

Non-negative Matrix Factorization via Archetypal Analysis

Hamid Javadi* and Andrea Montanari†

May 8, 2017

Abstract

Given a collection of data points, non-negative matrix factorization (NMF) suggests to express them as convex combinations of a small set of ‘archetypes’ with non-negative entries. This decomposition is unique only if the true archetypes are non-negative and sufficiently sparse (or the weights are sufficiently sparse), a regime that is captured by the separability condition and its generalizations.

In this paper, we study an approach to NMF that can be traced back to the work of Cutler and Breiman [CB94] and does not require the data to be separable, while providing a generally unique decomposition. We optimize the trade-off between two objectives: we minimize the distance of the data points from the convex envelope of the archetypes (which can be interpreted as an empirical risk), while minimizing the distance of the archetypes from the convex envelope of the data (which can be interpreted as a data-dependent regularization). The archetypal analysis method of [CB94] is recovered as the limiting case in which the last term is given infinite weight.

We introduce a ‘uniqueness condition’ on the data which is necessary for exactly recovering the archetypes from noiseless data. We prove that, under uniqueness (plus additional regularity conditions on the geometry of the archetypes), our estimator is robust. While our approach requires solving a non-convex optimization problem, we find that standard optimization methods succeed in finding good solutions both for real and synthetic data.

1 Introduction

Given a set of data points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, it is often useful to represent them as convex combinations of a small set of vectors (the ‘archetypes’ $\mathbf{h}_1, \dots, \mathbf{h}_\ell$):

$$\mathbf{x}_i \approx \sum_{\ell=1}^r w_{i,\ell} \mathbf{h}_\ell, \quad w_{i,\ell} \geq 0, \quad \sum_{\ell=1}^r w_{i,\ell} = 1. \quad (1.1)$$

Decompositions of this type have wide ranging applications, from chemometrics [PT94] to image processing [LS99] and topic modeling [XLG03]. As an example, Figure 1 displays the infrared reflection spectra¹ of four molecules (caffeine, sucrose, lactose and trioctanoin) for wavenumber between 1186 cm^{-1} and 1530 cm^{-1} . Each spectrum is a vector $\mathbf{h}_{0,\ell} \in \mathbb{R}^d$, with $d = 87$ and $\ell \in \{1, \dots, 4\}$. If a mixture of these substances is analyzed, the resulting spectrum will be a convex combination of the spectra of the four analytes. This situation arises in hyperspectral imaging [MBDC⁺14], where a main focus is to estimate spatially varying proportions of a certain

*Department of Electrical Engineering, Stanford University

†Department of Electrical Engineering and Statistics, Stanford University

¹Data were retrieved from the NIST Chemistry WebBook dataset [LM].

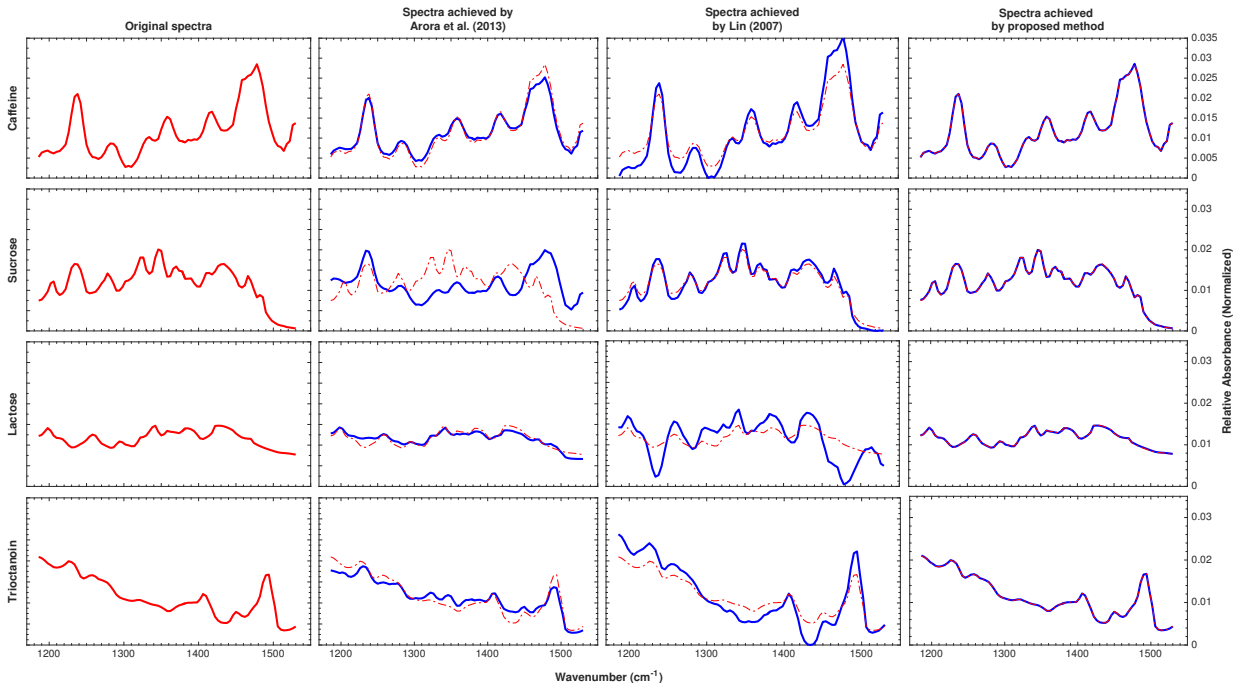


Figure 1: Left column: Infrared reflection spectra of four molecules. Subsequent columns: Spectra estimated from $n = 250$ spectra of mixtures of the four original substances (synthetic data generated by taking random convex combinations of the pure spectra, see Appendix A for details). Each column reports the results obtained with a different estimator: continuous blue lines correspond to the reconstructed spectra; dashed red lines correspond to the ground truth.

number of analytes. In order to mimic this setting, we generated $n = 250$ synthetic random convex combinations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ of the four spectra $\mathbf{h}_{0,1}, \dots, \mathbf{h}_{0,4}$, each containing two or more of these four analytes, and tried to reconstruct the pure spectra from the \mathbf{x}_i 's. Each column in Figure 1 displays the reconstruction obtained using a different procedure. We refer to Appendix A for further details.

