
Isometry pursuit

Anonymous Author(s)

Affiliation

Address

email

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 $X_{\mathcal{S}}$ where subset $\mathcal{S} \subset P$ satisfies $|\mathcal{S}| = D$ 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 found in Chen and Meila [2], Kohli et al. [3], Jones et al. [4] 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.

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 $X_{\mathcal{S}}$ is as orthonormal as possible in a computationally efficient way. To that end, 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 [5, 6, 7], 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 [8, 9, 10]. After functional independence [9], also known as feature decomposability [11], 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 ϕ between D dimensional submanifolds with inherited Euclidean metric $\mathcal{M} \subseteq \mathbb{R}^{B_\alpha}$ and $\mathcal{N} \subseteq \mathbb{R}^{B_\beta}$ is ϕ is an **isometry at a point** $\xi \in \mathcal{M}$ if

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

That is, ϕ is an isometry at ξ 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 Isometry Pursuit 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, Multidimensional Scaling and Isomap non-parametrically estimate embeddings that are as isometric as possible. Pointwise isometries selected from a dictionary may be stitched together to form global embeddings [3]. This approach constructs 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

$$\arg \min_{S \subseteq [P]: |S|=D} l(X, S) \quad (3)$$

75 Brute force algorithms consider all possible solutions. These algorithms are conceptually simple,
 76 but often have prohibitive time complexity $O(C_l P^D)$ where C_l is the cost of evaluating l . Greedy
 77 algorithms consist of iteratively adding one element at a time to S . These algorithms have time
 78 complexity $O(CPD)$ and so are computationally more efficient than brute force algorithms, but can
 79 get stuck in local minima. Please see Section 6.1 for additional information.

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

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

82 to the overcomplete regression problem $Y = X\beta$ are a classic example [12]. When $f(\beta) = \|\beta\|_0$, this
 83 problem is non-convex, and must be solved via greedy or brute algorithms, but when $f(\beta) = \|\beta\|_1$,
 84 the problem is convex, and may be solved efficiently via interior-point methods. When the equality
 85 constraint is relaxed, Lagrangian duality may be used to reformulate as a so-called Lasso problem,
 86 which leads to an even richer set of optimization algorithms.

87 The particular form of basis pursuit that we apply is inspired by the group basis pursuit approach
 88 in Koelle et al. [9]. In group basis pursuit (which we call multitask basis pursuit when grouping
 89 is dependent only on the structure of matrix-valued response variable y) the objective function is
 90 $f(\beta) = \|\beta\|_{1,2} := \sum_{p=1}^P \|\beta_p\|_2$ [13, 14, 15] This objective creates joint sparsity across entire rows
 91 of β_p , and was used in [9] to select between sets of interpretable features.

92 3 Method

93 We adapt the group lasso paradigm used to select independent dictionary elements in Koelle et al.
 94 [9, 1] to select pointwise isometries from a dictionary. In particular, we first will define a ground truth
 95 objective computable via brute and greedy algorithms that is uniquely minimized by orthonormal
 96 matrices. We then define the combination of normalization and multitask basis pursuit that approxi-
 97 mates this ground truth loss function. We finally give a brute post-processing method for ensuring
 98 that the solution is D sparse.

99 3.1 Ground truth

100 We'd like a ground truth objective to be minimized uniquely by orthonormal matrices, invariant under
 101 rotation, and depend on all changes in the matrix. Deformation [3] and nuclear norm [16] use only a
 102 subset of the differential's information and are not uniquely minimized at unitarity, respectively. We
 103 therefore introduce an alternative ground truth objective that satisfies the above desiderata and has
 104 convenient connections to Isometry Pursuit.

105 We define this objective as

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

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

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

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

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

107 Using Proposition 1, we can check that g is uniquely maximized by orthonormal matrices. Moreover,
 108 $g(X^{-1}) = g(X)$ and g is convex. A graph of g is given in Figure ??, which also shows how this loss
 109 specifically enables comparison with that produced by basis pursuit after normalization as in Section
 110 3.2.

111 Our ground truth program is therefore

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

112 where $\binom{[P]}{d} = \{A \subseteq [P] : |A| = d\}$. Regardless of the convexity of l_c , brute combinatorial search
113 over $[P]$ is inherently non-convex.

