\text{avg}}+1}^{T_{\text{stop}}} o_{(t)}^{(j)}$  for  $j \in [m_0 + k]$ 
▷ Optimal weights by averaging last  $T_{\text{avg}}$  iterations
33:      else
34:        update  $p_{(t)}^{(j)} \leftarrow p_{(t-1)}^{(j)}$  for  $j \in [k]$  ▷ No update for the last  $T - T_{\text{stop}}$  iterations
35:        update  $o_{(t)}^{(j)} \leftarrow o_{(t-1)}^{(j)}$  for  $j \in [m_0 + k]$  ▷ No update for the last  $T - T_{\text{stop}}$  iterations
36:      end if
37:      return  $R_{(T)}^{(k+1)}$  and  $\{p_{(T_R)}^{(j)}\}_{j \in [k+1]}$ 
38:    end for
39:  end procedure

```

---

Table 1: Hardware details

CPU	RAM	GPU
Intel i9 13900k	64GB	RTX 3090 Ti

$$= f^\top \Sigma_{\mathbf{x}} f - 2q^\top R^\top \Sigma_{\mathbf{x}} f + q^\top R^\top R q. \quad (181)$$

1134 The regret only depends on the feature distribution  $P_{\mathbf{x}}$  via  $\Sigma_{\mathbf{x}}$ . For each run of the algorithm, the co-  
 1135 variance matrix  $\Sigma_{\mathbf{x}}$  was chosen to be diagonal with elements drawn from a log-normal distribution,  
 1136 with parameters  $(0, \sigma_0)$ , and  $S = I_d$ .

1137 **Regret gradients** The gradient of the regret w.r.t. the function  $f$  is given by

$$\nabla_f \mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x})^2 \right] = 2f^\top \Sigma_{\mathbf{x}} - 2q^\top R^\top \Sigma_{\mathbf{x}} \quad (182)$$

1138 and the projection on  $\mathcal{F}_S$  is

$$\Pi_{\mathcal{F}}(f) = \begin{cases} \frac{f}{\|f\|_S}, & \|f\|_S \geq 1 \\ f, & \|f\|_S < 1 \end{cases}. \quad (183)$$

1139 However, we may choose to normalize by  $\frac{f}{\|f\|_S}$  even if  $\|f\|_S \leq 1$  since in this case the regret is  
 1140 always larger if  $f$  is replaced by  $\frac{f}{\|f\|_S}$  (in other words, the worst case function is obtained on the  
 1141 boundary of  $\mathcal{F}_S$ ). The gradient w.r.t. the predictor  $q$  is given by

$$\nabla_q \mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x})^2 \right] = [-2f^\top \Sigma_{\mathbf{x}} R + 2q^\top R^\top \Sigma_{\mathbf{x}} R]. \quad (184)$$

1142 Finally, to derive the gradient w.r.t.  $R$ , let us denote  $R := [R_1, R_2, \dots, R_r] \in \mathbb{R}^{d \times r}$  where  $R_i \in \mathbb{R}^d$   
 1143 is the  $i$ th column ( $i \in [r]$ ), and  $q^\top = (q_1, q_2, \dots, q_r)$ . Then,  $q^\top R^\top \mathbf{x} = \sum_{i \in [d]} q_i R_i^\top \mathbf{x}$  and the loss  
 1144 function is

$$\mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x})^2 \right] = \mathbb{E} \left[ \left( f^\top \mathbf{x} - \sum_{i \in [d]} q_i \mathbf{x}^\top R_i \right)^2 \right] \quad (185)$$

$$= f^\top \Sigma_{\mathbf{x}} f - 2q^\top R^\top \Sigma_{\mathbf{x}} f + q^\top R^\top \Sigma_{\mathbf{x}} R q. \quad (186)$$

1145 The gradient of the regret w.r.t.  $R_k$  is then given by

$$\nabla_{R_k} \left\{ \mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x})^2 \right] \right\} = -2\mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x}) \cdot q_k \mathbf{x}^\top \right] \quad (187)$$

$$= -2q_k (f^\top \Sigma_{\mathbf{x}} - q^\top R^\top \Sigma_{\mathbf{x}}), \quad (188)$$

1146 hence, more succinctly, the gradient w.r.t.  $R$  is

$$\nabla_R \left\{ \mathbb{E} \left[ (f^\top \mathbf{x} - q^\top R^\top \mathbf{x})^2 \right] \right\} = -2q (f^\top \Sigma_{\mathbf{x}} - q^\top R^\top \Sigma_{\mathbf{x}}). \quad (189)$$

1147 We remark that in the algorithm these gradients are multiplied by weights. We omit this term when-  
 1148 ever the weight is common to all terms in order to keep the effective step size constant.

