

---

# Isometry pursuit

---

**Samson Koelle**  
Amazon Web Services  
Seattle, WA

**Marina Meila**  
Department of Statistics  
University of Washington  
Seattle, WA

## Abstract

Isometry pursuit is a convex algorithm for identifying orthonormal column-submatrices of wide matrices. It consists of a novel normalization method followed by multitask basis pursuit. Applied to Jacobians of putative coordinate functions, it helps identify isometric embeddings from within interpretable dictionaries. We provide theoretical and experimental results justifying this method. It appears to be more accurate than greedy search and more efficient than brute force search.

## 1 Introduction

Many real-world problems may be abstracted as selecting a subset of the columns of a matrix representing stochastic observations or analytically exact data. This paper focuses on a simple such problem that appears in interpretable learning. Given a rank  $D$  matrix  $X \in \mathbb{R}^{D \times P}$  with  $P > D$ , select a square submatrix  $X_{\cdot S}$  where subset  $S \subset P$  satisfies  $|S| = D$  that is as orthonormal as possible.

This problem arises in interpretable learning because while the coordinate functions of a learned latent space may have no intrinsic meaning, it is sometimes possible to generate a dictionary of interpretable features which may be considered as potential parametrizing coordinates. When this is the case, selection of candidate interpretable features as coordinates can take the above form. While implementations vary across data and algorithmic domains, identification of such coordinates generally aids mechanistic understanding, generative control, and statistical efficiency.

This paper shows that an adapted version of the algorithm in Koelle et al. [1] leads to a convex procedure that helps improve upon greedy approaches such as those in Cai and Wang [2], Chen and Meila [3], Kohli et al. [4], Jones et al. [5] for finding isometries. The insight that leads to isometry pursuit is that  $D$  function solutions multitask basis pursuit applied to an appropriately normalized  $X$  selects orthonormal submatrices. In particular, the normalization log-symmetrizes length in the column-space of  $X$  and favors vectors closer to unit length, while basis pursuit favors vectors which are orthogonal. Our theoretical results formalize this intuition within a limited setting, while our experimental results show the usefulness of isometry pursuit as a trimming procedure prior to brute force search. Additionally, we introduce a novel ground truth objective function that we measure the success of our algorithm against.

## 2 Background

Our algorithm is motivated by spectral and convex analysis.

---

<sup>1</sup>Work conducted outside of Amazon.

<sup>2</sup>Code is available at <https://github.com/sjkoelle/convexlocalisometry>.

## 2.1 Problem

Our goal is, given a matrix  $X \in \mathbb{R}^{D \times P}$ , select a subset  $S \subset [P]$  with  $|S| = D$  such that  $X_{\cdot S}$  is as orthonormal as possible in a computationally efficient way. To that end, we define a ground truth loss function that measures orthonormalness, and then introduce a surrogate loss function that convexifies the problem so that it may be efficiently solved.

## 2.2 Interpretability and isometry

Our motivating example is the selection of data representations from within sets of putative coordinates. These putative coordinates are simply the columns of a provided wide matrix. The proposed method is thus even simpler than Sparse PCA [6, 7, 8], in which column-covariance is used to select low-dimensional projections from within the span of such a subset.

This method is specifically applicable with respect to interpretability, for which parsimony is at a premium. Interpretability arises through comparison of data with what is known to be important in the domain of the problem. This a priori knowledge often takes the form of a functional dictionary. Regardless of implementation details such as whether this dictionary is given or learned, core concepts like evaluation of independentness of dictionary features arise in numerous scenarios [9, 10, 11]. After functional independence [10], also known as feature decomposability [12], which only requires that the differential of sets of dictionary features be full rank, metric properties of such sets are of natural interest.