114 3.2 Normalization

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

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

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

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

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

119

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

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

$$n, c \mapsto \frac{n}{q(\|n\|_2, c)} \quad (14)$$

123 and matrices by

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

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

124 In particular, we choose q as follows.

$$q : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}^+ \quad (17)$$

$$t, c \mapsto \frac{e^{t^c} + e^{t^{-c}}}{2e}, \quad (18)$$

125 Besides satisfying the conditions in Definition 3, this normalization has some additional nice
126 properties. First, q is convex. Second, it grows asymptotically log-linearly. Third, while
127 $\exp(-|\log t|) = \exp(-\max(t, 1/t))$ is a seemingly natural choice for normalization, it is non
128 smooth, and the LogSumExp replacement of $\max(t, 1/t)$ with $\log(\exp(t) + \exp(1/t))$ simplifies
129 to 17 upon exponentiation [16]. Finally, the parameter c grants control over the width of the basin,
130 which is important in avoiding numerical issues arising close to 0 and ∞ .

131 3.3 Isometry pursuit

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

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

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

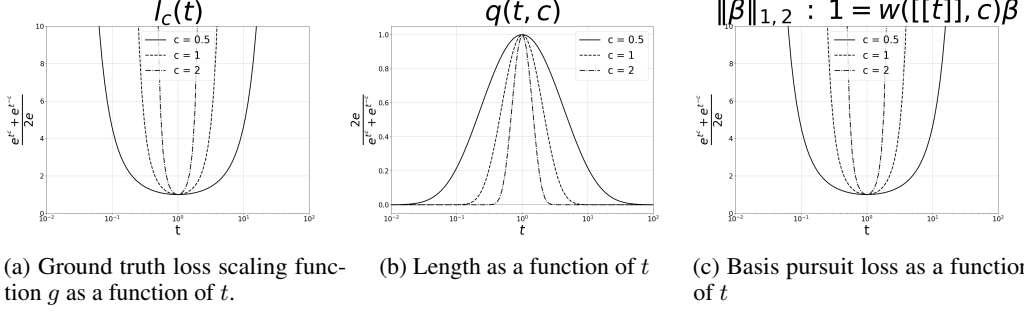


Figure 1: Plots of ground truth loss, vector length, and basis pursuit loss for different values of c in the one-dimensional case $D = 1$. The two losses are equivalent in the one-dimensional case.

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

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

136 Isometry pursuit is then given by

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

137 where I_D is the D dimensional identity matrix and recovered functions are the indices of the dictionary
 138 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** $\hat{S} = S(\hat{\beta})$
-

139 3.4 Theory

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

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

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

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

149 **Proposition 4 (Unitary preference)** Let w_c be a normalization satisfying the conditions in Def-
 150 inition 3. Then $\arg \min_{X, S \in \mathbb{R}^{D \times D}} \hat{\beta}_c^D(X)$ is orthonormal. Moreover when X is orthonormal,
 151 $(\min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} : I_D = w(X, c)\beta) = D$.

152 While this Proposition falls short of showing that an orthonormal submatrix will be selected should
 153 one be present, it provides intuition justifying the preferential efficacy of ISOMETRYPURSUIT on
 154 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 [17].

4 Experiments

We compare TWOSTAGEISOMETRYPURSUIT and GREEDYSEARCH on the Iris and Wine datasets [18, 19, 20], as well as the Ethanol dataset from Chmiela et al. [21], Koelle et al. [9]. 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. [9] 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 [22] from CVXPY [23, 24], which is able to push sparse values arbitrarily close to 0 [25]. 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 $	$l_1(X_{\widehat{S}})$	$P(l_1(X_{\widehat{S}_G}), l_1(X_{\widehat{S}}))$
Iris	4	75	25	13.4 ± 6.4	7 ± 1	8.0 ± 1.8	10^{-4}
Wine	5	89	25	$5.7 \pm .2$	12 ± 1	$5.6 \pm .1$	5×10^{-5}
Ethanol	2	756	100	$2.6 \pm .3$	90 ± 164	$2.5 \pm .2$	2×10^{-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

It could be used in the stitching step of an algorithm like the kohli one We leave aside the question of patch alignment The full gradient approach. In this case normalization prior to projection is subsummed by the larger coefficients needed to get the tangent space. Good news is tangent space estimation need not be performed. Let's compare the coefficients involved in projecting versus not projecting. We can perform regression in the high dimensional space instead of projecting on span of target variable.