1149 **Initialization** Algorithm 1 requires an initial representation  $R^{(1)}$  and an initial set of functions  
 1150  $\{f^{(j)}\}_{j \in [m_0]}$ . In the MSE setting, each function  $f \in \mathbb{R}^d$  is also a single column of a representation  
 1151 matrix  $R \in \mathbb{R}^{d \times r}$ . A plausible initialization matrix  $R^{(1)} \in \mathbb{R}^{d \times r}$  is therefore the worst  $r$  functions.  
 1152 These, in turn, can be found by running Algorithm (1) to obtain  $\tilde{m} = r$  functions, by setting  $\tilde{r} = 1$ . A  
 1153 proper initialization for this run is simply an all-zero representation  $\tilde{R}^{(1)} = 0 \in \mathbb{R}^{d \times 1}$ . The resulting  
 1154 output is then  $\{\tilde{R}_{(T)}^{(j)}\}_{j \in [r]}$  which can be placed as the  $r$  columns of  $R^{(1)}$ . This initialization is then  
 1155 used for Algorithm 1.

1156 **Algorithm parameters** The algorithm parameters used for Example 6 are shown in Table 2. The  
 1157 parameters were optimally tuned for  $\sigma_0 = 1$ .

Table 2: Parameters for linear MSE setting example

Parameter	$\beta_r$	$\beta_f$	$\eta_r$	$\eta_f$
Value	0.94	0.653	0.713	0.944
Parameter	$T_R$	$T_f$	$T_{\text{avg}}$	$T_{\text{stop}}$
Value	100	until convergence	10	80

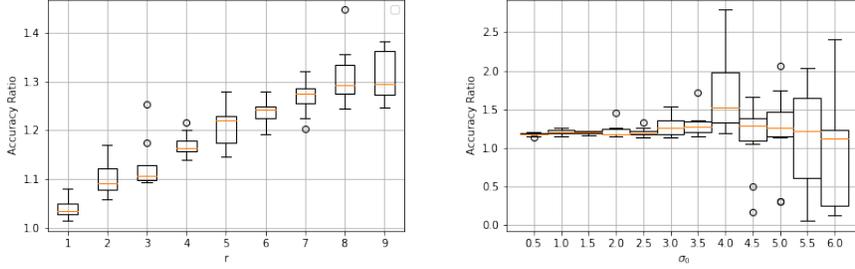


Figure 3: The ratio between the regret achieved by Algorithm 1 and the theoretical regret in the linear MSE setting. Left:  $d = 20, \sigma_0 = 1$ , varying  $r$ . Right:  $r = 5, d = 20$ , varying  $\sigma_0$ .