**Definition 1** The *differential* of a smooth map  $\phi : \mathcal{M} \rightarrow \mathcal{N}$  between  $D$  dimensional manifolds  $\mathcal{M} \subseteq \mathbb{R}^B$  and  $\mathcal{N} \subseteq \mathbb{R}^P$  is a map in tangent bases  $x_1 \dots x_D$  of  $T_\xi \mathcal{M}$  and  $y_1 \dots y_D$  of  $T_{\phi(\xi)} \mathcal{N}$  consisting of entries

$$D\phi(\xi) = \begin{bmatrix} \frac{\partial \phi_1}{\partial x_1}(\xi) & \dots & \frac{\partial \phi_1}{\partial x_D}(\xi) \\ \vdots & & \vdots \\ \frac{\partial \phi_D}{\partial x_1}(\xi) & \dots & \frac{\partial \phi_D}{\partial x_D}(\xi) \end{bmatrix}. \quad (1)$$

**Definition 2** A map  $\phi$  between  $D$  dimensional submanifolds with inherited Euclidean metric  $\mathcal{M} \subseteq \mathbb{R}^{B_\alpha}$  and  $\mathcal{N} \subseteq \mathbb{R}^{B_\beta}$  is  $\phi$  is an **isometry at a point**  $\xi \in \mathcal{M}$  if

$$D\phi(\xi)^T D\phi(\xi) = I_D. \quad (2)$$

That is,  $\phi$  is an isometry at  $\xi$  if  $D\phi(\xi)$  is orthonormal.

This property - that  $D\phi$  is orthonormal, has several equivalent formulations. The formulation that motivates our ground truth loss function comes from spectral analysis.

**Proposition 1** The singular values  $\sigma_1 \dots \sigma_D$  are equal to 1 if and only if  $U \in \mathbb{R}^{D \times D}$  is orthonormal.

On the other hand, the formulation that motivates our convex approach is that orthonormal matrices consist of  $D$  coordinate features whose gradients are orthogonal and evenly varying.

**Proposition 2** The component vectors  $u_1 \dots u_D \in \mathbb{R}^B$  form a orthonormal matrix if and only if, for all  $d_1, d_2 \in [D]$ ,  $\langle u_{d_1}, u_{d_2} \rangle = \begin{cases} 1 & d_1 = d_2 \\ 0 & d_1 \neq d_2 \end{cases}$ .

The applications of pointwise isometry are themselves manifold. Local Tangent Space Alignment [13], Multidimensional Scaling [14] and Isomap [15] non-parametrically estimate embeddings that are as isometric as possible. Pointwise isometries selected from a dictionary may be stitched together to form global embeddings [4]. This approach is particularly relevant since it constructs such isometries through greedy search, with putative dictionary features added one at a time.

Note that it is not necessary to explicitly estimate tangent spaces when applying the definition of isometry. The most commonly encountered manifolds are simply vector spaces, in which case the tangent spaces are trivial. This is the case for full-rank tabular data, as well as latent spaces of deep learning models. For example, the transformer residual stream at different tokens are analogous to tangent spaces of a non-linear manifold in the sense that the relative directions of dictionary vectors are not consistent between tokens.

### 2.3 Subset selection

Given a matrix  $X \in \mathbb{R}^{D \times P}$ , we compare algorithmic paradigms for solving problems of the form

$$\arg \min_{S \in \binom{[P]}{d}} l(X_S) \quad (3)$$

where  $\binom{[P]}{d} = \{A \subseteq [P] : |A| = d\}$ . Brute force algorithms consider all possible solutions. These algorithms are conceptually simple, but have the often prohibitive time complexity  $O(C_l P^D)$  where  $C_l$  is the cost of evaluating  $l$ . Greedy algorithms consist of iteratively adding one element at a time to  $S$ . These algorithms have time complexity  $O(C_l P D)$  and so are computationally more efficient than brute force algorithms, but can get stuck in local minima. Please see Section 6.1 for additional information.

Sometimes, it is possible to introduce an objective which convexifies problems of the above form. Solutions

$$\arg \min f(\beta) : Y = X\beta \quad (4)$$

to the overcomplete regression problem  $Y = X\beta$  are a classic example [16]. When  $f(\beta) = \|\beta\|_0$ , this problem is non-convex, and must be solved via greedy or brute algorithms, but when  $f(\beta) = \|\beta\|_1$ , the problem is convex, and may be solved efficiently via interior-point methods. When the equality constraint is relaxed, Lagrangian duality may be used to reformulate as a so-called Lasso problem, which leads to an even richer set of optimization algorithms.