With respect to pseudoinverse estimation, sparse methods have been applied in [26] Even though by Lagrangian duality, the basis pursuit solution corresponds to λ approaching 0, the solution is sparse [27].

While the sparse PCA problem is non-convex, our approach can be taken as a simpler version in the sense that the loadings are constrained to be the identify matrix. [6] gives a method for solving the sparse-PCA method more efficiently than the original greedy approach. Compared with the FISTA method used in [9, 1], coordinate descent [28, 29, 30] is faster [31, 32]. Compared with [33, 34], the sklearn multitask lasso is 2, 1 rather than ∞ , 1 regularized. This misses the utility of our

normalization for finding unitary matrices since isometry embeddings preserve important properties like distances between points.

Our notion of isometric recovery is distinct from the restricted isometry property [35, 36], which is used to show guaranteed recovery at fast convergence rates in supervised learning. In particular, our approach does not consider statistical error or the presence of a true underlying model. However, we note that disintegration of performance at high λ values in the lasso formulation may have some relation to these properties, as discussed in [9, 1].

A major area of comparison is in diversification in recommendation systems where greedy algorithms are used [37, 38], and also in document diversification for Retrieval Augmented Generation.

The most pressing piece of theoretical work which remains on this topic is the removal of the restriction $|S| = D$ on the conditions of Proposition. The resulting proposition, which seems almost obvious, is in fact more difficult to argue, and is seemingly violated by empirical results. Nevertheless, these violations are subtly non-dispositive since absence of sparsity and improvements of primal loss below D are accompanied by violations of the constraint of a similar magnitude, suggesting that we a more refined approach to optimization, substantial improvements in estimation accuracy may be possible. From a geometric perspective, we note that isometries may not always exist in the presence of curvature, and comparison of our loss with curvature could prove fertile. Finally, the speed increases garnered by the particularly simple form of our algorithm warrants comparison with other pseudoinverse estimators warrant further comparison. An extension of our estimator we omit for brevity is to use the multitask lasso formulation to trim the size of \hat{P} . Besides our new normalization, this results in a simpler procedure than in Koelle [39] and Koelle et al. [1] that is amenable to the more performant multitask lasso solver in sklearn.

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] 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.
- [3] Dhruv Kohli, Alexander Cloninger, and Gal Mishne. LDLE: Low distortion local eigenmaps. *J. Mach. Learn. Res.*, 22, 2021.
- [4] Peter W Jones, Mauro Maggioni, and Raanan Schul. Universal local parametrizations via heat kernels and eigenfunctions of the laplacian. September 2007.
- [5] 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.
- [6] D Bertsimas and Driss Lahlou Kitane. Sparse PCA: A geometric approach. *J. Mach. Learn. Res.*, 24:32:1–32:33, October 2022.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] 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.

- [11] 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. Summers, 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>.
- [12] Scott Shaobing Chen and David L. Donoho and Michael A. Saunders. Atomic decomposition by basis pursuit. *SIAM REVIEW*, 43(1):129, February 2001.
- [13] 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.
- [14] G Obozinski, B Taskar, and Michael I Jordan. Multi-task feature selection. 2006.
- [15] Dit-Yan Yeung and Yu Zhang. A probabilistic framework for learning task relationships in multi-task learning. 2011.
- [16] Stephen P Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, March 2004.
- [17] Tim Hesterberg, Nam Hee Choi, Lukas Meier, and Chris Fraley. Least angle and ℓ_1 penalized regression: A review. February 2008.
- [18] R. A. Fisher. Iris. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C56C76>.
- [19] Stefan Aeberhard and M. Forina. Wine. UCI Machine Learning Repository, 1991. DOI: <https://doi.org/10.24432/C5PC7J>.
- [20] 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.
- [21] 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.
- [22] 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>.
- [23] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [24] 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.
- [25] CVXPY Developers. Sparse solution with cvxpy. https://www.cvxpy.org/examples/applications/sparse_solution.html. Accessed: 2024-07-11.
- [26] Tingni Sun and Cun-Hui Zhang. Sparse matrix inversion with scaled lasso. February 2012.
- [27] Joel A. Tropp. Just relax: Convex programming methods for subset selection and sparse approximation. Technical Report 04-04, Institute for Computational Engineering and Sciences, The University of Texas at Austin, 2004.
- [28] Jerome Friedman, Trevor Hastie, Holger Höfling, and Robert Tibshirani. Pathwise coordinate optimization. *The Annals of Applied Statistics*, 1(2):302 – 332, 2007. doi: 10.1214/07-AOAS131. URL <https://doi.org/10.1214/07-AOAS131>.