1158 **Additional results** Additional results of the accuracy of the Algorithm 1 in the linear MSE setting  
 1159 are displayed in Figure 3. The left panel of Figure 3 shows that the algorithm output is accurate for  
 1160 small values of  $r$ , but deteriorates as  $r$  increases. This is because when  $r$  increases then so is  $\ell^*$  and  
 1161 so is the required number of matrices in the support of the representation rule (denoted by  $m$ ). Since  
 1162 the algorithm gradually adds representation matrices to the support, an inaccurate convergence at an  
 1163 early iteration significantly affects later iterations. One possible way to remedy this is to run each  
 1164 iteration multiple times, and choose the best one, before moving on to the next one. Another reason  
 1165 is that given large number of matrices in the support (large  $m$ ), it becomes increasingly difficult  
 1166 for the MWU to accurately converge. Since the iterations of the MWU do not converge to the  
 1167 equilibrium point, but rather their average (see discussion in Appendix B) this can only be remedied  
 1168 by allowing more iterations for convergence (in advance) for large values of  $m$ . The right panel of  
 1169 Figure 3 shows that the algorithm output is accurate for a wide range of the condition number of  
 1170 the covariance matrix. This condition number is determined by the choice of  $\sigma_0$ , where low values  
 1171 typically result covariance matrices with condition number that is close to 1, while high values will  
 1172 typically result large condition number. The right panel shows that while the hyperparameters were  
 1173 tuned for  $\sigma_0 = 1$ , the result is fairly accurate for a wide range of  $\sigma_0$  values, up to  $\sigma_0 \approx 5$ . Since for  
 1174  $Z \sim N(0, 1)$  (standard normal) it holds that  $\mathbb{P}[-2 < Z < 2] \approx 95\%$ , the typical condition number  
 1175 of a covariance matrix drawn with  $\sigma_0 = 5$  is roughly  $\frac{e^{2\sigma_0}}{e^{-2\sigma_0}} \approx 4.85 \cdot 10^8$ , which is a fairly large  
 1176 range.

## 1177 H.2 Details for Example 8: the linear cross-entropy setting

1178 In this setting,

$$\text{regret}(R, f \mid P_{\mathbf{x}}) = \min_{q \in \mathbb{R}^r} \mathbb{E} [D_{\text{KL}} ([1 + \exp(-f^\top \mathbf{x})]^{-1} \parallel [1 + \exp(-q^\top R^\top \mathbf{x})]^{-1})], \quad (190)$$

1179 and the expectation over the feature distribution typically cannot be carried out analytically. We thus  
 1180 tested Algorithm 1 on empirical distributions of samples drawn from a high-dimensional normal  
 1181 distribution. Specifically, for each run,  $B = 1000$  feature vectors were drawn from an isotropic  
 1182 normal distribution of dimension  $d = 15$ . The expectations of the regret and the corresponding  
 1183 gradients were then computed with respect to (w.r.t.) the resulting empirical distributions.

1184 **Regret gradients** We use the facts that

$$\frac{\partial}{\partial p_1} D_{\text{KL}}(p_1 \parallel p_2) = \log \frac{p_1(1 - p_2)}{p_2(1 - p_1)} \quad (191)$$

1185 and

$$\frac{\partial}{\partial p_2} D_{\text{KL}}(p_1 \parallel p_2) = \frac{p_2 - p_1}{p_2(1 - p_2)}. \quad (192)$$

1186 For brevity, let us next denote

$$p_1 := \frac{1}{1 + \exp(-f^\top \mathbf{x})} \quad (193)$$

1187 and

$$p_2 := \frac{1}{1 + \exp(-q^\top R^\top \mathbf{x})}. \quad (194)$$

1188 We next repeatedly use the chain rule for differentiation. First,