The form of basis pursuit that we apply is inspired by the group basis pursuit approach in Koelle et al. [10]. In group basis pursuit (which we call multitask basis pursuit when grouping is dependent only on the structure of matrix-valued response variable  $y$ ) the objective function is  $f(\beta) = \|\beta\|_{1,2} := \sum_{p=1}^P \|\beta_p\|_2$  [17, 18, 19] This objective creates joint sparsity across entire rows of  $\beta_p$ , and was used in [10] to select between sets of interpretable features.

## 3 Method

We adapt the group lasso paradigm used to select independent dictionary elements in Koelle et al. [10, 1] to select pointwise isometries from a dictionary. In particular, we first will define a ground truth objective computable via brute and greedy algorithms that is uniquely minimized by orthonormal matrices. We then define the combination of normalization and multitask basis pursuit that approximates this ground truth loss function. We finally give a brute post-processing method for ensuring that the solution is  $D$  sparse.

### 3.1 Ground truth

We'd like a ground truth objective to be minimized uniquely by orthonormal matrices, invariant under rotation, and depend on all changes in the matrix. Deformation [4] and nuclear norm [20] use only a subset of the differential's information and are not uniquely minimized at unitarity, respectively. We therefore introduce an alternative ground truth objective that satisfies the above desiderata and has convenient connections to Isometry Pursuit.

We define this objective as

$$l_c : \mathbb{R}^{D \times P} \rightarrow \mathbb{R}^+ \quad (5)$$

$$X \mapsto \sum_{d=1}^D g(\sigma_d(X), c) \quad (6)$$

where  $\sigma_d(X)$  is the  $d$ -th singular value of  $X$  and

$$g : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}^+ \quad (7)$$

$$t, c \mapsto \frac{e^{t^c} + e^{t^{-c}}}{2e}. \quad (8)$$

Using Proposition 1, we can check that  $g$  is uniquely maximized by orthonormal matrices. Moreover,  $g(X^{-1}) = g(X)$  and  $g$  is convex. Figure 1 gives a graph of  $l_c$  when  $D = 1$  and compares it with that produced by basis pursuit after normalization as in Section 3.2.

Our ground truth program is therefore

$$\hat{S}_{GT} = \arg \min_{S \in \binom{[P]}{d}} l_c(X.S). \quad (9)$$

Regardless of the convexity of  $l_c$ , brute combinatorial search over  $[P]$  is inherently non-convex.

### 3.2 Normalization

Since basis pursuit methods tend to select longer vectors, selection of orthonormal submatrices requires normalization such that both long and short candidate basis vectors are penalized in the subsequent regression. We introduce the following definition.

**Definition 3 (Symmetric normalization)** *A function  $q : \mathbb{R}^D \rightarrow \mathbb{R}^+$  is a symmetric normalization if*

$$\arg \max_{v \in \mathbb{R}^D} q(v) = \{v : \|v\|_2 = 1\} \quad (10)$$

$$q(v) = q\left(\frac{v}{\|v\|_2}\right) \quad (11)$$

$$q(v_1) = q(v_2) \quad \forall v_1, v_2 \in \mathbb{R}^D : \|v_1\|_2 = \|v_2\|_2. \quad (12)$$

We use such functions to normalize vector length in such a way that vectors of length 1 prior to normalization have longest length after normalization, vectors in general are shrunk proportionately to their deviation from 1. That is, we normalize vectors by

$$n : \mathbb{R}^D \rightarrow \mathbb{R}^D \quad (13)$$

$$v \mapsto q(v)v \quad (14)$$

and matrices by

$$w : \mathbb{R}^{D \times P} \rightarrow \mathbb{R}^D \quad (15)$$

$$X_{:,p} \mapsto n(X_{:,p}) \quad \forall p \in [P]. \quad (16)$$

In particular, given  $c > 0$ , we choose  $q$  as follows.

$$q_c : \mathbb{R}^D \rightarrow \mathbb{R}^+ \quad (17)$$

$$v \mapsto \frac{e^{\|v\|_2^c} + e^{\|v\|_2^{-c}}}{2e}. \quad (18)$$

Besides satisfying the conditions in Definition 3, this normalization has some additional nice properties. First,  $q$  is convex. Second, it grows asymptotically log-linearly. Third, while  $\exp(-|\log t|) = \exp(-\max(t, 1/t))$  is a seemingly natural choice for normalization, it is non smooth, and the LogSumExp replacement of  $\max(t, 1/t)$  with  $\log(\exp(t) + \exp(1/t))$  simplifies to 17 upon exponentiation [20]. Finally, the parameter  $c$  grants control over the width of the basin, which is important in avoiding numerical issues arising close to 0 and  $\infty$ .