Without further constraints, the decomposition (1.1) is dramatically underdetermined. Given a set of valid archetypes $\{\mathbf{h}_{0,\ell}\}_{\ell \leq r}$, any set $\{\mathbf{h}_\ell\}_{\ell \leq r}$ whose convex hull contains the $\{\mathbf{h}_{0,\ell}\}_{\ell \leq r}$ also satisfies Eq. (1.1). For instance, we can set $\mathbf{h}_\ell = \mathbf{h}_{0,\ell}$ for $\ell \leq r - 1$, and $\mathbf{h}_r = (1 + s)\mathbf{h}_{0,r} - s\mathbf{h}_{0,1}$ for any $s \geq 0$, and obtain an equally good representation of the data $\{\mathbf{x}_i\}_{i \leq n}$.

How should we constrain the decomposition (1.1) in such a way that it is generally unique (up to permutations of the r archetypes)? Since the seminal work of Paatero and Tapper [PT94, Paa97], and of Lee and Seung [LS99, LS01], an overwhelming amount of work has addressed this question by making the assumptions that the archetypes are componentwise non-negative $\mathbf{h}_\ell \geq 0$. Among other applications the non-negativity constraint is justified for chemometrics (reflection or absorption spectra are non-negative), and topic modeling (in this case archetypes correspond to topics, which are represented as probability distributions over words). This formulation has become popular as non-negative matrix factorization (NMF).

Under the non-negativity constraint $\mathbf{h}_\ell \geq 0$ the role of weights and archetypes becomes symmetric, and the decomposition (1.1) is unique provided that the archetypes or the weights are sufficiently sparse (without loss of generality one can assume $\sum_{\ell=1}^r h_{\ell,i} = 1$). This point was clar-

ified by Donoho and Stodden [DS03], introduced a separability condition that ensure uniqueness. The non-negative archetypes $\mathbf{h}_1, \dots, \mathbf{h}_r$ are separable if, for each $\ell \in [r]$ there exists an index $i(\ell) \in [d]$ such that $(\mathbf{h}_\ell)_{i(\ell)} = 1$, and $(\mathbf{h}_{\ell'})_{i(\ell)} = 0$ for all $\ell' \neq \ell$. If we exchange the roles of weights $\{w_{i,\ell}\}$ and archetypes $\{h_{\ell,i}\}$, separability requires that $\ell \in [r]$ there exists an index $i(\ell) \in [n]$ such that $w_{i(\ell),\ell} = 1$, and $w_{i(\ell),\ell'} = 0$ for all $\ell' \neq \ell$. This condition has a simple geometric interpretation: the data are separable if for each archetype \mathbf{h}_ℓ there is at least one data point \mathbf{x}_i such that $\mathbf{x}_i = \mathbf{h}_\ell$. A copious literature has developed algorithms for non-negative matrix factorization under separability condition or its generalizations [DS03, AGKM12, RRTB12, AGH⁺13, GZ15].

Of course this line of work has a drawback: in practice we do not know whether the data are separable. (We refer to the Section 5 for a comparison with [GZ15], which relaxes the separability assumption.) Further, there are many cases in which the archetypes $\mathbf{h}_1, \dots, \mathbf{h}_\ell$ are not necessarily non-negative. For instance, in spike sorting, the data are measurements of neural activity and the archetypes correspond to waveforms associated to different neurons [RCP09]. In other applications the archetypes \mathbf{h}_ℓ are non-negative, but –in order to reduce complexity– the data $\{\mathbf{x}_i\}_{i \leq n}$ are replaced by a random low-dimensional projection [KSD08, WL10]. The projected archetypes lose the non-negativity property. Finally, the decomposition (1.1) is generally non-unique, even under the constraint $\mathbf{h}_\ell \geq 0$. This is illustrated, again, in Figure 1: all the spectra are strictly positive, and hence we can find archetypes $\mathbf{h}_1, \dots, \mathbf{h}_4$ that are still non-negative and whose convex envelope contains $\mathbf{h}_{0,1}, \dots, \mathbf{h}_{0,4}$.

Since NMF is underdetermined, standard methods fail in such applications, as illustrated in Figure 1. We represent the data as a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ whose i -th row is the vector \mathbf{x}_i , the weights by a matrix $\mathbf{W} = (w_{i,\ell})_{i \leq n, \ell \leq r} \in \mathbb{R}^{n \times d}$ and the prototypes by a matrix $\mathbf{H} = (h_{\ell,j})_{\ell \leq r, j \leq d} \in \mathbb{R}^{r \times d}$. The third column of Figure 1 uses a projected gradient algorithm from [Lin07] to solve the problem

$$\begin{aligned} & \text{minimize} && \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2, \\ & \text{subject to} && \mathbf{W} \geq 0, \mathbf{H} \geq 0. \end{aligned} \tag{1.2}$$

Empirically, projected gradient converges to a point with very small fitting error $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$, but the reconstructed spectra (rows of \mathbf{H}) are inaccurate. The second column in the same figure shows the spectra reconstructed using an algorithm from [AGH⁺13], that assumes separability: as expected, the reconstruction is not accurate.

In a less widely known paper, Cutler and Breiman [CB94] addressed the same problem using what they call ‘archetypal analysis.’ Archetypal analysis presents two important differences with respect to standard NMF: (1) The archetypes \mathbf{h}_ℓ are not necessarily required to be non-negative (although this constraint can be easily incorporated); (2) The under-determination of the decomposition (1.1) is addressed by requiring that the archetypes belong to the convex hull of the data points: $\mathbf{h}_\ell \in \text{conv}(\{\mathbf{x}_i\}_{i \leq n})$.