- [29] L Meier, Sara van de Geer, and P Bühlmann. The group lasso for logistic regression. *J. R. Stat. Soc. Series B Stat. Methodol.*, 70, February 2008.
- [30] Zhiwei Qin, Katya Scheinberg, and Donald Goldfarb. Efficient block-coordinate descent algorithms for the group lasso. *Math. Program. Comput.*, 5(2):143–169, June 2013.
- [31] Alejandro Catalina, Carlos M Alaíz, and José R Dorronsoro. Revisiting FISTA for lasso: Acceleration strategies over the regularization path. *Eur Symp Artif Neural Netw*, 2018.
- [32] Yujie Zhao and Xiaoming Huo. A survey of numerical algorithms that can solve the lasso problems. March 2023.
- [33] Joel A. Tropp. Algorithms for simultaneous sparse approximation. part ii: Convex relaxation. *Signal Processing*, 86(3):589–602, 2006. ISSN 0165-1684. doi: <https://doi.org/10.1016/j.sigpro.2005.05.031>. URL <https://www.sciencedirect.com/science/article/pii/S0165168405002239>. Sparse Approximations in Signal and Image Processing.
- [34] Han Liu, Mark Palatucci, and Jian Zhang. Blockwise coordinate descent procedures for the multi-task lasso, with applications to neural semantic basis discovery. *ICML*, pages 649–656, June 2009.
- [35] Emmanuel Candes and Terence Tao. Decoding by linear programming. February 2005.
- [36] T Hastie, R Tibshirani, and M Wainwright. Statistical learning with sparsity: The lasso and generalizations. May 2015.
- [37] 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.
- [38] Qiong Wu, Yong Liu, Chunyan Miao, Yin Zhao, Lu Guan, and Haihong Tang. Recent advances in diversified recommendation. May 2019.
- [39] Samson Jonathan Koelle. *Geometric algorithms for interpretable manifold learning*. Phd thesis, University of Washington, 2022. URL <http://hdl.handle.net/1773/48559>. Statistics [108].

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  $\mathcal{X} \in \mathbb{R}^{D \times P}$ , objective  $f$ )
1: for each combination  $S \subseteq \{1, 2, \dots, P\}$  with  $|S| = D$  do
2:   Evaluate  $f(\mathcal{X}_{.S})$ 
3: end for
4: Output the combination  $S^*$  that minimizes  $f(\mathcal{X}_{.S})$ 

```

GREEDYSEARCH(Matrix $\mathcal{X} \in \mathbb{R}^{D \times P}$, 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  $f(\mathcal{X}_{(S \cup \{p\})})$ 
8:     if  $f(\mathcal{X}_{(S \cup \{p\})}) < f_{\text{best}}$  then
9:       Update  $S_{\text{best}} = S \cup \{p\}$ 
10:      Update  $f_{\text{best}} = f(\mathcal{X}_{(S \cup \{p\})})$ 
11:     end if
12:   end for
13:   Return Greedy( $\mathcal{X}, f, S_{\text{best}}, d + 1$ )
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 β .

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

Proof:

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

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

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

□

We then show that this implies that the resultant loss is unchanged by unitary transformation of \mathcal{X} .

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

Proof:

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

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

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

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

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

□

6.2.2 Proof of Proposition 4

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

325 **Proof:** The value of D is clearly obtained by β orthonormal, since by Proposition ??, for X
 326 orthogonal, without loss of generality

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

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

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

332 Therefore, it remains to show that X consisting of orthogonal vectors of length 1 have lower compared
 333 with X consisting of non-orthogonal vectors. Invertible matrices $X_{.S}$ admit QR decompositions
 334 $\tilde{X}_{.S} = QR$ where Q and R are orthonormal and upper-triangular matrices, respectively [?].
 335 Denoting Q to be composed of basis vectors $[e^1 \dots e^d]$, the matrix R has form

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

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

341

□