### 3.3 Isometry pursuit

Isometry pursuit is the application of multitask basis pursuit to the normalized design matrix  $w(X, c)$  to identify submatrices of  $X$  that are as orthonormal as possible. Define the multitask basis pursuit penalty

$$\|\cdot\|_{1,2} : \mathbb{R}^{P \times D} \rightarrow \mathbb{R}^+ \quad (19)$$

$$\beta \mapsto \sum_{p=1}^P \|\beta_p\|_2. \quad (20)$$

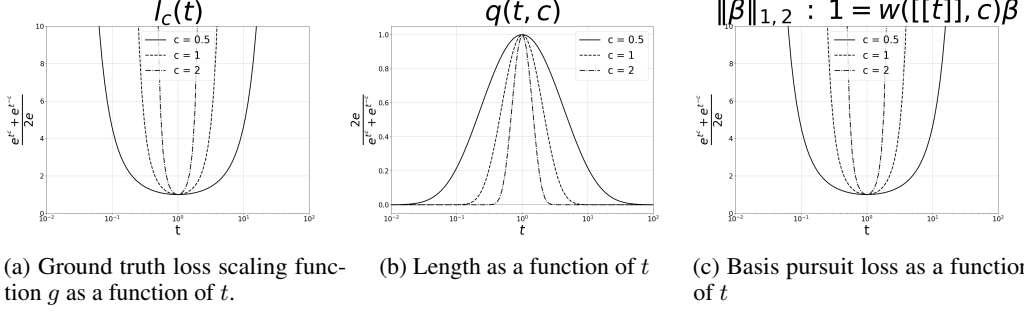


Figure 1: Plots of ground truth loss, normalized length, and basis pursuit loss for different values of  $c$  in the one-dimensional case  $D = 1$ . The two losses are equivalent in the one-dimensional case.

Given a matrix  $Y \in \mathbb{R}^{D \times E}$ , the multitask basis pursuit solution is

$$\hat{\beta}_{MBP}(X, Y) := \arg \min_{\beta \in \mathbb{R}^{P \times E}} \|\beta\|_{1,2} : Y = X\beta. \quad (21)$$

Isometry pursuit is then given by

$$\hat{\beta}_c(X) := \hat{\beta}_{MBP}(w(X, c), I_D) \quad (22)$$

where  $I_D$  is the  $D$  dimensional identity matrix and recovered functions are the indices of the dictionary elements with non-zero coefficients. That is, they are given by  $S(\beta)$  where

$$S : \mathbb{R}^{p \times d} \rightarrow \binom{[P]}{d} \quad (23)$$

$$\beta \mapsto \{p \in [P] : \|\beta_p\| > 0\}. \quad (24)$$

---

ISOMETRYPURSUIT(Matrix  $X \in \mathbb{R}^{D \times P}$ , scaling constant  $c$ )

---

- 1: Normalize  $X_c = w(X, c)$
  - 2: Optimize  $\hat{\beta} = \hat{\beta}_{MBP}(X_c, I_D)$
  - 3: **Output**  $\hat{S} = S(\hat{\beta})$
- 

### 3.4 Theory

The intuition behind our application of multitask basis pursuit in our setting is that submatrices consisting of vectors which are closer to 1 in length and more orthogonal will have smaller loss. A key initial theoretical assertion is that ISOMETRYPURSUIT is invariant to choice of basis for  $X$ .

**Proposition 3** Let  $U \in \mathbb{R}^{D \times D}$  be orthonormal. Then  $S(\hat{\beta}(UX)) = S(\hat{\beta}(X))$ .

A proof is given in Section 6.2.1. This has as an immediate corollary that we may replace  $I_D$  in the constraint by any orthonormal  $D \times D$  matrix.