$$\nabla_f p_1 = \nabla_f \left[ \frac{1}{1 + \exp(-f^\top \mathbf{x})} \right] = \frac{\exp(-f^\top \mathbf{x}) \cdot \mathbf{x}}{[1 + \exp(-f^\top \mathbf{x})]^2} = p_1(1 - p_1) \cdot \mathbf{x}. \quad (195)$$

1189 and

$$\nabla_q p_2 = \nabla_q \left[ \frac{1}{1 + \exp(-q^\top R^\top \mathbf{x})} \right] = \frac{\exp(-q^\top R^\top \mathbf{x}) \cdot R^\top \mathbf{x}}{[1 + \exp(-q^\top R^\top \mathbf{x})]^2} = p_2(1 - p_2) \cdot R^\top \mathbf{x}. \quad (196)$$

1190 So, assuming that  $P_{\mathbf{x}}$  is such that the order of differentiation and expectation may be interchanged  
 1191 (this can be guaranteed using dominated/monotone convergence theorems), the gradient of the regret  
 1192 w.r.t.  $f$  is

$$\nabla_f \text{regret}(R, f \mid P_{\mathbf{x}}) = \mathbb{E} \left[ \frac{\partial}{\partial p_1} D_{\text{KL}}(p_1 \parallel p_2) \times \nabla_f p_1 \right] \quad (197)$$

$$= \mathbb{E} \left[ \log \left( \frac{p_1(1 - p_2)}{p_2(1 - p_1)} \right) \cdot p_1(1 - p_1) \cdot \mathbf{x} \right] \quad (198)$$

$$= \mathbb{E} \left[ (f^\top - q^\top R^\top) \mathbf{x} \frac{\exp(-f^\top \mathbf{x})}{[1 + \exp(-f^\top \mathbf{x})]^2} \cdot \mathbf{x} \right] \quad (199)$$

$$= \mathbb{E} \left[ \frac{\exp(-f^\top \mathbf{x})}{[1 + \exp(-f^\top \mathbf{x})]^2} \cdot \mathbf{x}^\top (f - Rq) \mathbf{x} \right]. \quad (200)$$

1193 Next, under similar assumptions, the gradient of the regret w.r.t. the predictor  $q$  is

$$\nabla_q \text{regret}(R, f \mid P_{\mathbf{x}}) = \mathbb{E} \left[ \frac{\partial}{\partial p_2} D_{\text{KL}}(p_1 \parallel p_2) \times \nabla_q p_2 \right] \quad (201)$$

$$= \mathbb{E} \left[ \left( \frac{1}{1 + \exp(-q^\top R^\top \mathbf{x})} - \frac{1}{1 + \exp(-f^\top \mathbf{x})} \right) \cdot R^\top \mathbf{x} \right]. \quad (202)$$

1194 Finally, as for the MSE case, to derive the gradient w.r.t.  $R$ , we denote  $R := [R_1, R_2, \dots, R_r] \in$   
 1195  $\mathbb{R}^{d \times r}$  where  $R_i \in \mathbb{R}^d$  is the  $i$ th column ( $i \in [r]$ ), and  $q^\top = (q_1, q_2, \dots, q_r)$ . Then,  $q^\top R^\top \mathbf{x} =$   
 1196  $\sum_{i \in [d]} q_i R_i^\top \mathbf{x}$  and

$$p_2 = \frac{1}{1 + \exp(-\sum_{i \in [d]} q_i R_i^\top \mathbf{x})}. \quad (203)$$

1197 Then, the gradient of  $p_2$  w.r.t.  $R_k$  is then given by

$$\nabla_{R_k} p_2 = p_2(1 - p_2) \cdot q_k \mathbf{x}, \quad (204)$$

1198 hence, more succinctly, the gradient w.r.t.  $R$  is

$$\nabla_R p_2 = p_2(1 - p_2) \cdot \mathbf{x} q^\top. \quad (205)$$

1199 Hence,

$$\nabla_R \text{regret}(R, f \mid P_{\mathbf{x}}) = \mathbb{E} \left[ \frac{\partial}{\partial p_2} D_{\text{KL}}(p_1 \parallel p_2) \times \nabla_R p_2 \right] \quad (206)$$

$$= \mathbb{E} [(p_2 - p_1) \cdot \mathbf{x} q^\top] \quad (207)$$

$$= \mathbb{E} \left[ \left( \frac{1}{1 + \exp(-q^\top R^\top \mathbf{x})} - \frac{1}{1 + \exp(-f^\top \mathbf{x})} \right) \cdot \mathbf{x} q^\top \right]. \quad (208)$$

Table 3: Parameters for linear cross entropy setting example

Parameter	$\beta_r$	$\beta_f$	$\eta_r$	$\eta_f$
Value	0.9	0.9	$10^{-3}$	$10^{-1}$
Parameter	$T_R$	$T_f$	$T_{\text{avg}}$	$T_{\text{stop}}$
Value	100	1000	25	50

1200 **Initialization** Here the initialization is similar to the linear MSE setting, except that since a column  
 1201 of the representation cannot ideally capture even a single adversarial function, the initialization  
 1202 algorithm only searches for a single adversarial function ( $\tilde{m} = 1$ ). This single function is then used  
 1203 to produce  $R^{(1)}$  as the initialization of Algorithm 1.