In applications the condition $\mathbf{h}_\ell \in \text{conv}(\{\mathbf{x}_i\}_{i \leq n})$ is too strict. This paper builds on the ideas of [CB94] to propose a formulation of NMF that is uniquely defined (barring degenerate cases) and provides a useful notion of optimality. In particular, we present the following contributions.

Archetypal reconstruction. We propose to reconstruct the archetypes $\mathbf{h}_1, \dots, \mathbf{h}_r$ by optimizing a combination of two objectives. On one hand, we minimize the error in the decomposition (1.1). This amounts to minimizing the distance between the data points and the convex hull of the archetypes. On the other hand, we minimize the distance of the archetypes from the convex hull of data points. This relaxes the original condition imposed in [CB94] which required the archetypes to lie in $\text{conv}(\{\mathbf{x}_i\})$, and allows to treat non-separable data.

Robustness guarantee. We next assume that the decomposition (1.1) approximately hold for some ‘true’ archetypes \mathbf{h}_ℓ^0 and weights $w_{i,\ell}^0$, namely $\mathbf{x}_i = \mathbf{x}_i^0 + \mathbf{z}_i$, where $\mathbf{x}_i^0 = \sum_{\ell=1}^r w_{i,\ell}^0 \mathbf{h}_\ell^0$ and

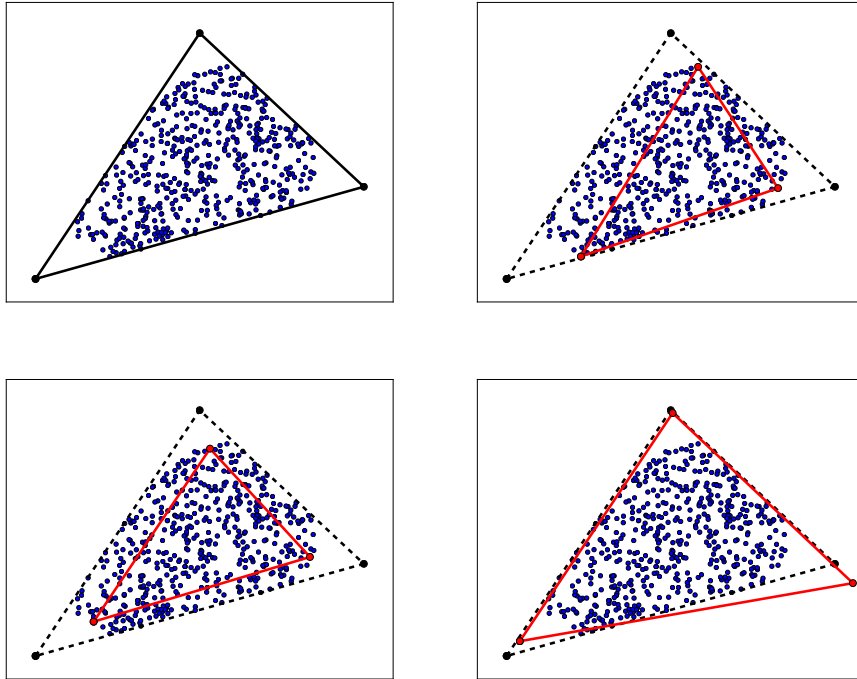


Figure 2: Toy example of archetype reconstruction. Top left: data points (blue) are generated as random linear combinations of $r = 3$ archetypes in $d = 2$ dimensions (red, see Appendix A for details). Top right: Initialization using the algorithm of [AGH⁺13]. Bottom left: Output of the alternate minimization algorithm of [CB94] with initialization from the previous frame. Bottom right: Alternate minimization algorithm to compute the estimator (2.4), with $\lambda = 0.0166$.

\mathbf{z}_i captures unexplained effects. We introduce a ‘uniqueness condition’ on the data $\{\mathbf{x}_i^0\}_{i \leq n}$ which is necessary for exactly recovering the archetypes from the noiseless data. We prove that, under uniqueness (plus additional regularity conditions on the geometry of the archetypes), our estimator is robust. Namely it outputs archetypes $\{\hat{\mathbf{h}}_\ell\}_{\ell \leq r}$ whose distance from the true ones $\{\mathbf{h}_\ell^0\}_{\ell \leq r}$ (in a suitable metric) is controlled by $\sup_{i \leq n} \|\mathbf{z}_i\|_2$.

Algorithms. Our approach reconstructs the archetypes $\mathbf{h}_1, \dots, \mathbf{h}_r$ by minimizing a non-convex risk function $\mathcal{R}_\lambda(\mathbf{H})$. We propose three descent algorithms that appear to perform well on realistic instances of the problem. In particular, Section 4 introduces a proximal alternating linearized minimization algorithm (PALM) that is guaranteed to converge to critical points of the risk function. Appendix E discusses two alternative approaches. One possible explanation for the success of such descent algorithms is that reasonably good initializations can be constructed using spectral methods, or approximating the data as separable, cf. Section 4.1. We defer a study of global convergence of this two-stages approach to future work.