We also claim that the conditions of the consequent of Proposition 2 are satisfied by minimizers of the multitask basis pursuit objective applied to suitably normalized matrices in the special case where both such a submatrix exists and  $|S| = D$ .

**Proposition 4** Let  $w_c$  be a normalization satisfying the conditions in Definition 3. Then  $\arg \min_{X, S \in \mathbb{R}^{D \times D}} \hat{\beta}_c^D(X)$  is orthonormal. Moreover when  $X$  is orthonormal,  $(\min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : I_D = w(\hat{X}, c)\beta) = D$ .

While this Proposition falls short of showing that an orthonormal submatrix will be selected should one be present, it provides intuition justifying the preferential efficacy of ISOMETRYPURSUIT on real data. A proof is given in Section 6.2.2.

### 3.5 Two-stage isometry pursuit

Since cannot in general ensure either that  $|\widehat{S}| = D$  or that a orthonormal submatrix  $X_{\widehat{S}}$  exists, we first use the convex problem to prune and then apply brute search upon the substantially reduced feature set.

---

TWOSTAGEISOMETRYPURSUIT(Matrix  $X \in \mathbb{R}^{D \times P}$ , scaling constant  $c$ )

---

- 1:  $\widehat{S}_1 = \text{ISOMETRYPURSUIT}(X, c)$
  - 2:  $\widehat{S} = \text{BRUTESEARCH}(X_{\widehat{S}_1}, l_c)$
  - 3: **Output**  $\widehat{S}$
- 

Similar two-stage approaches are standard in the Lasso literature [21].

## 4 Experiments

We compare TWOSTAGEISOMETRYPURSUIT and GREEDYSEARCH on the Iris and Wine datasets [22, 23, 24], as well as the Ethanol dataset from Chmiela et al. [25], Koelle et al. [10]. For the latter, a dictionary of interpretable features are evaluated for their ability to parameterize the data manifold through computation of their Jacobian matrices and projection onto estimated tangent spaces (see Koelle et al. [10] for preprocessing details). Statistical replicas for Wine and Iris are created by resampling across  $P$ , while for Ethanol they are created by sampling from data point and their corresponding tangent spaces. For basis pursuit, we use the SCS interior point solver [26] from CVXPY [27, 28], which is able to push sparse values arbitrarily close to 0 [29]. Table 1 presents results showing that the  $l_c$  accrued by the subset  $\widehat{S}_G$  estimated using GREEDYSEARCH with objective  $l_c$  is higher than that for the subset estimated by TWOSTAGEISOMETRYPURSUIT. We also evaluated second stage BRUTESEARCH selection after random selection of  $\widehat{S}_1$  but do not report it since it often lead to catastrophic failure to satisfy the basis pursuit constraint.

Name	$D$	$P$	$R$	$l_1(X_{\widehat{S}_G})$	$ \widehat{S}_1 $	$l_1(X_{\widehat{S}})$	$P(\bar{l}_1(X_{\widehat{S}_G}) > \bar{l}_1(X_{\widehat{S}}))$
Iris	4	75	25	$13.4 \pm 6.4$	$7 \pm 1$	$8.0 \pm 1.8$	$10^{-4}$
Wine	5	89	25	$5.7 \pm .2$	$12 \pm 1$	$5.6 \pm .1$	$5 \times 10^{-5}$
Ethanol	2	756	100	$2.6 \pm .3$	$90 \pm 164$	$2.5 \pm .2$	$2 \times 10^{-5}$

Table 1: Experimental parameters and results. For the Wine dataset, even BRUTESEARCH on  $\widehat{S}_1$  is prohibitive in  $D = 13$ , and so we truncate our inputs to  $D = 5$ . For Iris and Wine,  $P$  is randomly downsampled by a factor of 2 to create replicates. P-values are computed by paired two-sample T-test on  $l_1(X_{\widehat{S}})$  and  $l_1(X_{\widehat{S}_G})$ .

## 5 Discussion

We have shown the efficacy of a convex multitask basis pursuit approach for selecting isometric submatrices from wide-matrices. This approach - ISOMETRYPURSUIT - should be considered as an alternative to greedy methods for selection of such features from within a dictionary. In a sense it is a simpler simpler version of Sparse PCA.