1204 **Algorithm parameters** The algorithm parameters used for Example 8 are shown in Table 3.

1205

## 1206 I An experiment with a NN architecture

1207 In the analysis and the experiments above we have considered basic linear functions. As mentioned,  
 1208 since the operation of Algorithm 1 only depends on the gradients of the loss function, it can be easily  
 1209 generalized to representations, response functions and predictors for which such gradients (or sub-  
 1210 gradients) can be provided. In this section, we exemplify this idea with a simple NN architecture.  
 1211 For  $x \in \mathbb{R}^d$ , we let the rectifier linear unit (ReLU) be denoted as  $(x)_+$ .

1212 **Definition 22** (The NN setting). Assume the same setting as in Definitions 1 and 7, except that the  
 1213 class of representation, response and predictors are NN with  $c$  hidden layers of sizes  $h_R, h_f, h_q \in$   
 1214  $\mathbb{N}_+$ , respectively, instead of linear functions. Specifically: (1) The representation is

$$R(x) = R_c^\top \left( \cdots \left( R_1^\top (R_0^\top x)_+ \right)_+ \right)_+ \quad (209)$$

1215 for some  $(R_0, R_1, \dots, R_c) \in \mathcal{R} := \{\mathbb{R}^{d \times h_R} \times \mathbb{R}^{h_R \times h_R} \dots \mathbb{R}^{h_R \times h_R} \times \mathbb{R}^{h_R \times r}\}$  where  $d > r$ . (2)  
 1216 The response is determined by

$$f(x) = f_c^\top \left( \cdots \left( F_1^\top (F_0^\top x)_+ \right)_+ \right)_+ \quad (210)$$

1217 where  $(F_0, F_1, \dots, f_c) \in \mathcal{F} := \{\mathbb{R}^{d \times h_f} \times \mathbb{R}^{h_f \times h_f} \dots \mathbb{R}^{h_f \times h_f} \times \mathbb{R}^{h_f}\}$ . (3) The predictor is deter-  
 1218 mined by for some

$$q(z) = q_c^\top \left( \cdots \left( Q_1^\top (Q_0^\top z)_+ \right)_+ \right)_+ \quad (211)$$

1219 where  $(Q_0, Q_1, \dots, q_c) \in \mathcal{Q} := \{\mathbb{R}^{r \times h_q} \times \mathbb{R}^{h_q \times h_q} \dots \mathbb{R}^{h_q \times h_q} \times \mathbb{R}^{h_q}\}$ .

1220 **Regret gradients** Gradients were computed using PyTorch with standard gradients computation  
 1221 using backpropagation for an SGD optimizer.

1222 **Initialization** The initialization algorithm is similar to the initialization algorithm used in the lin-  
 1223 ear cross-entropy setting.

1224 **Algorithm parameters** The algorithm parameters used for the example are shown in Table 4.

1225 **Results** For a single hidden layer, Figure 4 shows the reduction of the regret with the iteration for  
 1226 the cross-entropy loss.

Table 4: Parameters for the NN cross-entropy setting.

Parameter	$c$	$h_R$	$h_f$	$h_q$	
Value	1	$d$	$d$	$d$	
Parameter	$\beta_r$	$\beta_f$	$\eta_r$	$\eta_f$	$\eta_q$
Value	0.9	0.9	$10^{-3}$	$10^{-1}$	$10^{-1}$
Parameter	$T_R$	$T_f$	$T_Q$	$T_{\text{avg}}$	$T_{\text{stop}}$
Value	100	1000	100	10	80

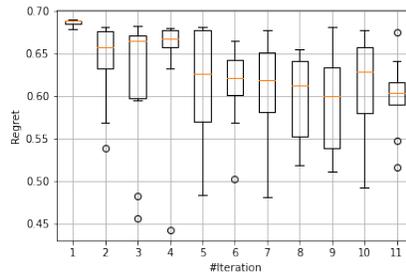


Figure 4: The regret achieved by Algorithm 1 in the NN cross-entropy setting as a function of the iteration  $m$ .