2 An archetypal reconstruction approach

Let $\mathcal{Q} \subseteq \mathbb{R}^d$ be a convex set and $D : \mathcal{Q} \times \mathcal{Q} \rightarrow \mathbb{R}$, $(\mathbf{x}, \mathbf{y}) \mapsto D(\mathbf{x}; \mathbf{y})$ a loss function on \mathcal{Q} . For a point $\mathbf{u} \in \mathcal{Q}$, and a matrix $\mathbf{V} \in \mathbb{R}^{m \times d}$, with rows $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathcal{Q}$, we let

$$\mathcal{D}(\mathbf{u}; \mathbf{V}) \equiv \min \left\{ D(\mathbf{u}; \mathbf{V}^\top \boldsymbol{\pi}) : \boldsymbol{\pi} \in \Delta^m \right\}, \quad (2.1)$$

$$\Delta^m \equiv \left\{ \mathbf{x} \in \mathbb{R}_{\geq 0}^m : \langle \mathbf{x}, \mathbf{1} \rangle = 1 \right\}. \quad (2.2)$$

In other words, denoting by $\text{conv}(\mathbf{V}) = \text{conv}(\{\mathbf{v}_1, \dots, \mathbf{v}_m\})$ the convex hull of the rows of matrix \mathbf{V} , $\mathcal{D}(\mathbf{u}; \mathbf{V})$ is the minimum loss between \mathbf{x} and any point in $\text{conv}(\mathbf{V})$. If $\mathbf{U} \in \mathbb{R}^{k \times d}$ is a matrix with rows $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathcal{Q}$, we generalize this definition by letting

$$\mathcal{D}(\mathbf{U}; \mathbf{V}) \equiv \sum_{\ell=1}^k \mathcal{D}(\mathbf{u}_\ell; \mathbf{V}). \quad (2.3)$$

While this definition makes sense more generally, we have in mind two specific examples in which $D(\mathbf{x}; \mathbf{y})$ is actually separately convex in its arguments \mathbf{x} and \mathbf{y} . (Most of our results will concern the first example.)

Example 2.1 (Square loss). In this case $\mathcal{Q} = \mathbb{R}^d$, and $D(\mathbf{x}; \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$. This is the case originally studied by Cutler and Breiman [CB94].

Example 2.2 (KL divergence). We take $\mathcal{Q} = \Delta^d$, the d -dimensional simplex, and $D(\mathbf{x}; \mathbf{y})$ to be the Kullback-Leibler divergence between probability distributions \mathbf{x} and \mathbf{y} , namely $D(\mathbf{x}; \mathbf{y}) \equiv \sum_{i=1}^d x_i \log(x_i/y_i)$.

Given data $\mathbf{x}_1, \dots, \mathbf{x}_n$ organized in the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, we estimate the archetypes by solving the problem²

$$\widehat{\mathbf{H}}_\lambda \in \arg \min \left\{ \mathcal{D}(\mathbf{X}; \mathbf{H}) + \lambda \mathcal{D}(\mathbf{H}; \mathbf{X}) : \mathbf{H} \in \mathcal{Q}^r \right\}, \quad (2.4)$$

where we denote by \mathcal{Q}^r the set of matrices $\mathbf{H} \in \mathbb{R}^{r \times d}$ with rows $\mathbf{h}_1, \dots, \mathbf{h}_r \in \mathcal{Q}$. A few values of λ are of special significance. If we set $\lambda = 0$, and $\mathcal{Q} = \Delta^d$, we recover the standard NMF objective (1.2), with a more general distance function $D(\cdot, \cdot)$. As pointed out above, in general this objective has no unique minimum. If we let $\lambda \rightarrow 0+$ after the minimum is evaluated, $\widehat{\mathbf{H}}_\lambda$ converges to the minimizer of $\mathcal{D}(\mathbf{X}; \mathbf{H})$ which is the ‘closest’ to the convex envelope of the data $\text{conv}(\mathbf{X})$ (in the sense of minimizing $\mathcal{D}(\mathbf{H}; \mathbf{X})$). Finally as $\lambda \rightarrow \infty$, the archetypes \mathbf{h}_ℓ are forced to lie in $\text{conv}(\mathbf{X})$ and hence we recover the method of [CB94].

Figure 2 illustrates the advantages of the estimator (2.4) on a small synthetic example, with $d = 2$, $r = 3$, $n = 500$: in this case the data are non separable. We first use the successive projections algorithm of [AGH⁺13] (that is designed to deal with separable data) in order to estimate the archetypes. As expected, the reconstruction is not accurate because this algorithm assumes separability and hence estimates the archetypes with a subset of the data points. We then use these estimates as initialization in the alternate minimization algorithm of [CB94], which optimizes the objective (2.4) with $\lambda = \infty$. The estimates improve but not substantially: they are still constrained to lie in $\text{conv}(\mathbf{X})$. A significant improvement is obtained by setting λ to a small value. We

²This problem can have multiple global minima if $\lambda = 0$ or in degenerate settings. One minimizer is selected arbitrarily when this happens.

(approximately) minimize the cost function (2.4) by generalizing the alternate minimization algorithm, cf. Section 4. The optimal archetypes are no longer constrained to $\text{conv}(\mathbf{X})$, and provide a better estimate of the true archetypes. The last column in Figure 1 uses the same estimator, and approximately solves problem (2.4) by gradient descent algorithm.

In our analysis we will consider a slightly different formulation in which the Lagrangian of Eq. (2.4) is replaced by a hard constraint:

$$\begin{aligned} & \text{minimize} && \mathcal{D}(\mathbf{H}; \mathbf{X}), \\ & \text{subject to} && \mathcal{D}(\mathbf{x}_i; \mathbf{H}) \leq \delta^2 \quad \text{for all } i \in \{1, \dots, n\}. \end{aligned} \tag{2.5}$$

We will use this version in the analysis presented in the next section, and denote the corresponding estimator by $\widehat{\mathbf{H}}$.

3 Robustness

In order to analyze the robustness properties of estimator $\widehat{\mathbf{H}}$, we assume that there exists an approximate factorization

$$\mathbf{X} = \mathbf{W}_0 \mathbf{H}_0 + \mathbf{Z}, \tag{3.1}$$

where $\mathbf{W}_0 \in \mathbb{R}^{n \times r}$ is a matrix of weights (with rows $\mathbf{w}_{0,i} \in \Delta^r$), $\mathbf{H}_0 \in \mathbb{R}^{r \times d}$ is a matrix of archetypes (with rows $\mathbf{h}_{0,\ell}$), and we set $\mathbf{X}_0 = \mathbf{W}_0 \mathbf{H}_0$. The deviation \mathbf{Z} is arbitrary, with rows \mathbf{z}_i satisfying $\max_{i \leq n} \|\mathbf{z}_i\|_2 \leq \delta$. We will assume throughout r to be known.