Algorithmic variants of interest include the multitask lasso extension of our estimator. Applications of interest include greedy diversification in recommendation systems [30, 31] and other retrieval-type systems, as well as decomposing interpretable yet overcomplete dictionaries in transformer residual streams. The most pressing piece of remaining theoretical work removal of the restriction  $|S| = D$  on the conditions of Proposition 4 and investigating of the necessity of the second selection step. The resulting proposition is intuitive but is more difficult to prove and seemingly violated by (omitted) empirical results that are subtly non-dispositive in the sense that improvements of primal loss below  $D$  are accompanied by violations of the constraint of a similar magnitude. We also note that isometries may not always exist in the presence of curvature, so comparison of our loss with curvature could prove fertile, as could comparison with the so-called restricted isometry property used to show guaranteed recovery at fast convergence rates in supervised learning [32, 33].

## References

- [1] Samson J Koelle, Hanyu Zhang, Octavian-Vlad Murad, and Marina Meila. Consistency of dictionary-based manifold learning. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 4348–4356. PMLR, 2024.
- [2] T. Tony Cai and Lie Wang. Orthogonal matching pursuit for sparse signal recovery with noise. *IEEE Transactions on Information Theory*, 57(7):4680–4688, 2011. doi: 10.1109/TIT.2011.2146090.
- [3] Yu-Chia Chen and Marina Meila. Selecting the independent coordinates of manifolds with large aspect ratios. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/6a10bbd480e4c5573d8f3af73ae0454b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/6a10bbd480e4c5573d8f3af73ae0454b-Paper.pdf).
- [4] Dhruv Kohli, Alexander Cloninger, and Gal Mishne. LDLE: Low distortion local eigenmaps. *J. Mach. Learn. Res.*, 22, 2021.
- [5] Peter W Jones, Mauro Maggioni, and Raanan Schul. Universal local parametrizations via heat kernels and eigenfunctions of the laplacian. September 2007.
- [6] Santanu S Dey, R Mazumder, M Molinaro, and Guanyi Wang. Sparse principal component analysis and its  $l_1$ -relaxation. *arXiv: Optimization and Control*, December 2017.
- [7] D Bertsimas and Driss Lahlou Kitane. Sparse PCA: A geometric approach. *J. Mach. Learn. Res.*, 24:32:1–32:33, October 2022.
- [8] Dimitris Bertsimas, Ryan Cory-Wright, and Jean Pauphilet. Solving Large-Scale sparse PCA to certifiable (near) optimality. *J. Mach. Learn. Res.*, 23(13):1–35, 2022.
- [9] Yu-Chia Chen and M Meilă. Selecting the independent coordinates of manifolds with large aspect ratios. *Adv. Neural Inf. Process. Syst.*, abs/1907.01651, July 2019.
- [10] Samson J Koelle, Hanyu Zhang, Marina Meila, and Yu-Chia Chen. Manifold coordinates with physical meaning. *J. Mach. Learn. Res.*, 23(133):1–57, 2022.
- [11] Jesse He, Tristan Brugère, and Gal Mishne. Product manifold learning with independent coordinate selection. In *Proceedings of the 2nd Annual Workshop on Topology, Algebra, and Geometry in Machine Learning (TAG-ML) at ICML*, June 2023.
- [12] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermy, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- [13] Zhenyue Zhang and Hongyuan Zha. Principal manifolds and nonlinear dimensionality reduction via tangent space alignment. *SIAM J. Scientific Computing*, 26(1):313–338, 2004.
- [14] Lisha Chen and Andreas Buja. Local Multidimensional Scaling for nonlinear dimension reduction, graph drawing and proximity analysis. *Journal of the American Statistical Association*, 104(485):209–219, March 2009.
- [15] J.B. Tenenbaum, V. Silva, and J.C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [16] Scott Shaobing Chen and David L. Donoho and Michael A. Saunders. Atomic decomposition by basis pursuit. *SIAM REVIEW*, 43(1):129, February 2001.

