

# Isometry pursuit

Samson Koelle, Marina Meila

July 25, 2024

## **Abstract**

Isometry pursuit is an algorithm for identifying unitary column-submatrices of wide matrices in polynomial time. It achieves sparsity via use of the group lasso norm, and therefore has constrained and penalized formulations. Applied to tabular data, it selects a subset of columns that maximize diversity. Applied to Jacobians of putative coordinate functions, it identifies isometric embeddings from within dictionaries. It therefore has relevance to interpretability of learned representations.

# 1 Method

Recall that our objective is to, given a rank  $D$  matrix  $\mathcal{X} \in \mathbb{R}^{D \times P}$  with  $P > D$ , select a square submatrix  $\mathcal{X}_{\mathcal{S}}$  where subset  $\mathcal{S} \subset P$  satisfies  $|\mathcal{S}| = D$  that is as unitary as possible. Thus, we first will define a function that is uniquely minimized by unitary matrices and some favorable properties for optimization that will be the ground truth we evaluate the success of our method against. We then define the combination of normalization and multitask basis pursuit that approximates this ground truth loss function. We include claims that ground truth and convex loss values are the same for all diagonalizable matrices, and that the convex basis pursuit program recovers to optimum in a deterministic manner should it exist; proofs are given in Section ?? . We finally define the lasso dual to the basis pursuit program and a post processing method for ensuring that the solution is  $D$  sparse. Experimental results using these methods will then be given in Section 3

## 1.1 Ground truth

The main goal of isometry pursuit is to expediate the selection of unitary submatrices. More traditional measures of unitariness which use the singular values of a matrix like the log operator norm (i.e. log deformation) and nuclear norm are poorly suited for optimization since they use a subset of the matrix's information and are not uniquely minimized at unitarity, respectively. Thus, we define the loss

$$l_c : \mathbb{R}^{D \times P} \rightarrow \mathbb{R}^+ \quad (1)$$

$$\mathcal{X} \mapsto \sum_{d=1}^D g(\sigma^d(\mathcal{X}), c) \quad (2)$$

where  $\sigma^d((X))$  is the  $d$ -th singular value of  $\mathcal{X}$  and

$$g : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}^+ \quad (3)$$

$$t, c \mapsto e^{t^c} + e^{t^{-c}}. \quad (4)$$

Plainly,  $g$  is uniquely maximized by unitary matrices, and  $g(\mathcal{X}^\dagger) = g(\mathcal{X})^{-1}$ . The former condition is necessary for success of the method, while the latter, as well as the convexity of  $g$ , are somewhat aesthetic choices. A graph of  $g$  is given in Figure 1. Most importantly, this loss enables comparison with produced after normalization as in Section 1.2.

The overall algorithm we seek to improve upon is

$$\widehat{\mathcal{S}}_{GT} = \arg \min_{\mathcal{S} \subseteq [P]: |\mathcal{S}|=D} l_c(\mathcal{X}_{\mathcal{S}}) \quad (5)$$

In practice, non-convexity occurs in two places, but only one is essential. The inessential non-convexity is in the computation of  $l_c$ . While this function is in fact convex, computation of the individual singular values prior to summation is not, and our experiments rely on such piecemeal computation rather than implementing an end-to-end method. However, the combinatorial search over  $[P]$  is inherently non-convex, and requires combinatorial search over all combinations.

## 1.2 Normalization

We will propose a basis pursuit method which approximates the results of Program 5. Since basis pursuit methods tend to select longer vectors, selection of unitary submatrices requires normalization such that long and short candidate basis vectors are penalized in the subsequent regression. This calls for a "normalization" method that differs from other forms in its requirements, and we can't yet prove that these conditions relate it to any sort of norm, even on an appropriately chosen space. This normalization is Now establish some basic conditions for normalization of vectors  $v \in \mathbb{R}^D$ .

**Definition 1 (Symmetric normalization)** *A function  $q : \mathbb{R}^D \rightarrow \mathbb{R}^+$  is a symmetric normalization if*

$$\arg \max_{v \in \mathbb{R}^D} q(v) = \{v : \|v\| = 1\} \quad (6)$$

$$q(v) = q\left(\frac{v}{\|v\|^2}\right) \quad (7)$$

$$q(v^1) = q(v^2) \quad \forall v^1, v^2 : \|v^1\| = \|v^2\| \quad (8)$$

Note that requiring the full structure of a multiplicative norm here is unnecessary for basic success of the algorithm, but certain characteristics such as  $q(v^{-1}) = q(v)$  seem desirable, provided one can give a reasonable way to compute  $v^{-1}$ , such as by considering each vector as a scaled rotation subgroup of the general linear group. Mindful of this opportunity, and also of the desire to compare with the ground truth and provide computational expediency, consider the normalization by

$$q : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}^+ \quad (9)$$

$$t, c \mapsto \frac{e^{t^c} + e^{t^{-c}}}{2e}, \quad (10)$$

and use this to define the vector normalization

$$n : \mathbb{R}^D \times \mathbb{R}^+ \rightarrow \mathbb{R}^D \quad (11)$$

$$n, c \mapsto \frac{n}{q(\|n\|_2, c)} \quad (12)$$

and matrix normalization

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

$$\mathcal{X}_{.p}, c \mapsto n(\mathcal{X}_{.p}, c) \quad \forall p \in [P]. \quad (14)$$

While this normalization satisfies 1, it also 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 9 upon exponentiation. Finally, the parameter  $c$  grants control over the width of the basin, which is important in avoiding numerical issues arising close to 0 and  $\infty$ . This completes the deterministic data preprocessing.