We will quantify estimation error by the sum of distances between the true archetypes and the closest estimated archetypes

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \equiv \sum_{\ell=1}^r \min_{\ell' \leq r} D(\mathbf{h}_{0,\ell}, \hat{\mathbf{h}}_{\ell'}). \tag{3.2}$$

In words, if $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})$ is small, then for each true archetype $\mathbf{h}_{0,\ell}$ there exists an estimated archetype $\hat{\mathbf{h}}_{\ell'}$ that is close to it in D -loss. Unless two or more of the true archetypes are close to each other, this means that there is a one-to-one correspondence between estimated archetypes and true archetypes, with small errors.

Assumption (Uniqueness). *We say that the factorization $\mathbf{X}_0 = \mathbf{W}_0 \mathbf{H}_0$ satisfies uniqueness with parameter $\alpha > 0$ (equivalently, is α -unique) if for all $\mathbf{H} \in \mathcal{Q}^r$ with $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\mathbf{H})$, we have*

$$\mathcal{D}(\mathbf{H}, \mathbf{X}_0)^{1/2} \geq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \alpha \{ \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \}. \tag{3.3}$$

The rationale for this assumption is quite clear. Assume that the data lie in the convex hull of the true archetypes \mathbf{H}_0 , and hence Eq. (3.1) holds without error term $\mathbf{Z} = 0$, i.e. $\mathbf{X} = \mathbf{X}_0$. We reconstruct the archetypes by demanding $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\mathbf{H})$: any such \mathbf{H} is a plausible explanation of the data. In order to make the problem well specified, we define \mathbf{H}_0 to be the matrix of archetypes that are the closest to \mathbf{X}_0 , and hence $\mathcal{D}(\mathbf{H}, \mathbf{X}_0) \geq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)$ for all \mathbf{H} . In order for the reconstruction to be unique (and hence for the problem to be identifiable) we need to assume $\mathcal{D}(\mathbf{H}, \mathbf{X}_0) > \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)$ strictly for $\mathbf{H} \neq \mathbf{H}_0$. The uniqueness assumption provides a quantitative version of this condition.

Remark 3.1. Given $\mathbf{X}_0, \mathbf{H}_0$, the best constant α such that Eq. (3.3) holds for all \mathbf{H} is a geometric property that depend on \mathbf{X}_0 only through $\text{conv}(\mathbf{X}_0)$. In particular, if $\mathbf{X}_0 = \mathbf{W}_0\mathbf{H}_0$ is a separable factorization, then it satisfies uniqueness with parameter $\alpha = 1$. Indeed in this case $\text{conv}(\mathbf{H}_0) = \text{conv}(\mathbf{X}_0)$, whence $\mathcal{D}(\mathbf{H}, \mathbf{X}_0) = \mathcal{D}(\mathbf{H}, \mathbf{H}_0)$ and $\mathcal{D}(\mathbf{H}_0, \mathbf{X}_0) = \mathcal{D}(\mathbf{H}_0, \mathbf{H}) = 0$.

It is further possible to show that $\alpha \in [0, 1]$ for all $\mathbf{H}_0, \mathbf{X}_0$. Indeed, we took \mathbf{H}_0 to be the matrix of archetypes that are closest to \mathbf{X}_0 . In other words, $\mathcal{D}(\mathbf{H}, \mathbf{X}_0) \geq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)$ and hence, $\alpha \geq 0$. In addition, since $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\mathbf{H}_0)$, for \mathbf{h}_i an arbitrary row of \mathbf{H} we have

$$\mathcal{D}(\mathbf{h}_i, \mathbf{X}_0) \leq \mathcal{D}(\mathbf{h}_i, \mathbf{H}_0). \quad (3.4)$$

Hence, $\mathcal{D}(\mathbf{H}, \mathbf{X}_0) \leq \mathcal{D}(\mathbf{H}, \mathbf{H}_0)$ and therefore

$$\mathcal{D}(\mathbf{H}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \{\mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2}\}. \quad (3.5)$$

Thus, $\alpha \leq 1$.

We say that the convex hull $\text{conv}(\mathbf{X}_0)$ has *internal radius* (at least) μ if it contains an $r - 1$ -dimensional ball of radius μ , i.e. if there exists $\mathbf{z}_0 \in \mathbb{R}^d$, $\mathbf{U} \in \mathbb{R}^{d \times (r-1)}$, with $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_d$, such that $\mathbf{z}_0 + \mathbf{U}\mathbf{B}_{r-1}(\mu) \subseteq \text{conv}(\mathbf{X}_0)$. We further denote by $\kappa(\mathbf{M})$ the condition number of matrix \mathbf{M} .

Theorem 1. Assume $\mathbf{X} = \mathbf{W}_0\mathbf{H}_0 + \mathbf{Z}$ where the factorization $\mathbf{X}_0 = \mathbf{W}_0\mathbf{H}_0$ satisfies the uniqueness assumption with parameter $\alpha > 0$, and that $\text{conv}(\mathbf{X}_0)$ has internal radius $\mu > 0$. Consider the estimator $\widehat{\mathbf{H}}$ defined by Eq. (2.5), with $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ (square loss) and $\delta = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2$. If

$$\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha\mu}{30r^{3/2}}, \quad (3.6)$$

then, we have

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_*^2 r^5}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (3.7)$$

where C_* is a coefficient that depends uniquely on the geometry of $\mathbf{H}_0, \mathbf{X}_0$, namely $C_* = 120(\sigma_{\max}(\mathbf{H}_0)/\mu) \cdot \max(1, \kappa(\mathbf{H}_0)/\sqrt{r})$.