- [17] Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. Series B Stat. Methodol.*, 68(1):49–67, February 2006.
- [18] G Obozinski, B Taskar, and Michael I Jordan. Multi-task feature selection. 2006.
- [19] Dit-Yan Yeung and Yu Zhang. A probabilistic framework for learning task relationships in multi-task learning. 2011.
- [20] Stephen P Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, March 2004.
- [21] Tim Hesterberg, Nam Hee Choi, Lukas Meier, and Chris Fraley. Least angle and  $\ell_1$  penalized regression: A review. February 2008.
- [22] R. A. Fisher. Iris. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C56C76>.
- [23] Stefan Aeberhard and M. Forina. Wine. UCI Machine Learning Repository, 1991. DOI: <https://doi.org/10.24432/C5PC7J>.
- [24] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [25] Stefan Chmiela, Huziel E Sauceda, Klaus-Robert Müller, and Alexandre Tkatchenko. Towards exact molecular dynamics simulations with machine-learned force fields. *Nat. Commun.*, 9(1): 3887, September 2018.
- [26] Brendan O’Donoghue, Eric Chu, Neal Parikh, and Stephen Boyd. Conic optimization via operator splitting and homogeneous self-dual embedding. *Journal of Optimization Theory and Applications*, 169(3):1042–1068, June 2016. URL <http://stanford.edu/~boyd/papers/scs.html>.
- [27] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [28] Akshay Agrawal, Robin Verschueren, Steven Diamond, and Stephen Boyd. A rewriting system for convex optimization problems. *Journal of Control and Decision*, 5(1):42–60, 2018.
- [29] CVXPY Developers. Sparse solution with cvxpy. [https://www.cvxpy.org/examples/applications/sparse\\_solution.html](https://www.cvxpy.org/examples/applications/sparse_solution.html). Accessed: 2024-07-11.
- [30] Jaime Carbonell and Jade Goldstein. The use of MMR, diversity-based reranking for reordering documents and producing summaries. *SIGIR Forum*, 51(2):209–210, August 1998.
- [31] Qiong Wu, Yong Liu, Chunyan Miao, Yin Zhao, Lu Guan, and Haihong Tang. Recent advances in diversified recommendation. May 2019.
- [32] Emmanuel Candes and Terence Tao. Decoding by linear programming. February 2005.
- [33] T Hastie, R Tibshirani, and M Wainwright. Statistical learning with sparsity: The lasso and generalizations. May 2015.
- [34] E Anderson, Z Bai, and J Dongarra. Generalized qr factorization and its applications. *Linear Algebra Appl.*, 162-164:243–271, February 1992.

## 6 Supplement

This section contains algorithms and proofs in support of the main text.



## 6.1 Algorithms

We give definitions of the brute and greedy algorithms used in this paper.

---

**BRUTESearch**(Matrix  $X \in \mathbb{R}^{D \times P}$ , objective  $f$ )

---

```

1: for each combination  $S \subseteq \{1, 2, \dots, P\}$  with  $|S| = D$  do
2:   Evaluate  $f(X_{.S})$ 
3: end for
4: Output the combination  $S^*$  that minimizes  $f(X_{.S})$ 

```

---



---

**GREEDYSEARCH**(Matrix  $X \in \mathbb{R}^{D \times P}$ , objective  $f$ , selected set  $S = \emptyset$ , current size  $d = 0$ )

---

```

1: if  $d = D$  then
2:   Return  $S$ 
3: else
4:   Initialize  $S_{\text{best}} = S$ 
5:   Initialize  $f_{\text{best}} = \infty$ 
6:   for each  $p \in \{1, 2, \dots, P\} \setminus S$  do
7:     Evaluate  $f(X_{.(S \cup \{p\})})$ 
8:     if  $f(X_{.(S \cup \{p\})}) < f_{\text{best}}$  then
9:       Update  $S_{\text{best}} = S \cup \{p\}$ 
10:      Update  $f_{\text{best}} = f(X_{.(S \cup \{p\})})$ 
11:    end if
12:  end for
13:  Return GREEDYSEARCH( $X, f, S_{\text{best}}, d + 1$ )
14: end if

```

---

## 6.2 Proofs

### 6.2.1 Proof of Proposition 3

In this proof we first show that the penalty  $\|\beta\|_{1,2}$  is unchanged by unitary transformation of  $\beta$ .