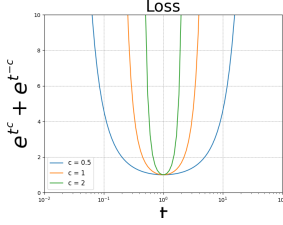


Figure 1: Ground truth loss scaling function  $g$  as a function of  $t$

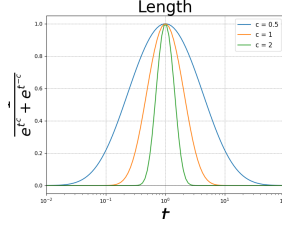


Figure 2: Length as a function of  $t$

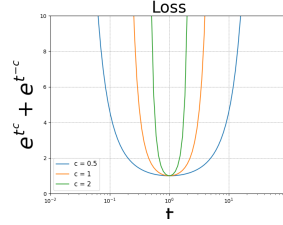


Figure 3: Basis pursuit losses as a function of  $t$

Figure 4: Plots of Length and Loss for different values of  $c$ . Since  $t$  is one dimensional and therefore diagonalizable, basis pursuit and ground truth give identical loss values.

### 1.3 Isometry pursuit

We will show how to use an appropriate normalized matrix  $w(\mathcal{X})$  in multitask basis pursuit to identify submatrices of  $\mathcal{X}$  that are as unitary as possible. Multitask basis pursuit is a method for identifying sparse signals from overcomplete dictionaries, and the intuition behind its application in our setting is that submatrices consisting of vectors which are closer to 1 in length and more orthogonal will have smaller loss. In contrast to the typical statistical setting, these features correspond to individual observations in our diversification example, and basis vectors of data manifold tangent spaces in our non-linear dimension reduction example.

The multitask basis pursuit penalty is

$$\|\cdot\|_{1,2} : \mathbb{R}^{P \times D} \rightarrow \mathbb{R}^+ \quad (15)$$

$$\beta \mapsto \sum_{p=1}^P \|\beta_p\|_2. \quad (16)$$

The isometry pursuit program is then

$$\hat{\beta}_c^P(\mathcal{X}) = \arg \min_{\beta \in \mathbb{R}^{P \times D}} \|\beta\|_{1,2} \text{ s.t. } I_D = w(\mathcal{X}, c)\beta. \quad (17)$$

The recovered supports are then given by  $S(\beta)$  where

$$S : \mathbb{R}^{P \times d} \rightarrow \binom{\{1, 2, \dots, P\}}{d} \quad (18)$$

$$\beta \mapsto \{p \in \{1, 2, \dots, P\} : \|\beta_p\| > 0\} \quad (19)$$

and  $\binom{\{1, 2, \dots, P\}}{d} = \{A \subseteq \{1, 2, \dots, P\} : |A| = d\}$  are the indices of the dictionary elements with non-zero coefficients.

---

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

---

1: **Output**  $\hat{S} = S(\hat{\beta}_P(w_c(\mathcal{X})))$

---

A key theoretical assertion is that selection methods  $S(\hat{\beta}(\mathcal{X}))$  are invariant to choice of basis for  $\mathcal{X}$ .

**Proposition 1 (Basis pursuit selection equivalence)** *Let  $U \in \mathbb{R}^{D \times D}$  be unitary. Then  $S(\hat{\beta}(U\mathcal{X})) = S(\hat{\beta}(\mathcal{X}))$ .*

With these preliminaries, we may state our main result.

**Proposition 2 (Unitary selection)** *Given a matrix  $\mathcal{X} \in \mathbb{R}^{D \times P}$  with a rank  $D$  submatrix  $\mathcal{X}_S \in \mathbb{R}^{D \times D}$  that is unitary,  $S = S(\hat{\beta}(\mathcal{X}))$*

This proof admits two immediate generalizations. First, any normalization function that satisfies the normalization conditions will do. Second, the ground truth and convex losses are equivalent for diagonalizable matrices.

#### 1.4 Isometric lasso

The convex loss function 15 and linear constraint in 17 admit a Lagrangian dual which we shall call Isometric Lasso. The Isometric Lasso loss is

$$l_\lambda(\mathcal{X}, \beta) = \|I_D - \tilde{\mathcal{X}}_c \beta\|_2^2 + \lambda \|\beta\|_{1,2} \quad (20)$$

which can be optimized as

$$\hat{\beta}_\lambda(\mathcal{X}) = \arg \min_{\beta \in \mathbb{R}^{P \times D}} l_\lambda(\mathcal{X}, \beta) \quad (21)$$

Similarly to Section ??, we assert that  $S(\hat{\beta}_\lambda(\mathcal{X}))$ .

**Proposition 3 (Lasso selection equivalence)** *Let  $U \in \mathbb{R}^{D \times D}$  be unitary. Then  $S(\hat{\beta}_\lambda(U\mathcal{X})) = S(\hat{\beta}_\lambda(\mathcal{X}))$ .*

### 1.5 Extension to non-linear spaces

**Proposition 4 (Local isometry selection)** *Given a set of functions  $G$  that contains a subset that defines a locally isometric embedding at a point  $\xi$ , then these will be selected as  $\arg \min_{\beta}$ .*

A proof is given in Section ??.

### 1.6 Two-stage isometry pursuit

#### 1.7 Implementation

We use the multitask lasso from sklearn and the cvxpy package for basis pursuit. We use the SCS interior point solver from CVXPY, which is able to push sparse values arbitrarily close to 0 **cvxpy's sparse solution**. Data is IRIS and Wine, as well as flat torus from ldle.

### 1.8 Computational complexity

## 2 Experiments

Comparison with isometry loss.

### 3 Experiments

Comparison with isometry loss.