4 Algorithms

While our main focus is on structural properties of non-negative matrix factorization, we provide evidence that the optimization problem we defined can be solved in practical scenarios. A more detailed study is left to future work.

From a computational point of view, the Lagrangian formulation (2.4) is more appealing. For the sake of simplicity, we denote the regularized risk by

$$\mathcal{R}_\lambda(\mathbf{H}) \equiv \mathcal{D}(\mathbf{X}; \mathbf{H}) + \lambda \mathcal{D}(\mathbf{H}; \mathbf{X}), \quad (4.1)$$

and leave implicit the dependence on the data \mathbf{X} . Notice that this function is non-convex and indeed has multiple global minima: in particular, permuting the rows of a minimizer \mathbf{H} yields other minimizers. We will describe two greedy optimization algorithms: one based on gradient descent, and one on alternating minimization, which generalizes the algorithm of [CB94]. In both cases it is helpful to use a good initialization: two initialization methods are introduced in the next section.

4.1 Initialization

We experimented with two initialization methods, described below.

(1) *Spectral initialization.* Under the assumption that the archetypes $\{\mathbf{h}_{0,\ell}\}_{\ell \leq r}$ are linearly independent (and for non-degenerate weights \mathbf{W}), the ‘noiseless’ matrix \mathbf{X}_0 has rank exactly r . This motivates the following approach. We compute the singular value decomposition $\mathbf{X} = \sum_{i=1}^{n \wedge d} \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{n \wedge d}$, and initialize $\widehat{\mathbf{H}}$ as the matrix $\widehat{\mathbf{H}}^{(0)}$ with rows $\hat{\mathbf{h}}_1^{(0)} = \mathbf{v}_1, \dots, \hat{\mathbf{h}}_r^{(0)} = \mathbf{v}_r$.

(2) *Successive projections initialization.* We initialize $\widehat{\mathbf{H}}^{(0)}$ by choosing archetypes $\{\hat{\mathbf{h}}_\ell^{(0)}\}_{1 \leq \ell \leq r}$ that are a subset of the data $\{\mathbf{x}_i\}_{1 \leq i \leq n}$, selected as follows. The first archetype $\hat{\mathbf{h}}_1^{(0)}$ is the data point which is farthest from the origin. For each subsequent archetype, we choose the point that is farthest from the affine subspace spanned by the previous ones.

ARCHETYPE INITIALIZATION ALGORITHM

Input : Data $\{\mathbf{x}_i\}_{i \leq n}$, $\mathbf{x}_i \in \mathbb{R}^d$; integer r ;

Output : Initial archetypes $\{\hat{\mathbf{h}}_\ell^{(0)}\}_{1 \leq \ell \leq r}$;

1: Set $i(1) = \arg \max\{D(\mathbf{x}_i; \mathbf{0}) : i \leq n\}$;

2: Set $\hat{\mathbf{h}}_1^{(0)} = \mathbf{x}_{i(1)}$;

3: For $\ell \in \{1, \dots, r\}$

4: Define $V_\ell \equiv \text{aff}(\hat{\mathbf{h}}_1^{(0)}, \hat{\mathbf{h}}_2^{(0)}, \dots, \hat{\mathbf{h}}_\ell^{(0)})$;

5: Set $i(\ell + 1) = \arg \max\{\mathcal{D}(\mathbf{x}_i; V_\ell) : i \leq n\}$;

6: Set $\hat{\mathbf{h}}_{\ell+1}^{(0)} = \mathbf{x}_{i(\ell+1)}$;

7: End For;

8: Return $\{\hat{\mathbf{h}}_\ell^{(0)}\}_{1 \leq \ell \leq r}$

This coincides with the successive projections algorithm of [ASG⁺01], with the minor difference that V_ℓ is the affine subspace spanned by the first ℓ vectors, instead of the linear subspace³ This method can be proved to return the exact archetypes if data are separable the archetypes are affine independent [AGH⁺13, GV14]. When data are not separable it provides nevertheless a good initial assignment.

4.2 Proximal alternating linearized minimization

The authors of [BST14] develop a proximal alternating linearized minimization algorithm (PALM) to solve the problems of the form

$$\text{minimize} \quad \Psi(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + g(\mathbf{y}) + h(\mathbf{x}, \mathbf{y}) \quad (4.2)$$

where $f : \mathbb{R}^m \rightarrow (-\infty, +\infty]$ and $g : \mathbb{R}^n \rightarrow (-\infty, \infty]$ are lower semicontinuous and $h \in C^1(\mathbb{R}^m \times \mathbb{R}^n)$. PALM is guaranteed to converge to critical points of the function Ψ [BST14].

We apply this algorithm to minimize the cost function (4.1), with $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ which we write as

$$\mathcal{R}_\lambda(\mathbf{H}) = \min_{\mathbf{W}} \Psi(\mathbf{H}, \mathbf{W}) = f(\mathbf{H}) + g(\mathbf{W}) + h(\mathbf{H}, \mathbf{W}). \quad (4.3)$$

³The same modification is also used in [AGH⁺13], but we do not apply the full algorithm of this paper.

where,

$$f(\mathbf{H}) = \lambda \mathcal{D}(\mathbf{H}, \mathbf{X}), \quad (4.4)$$

$$g(\mathbf{W}) = \sum_{i=1}^n \mathbf{I}(\mathbf{w}_i \in \Delta^r), \quad (4.5)$$

$$h(\mathbf{H}, \mathbf{W}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2. \quad (4.6)$$

In above equations \mathbf{w}_i are the rows of \mathbf{W} and the indicator function $\mathbf{I}(\mathbf{x} \in \Delta^r)$ is equal to zero if $\mathbf{x} \in \Delta^r$ and is equal to infinity otherwise.