**Proposition 5** *Let  $U \in \mathbb{R}^{D \times D}$  be unitary. Then  $\|\beta\|_{1,2} = \|\beta U\|$ .*

**Proof:**

$$\|\beta U\|_{1,2} = \sum_{p=1}^P \|\beta_p \cdot U\| \quad (25)$$

$$= \sum_{p=1}^P \|\beta_p \cdot\| \quad (26)$$

$$= \|\beta\|_{1,2} \quad (27)$$

□

We then show that this implies that the resultant loss is unchanged by unitary transformation of  $X$ .

**Proposition 6** *Let  $U \in \mathbb{R}^{D \times D}$  be unitary. Then  $\widehat{\beta}(UX) = \widehat{\beta}(X)U$ .*

**Proof:**

$$\hat{\beta}(UX) = \arg \min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : I_D = UX\beta \quad (28)$$

$$= \arg \min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : U^{-1}U = U^{-1}UX\beta U \quad (29)$$

$$= \arg \min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : I_D = X\beta U \quad (30)$$

$$= \arg \min_{\beta \in \mathbb{R}^{P \times D}} \|\beta U\|_{1,2} : I_D = X\beta U \quad (31)$$

$$= \arg \min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : I_D = X\beta. \quad (32)$$

□

### 6.2.2 Proof of Proposition 4

**Proposition 7** Let  $w_c$  be a normalization satisfying the conditions in Definition 3. Then  $\arg \min_{X, S \in \mathbb{R}^{D \times D}} \hat{\beta}_c^D(X, S)$  is orthonormal. Moreover when  $X$  is orthonormal,  $\min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : I_D = w(X, c)\beta = D$ .

**Proof:** The value of  $D$  is clearly obtained by  $\beta$  orthonormal, since by Proposition 3, for  $X$  orthogonal, without loss of generality

$$\beta_{dd'} = \begin{cases} 1 & d = d' \in \{1 \dots D\} \\ 0 & \text{otherwise} \end{cases}. \quad (33)$$

Thus, we need to show that this is a lower bound on the obtained loss.

From the conditions in Definition 3, normalized matrices will consist of vectors of maximum length (i.e. 1) if and only if the original matrix also consists of vectors of length 1. Such vectors will clearly result in lower basis pursuit loss, since longer vectors in  $X$  require smaller corresponding covectors in  $\beta$  to equal the same result.

Therefore, it remains to show that  $X$  consisting of orthogonal vectors of length 1 have lower compared with  $X$  consisting of non-orthogonal vectors. Invertible matrices  $X_{\cdot S}$  admit QR decompositions  $\tilde{X}_{\cdot S} = QR$  where  $Q$  and  $R$  are orthonormal and upper-triangular matrices, respectively [34]. Denoting  $Q$  to be composed of basis vectors  $[e^1 \dots e^d]$ , the matrix  $R$  has form

$$R = \begin{bmatrix} \langle e^1, X_{\cdot S_1} \rangle & \langle e^1, X_{\cdot S_2} \rangle & \dots & \langle e^1, X_{\cdot S_D} \rangle \\ 0 & \langle e^2, X_{\cdot S_2} \rangle & \dots & \langle e^2, X_{\cdot S_D} \rangle \\ 0 & 0 & \dots & \dots \\ 0 & 0 & \dots & \langle e^d, X_{\cdot S_D} \rangle \end{bmatrix}. \quad (34)$$

Thus,  $|R_{dd}| \leq \|X_{\cdot S_d}\|_2$ , with equality obtained across  $d$  only by orthonormal matrices. On the other hand, by Proposition 3,  $l(X) = l(R)$  and so  $\|\beta\|_{1,2} = \|R^{-1}\|_{1,2}$ . Since  $R$  is upper triangular it has diagonal elements  $\beta_{dd} = R_{dd}^{-1}$  and so  $\|\beta_d\| \geq \|X_{\cdot S_d}\|^{-1} = 1$ . That is, the penalty accrued by a particular covector in  $\beta$  is bounded from below by 1 - the inverse of the length of the corresponding vector in  $X_{\cdot S}$  - with equality occurring only when  $X_{\cdot S}$  is orthonormal.

□