By using this decomposition, the iterations of the PALM iteration reads

$$\widetilde{\mathbf{H}}^k = \mathbf{H}^k - \frac{1}{\gamma_1^k} (\mathbf{W}^k)^\top (\mathbf{W}^k \mathbf{H}^k - \mathbf{X}), \quad (4.7)$$

$$\mathbf{H}^{k+1} = \widetilde{\mathbf{H}}^k - \frac{\lambda}{\lambda + \gamma_1^k} \left(\widetilde{\mathbf{H}}^k - \Pi_{\text{conv}(\mathbf{X})}(\widetilde{\mathbf{H}}^k) \right), \quad (4.8)$$

$$\mathbf{W}^{k+1} = \Pi_{\Delta^r} \left(\mathbf{W}^k - \frac{1}{\gamma_2^k} (\mathbf{W}^k \mathbf{H}^{k+1} - \mathbf{X}) (\mathbf{H}^{k+1})^\top \right), \quad (4.9)$$

where γ_1^k, γ_2^k are step sizes and, for $\mathbf{M} \in \mathbb{R}^{m_1 \times m_2}$, and $\mathcal{S} \subseteq \mathbb{R}^{m_2}$ a closed convex set, $\Pi_{\mathcal{S}}(\mathbf{M})$ is the matrix obtained by projecting the rows of \mathbf{M} onto the simplex \mathcal{S} .

Proposition 4.1. *Consider the risk (4.1), with loss $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$, and the corresponding cost function $\Psi(\mathbf{H}, \mathbf{W})$. If the step sizes are chosen such that $\gamma_1^k > \|\mathbf{W}^{k\top} \mathbf{W}^k\|_F$, $\gamma_2^k > \max \left\{ \|\mathbf{H}^{k+1} \mathbf{H}^{k+1\top}\|_F, \varepsilon \right\}$ for some constant $\varepsilon > 0$, then $(\mathbf{H}^k, \mathbf{W}^k)$ converges to a stationary point of the function $\Psi(\mathbf{H}, \mathbf{W})$.*

The proof of this statement is deferred to Appendix D.

It is also useful to notice that the gradient of $\mathcal{R}_\lambda(\mathbf{H})$ can be computed explicitly (this can be useful to devise a stopping criterion).

Proposition 4.2. *Consider the risk (4.1), with loss $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$, and assume that the rows of \mathbf{H} are affine independent. Then, \mathcal{R}_λ is differentiable at \mathbf{H} with gradient*

$$\nabla \mathcal{R}_\lambda(\mathbf{H}) = 2 \sum_{i=1}^n \alpha_i^* (\Pi_{\text{conv}(\mathbf{H})}(\mathbf{x}_i) - \mathbf{x}_i) + 2\lambda (\mathbf{H} - \Pi_{\text{conv}(\mathbf{X})}(\mathbf{H})), \quad (4.10)$$

$$\alpha_i^* = \arg \min_{\alpha \in \Delta^r} \left\| \mathbf{H}^\top \alpha - \mathbf{x}_i^\top \right\|_2. \quad (4.11)$$

where we recall that $\Pi_{\text{conv}(\mathbf{X})}(\mathbf{H})$ denotes the matrix with rows $\Pi_{\text{conv}(\mathbf{X})}(\mathbf{H}_{1,\cdot}), \dots, \Pi_{\text{conv}(\mathbf{X})}(\mathbf{H}_{r,\cdot})$.

The proof of this proposition is given in Appendix C. Appendix E also discusses two alternative algorithms.

4.3 Numerical experiments

We implemented both the PALM algorithm described in the previous section, and the two algorithms described in Appendix E. The outcomes are generally similar.

Figures 3 and 4 repeat the experiment already described in the introduction. We generate $n = 250$ convex combinations of $r = 4$ spectra $\mathbf{h}_{0,1}, \dots, \mathbf{h}_{0,4} \in \mathbb{R}^d$, $d = 87$, this time adding white

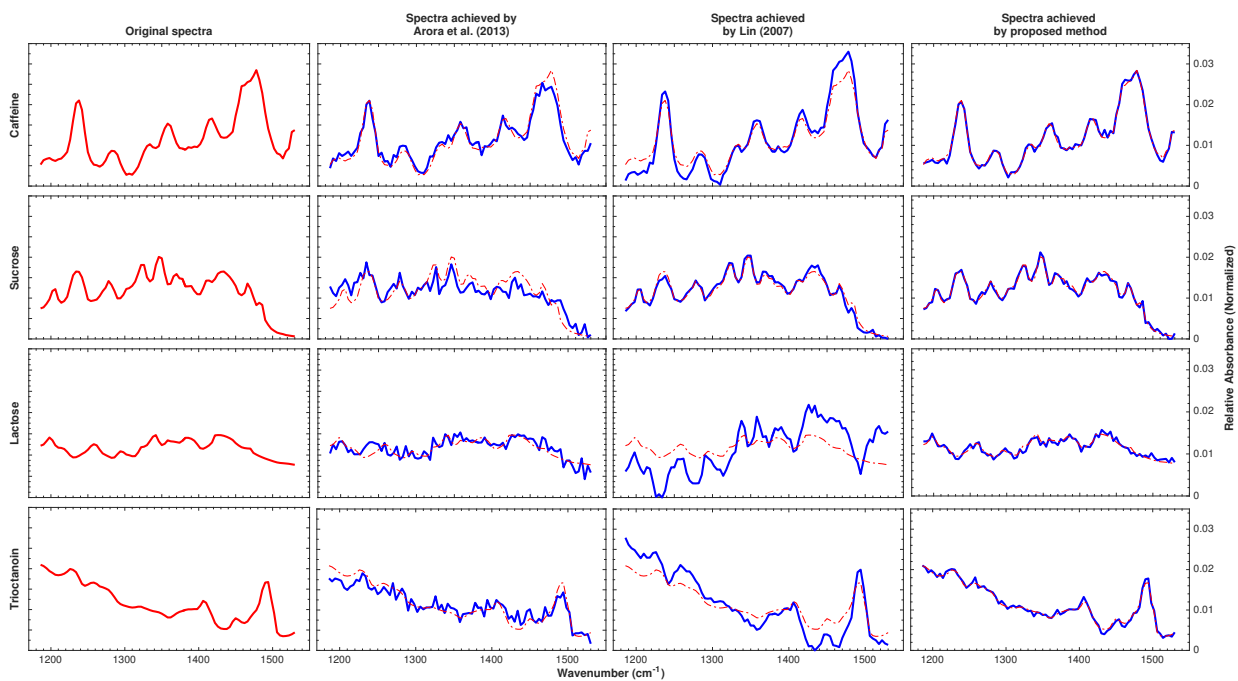


Figure 3: Reconstructing infrared spectra of four molecules, from noisy random convex combinations. Noise level $\sigma = 10^{-3}$. Left column: original spectra. The other columns correspond to different reconstruction methods.

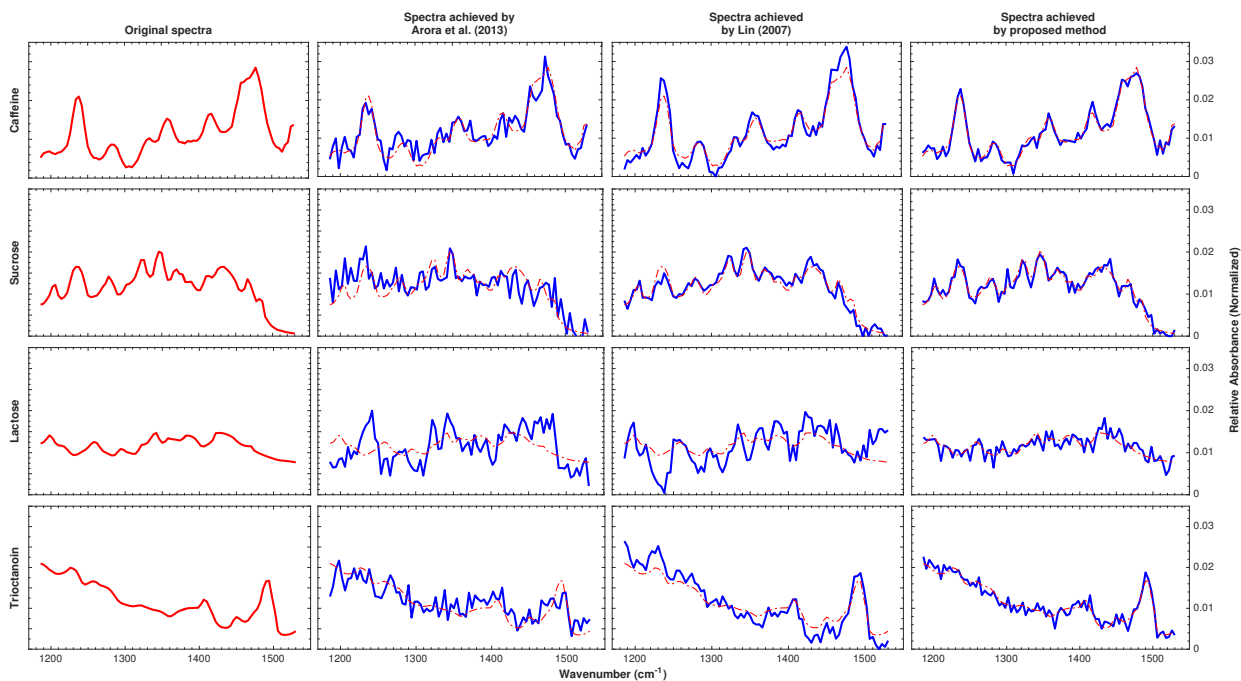


Figure 4: As in Figure 3, with $\sigma = 2 \cdot 10^{-3}$ (in blue).

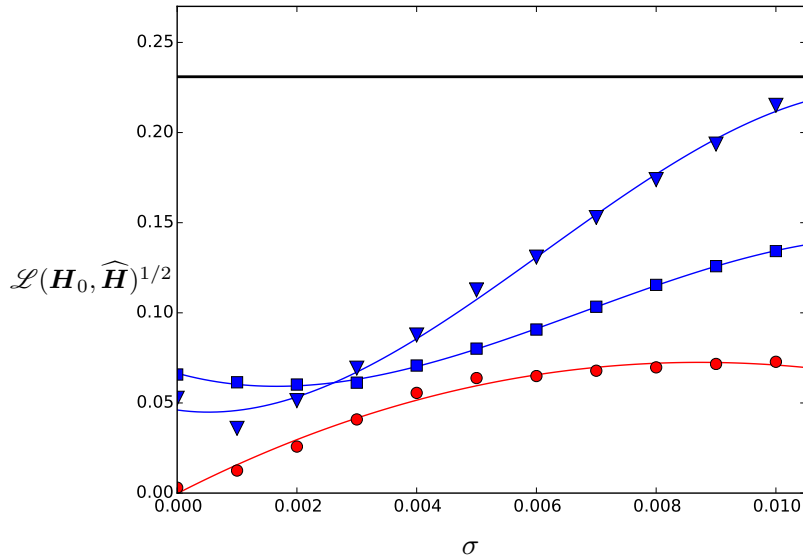


Figure 5: Risk $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}$ vs σ for different reconstruction methods. Triangles (blue): anchor words algorithm from [AGH⁺13]. Squares (blue): minimizing the objective function (1.2) using the projected gradient algorithm of [Lin07]. Circles (red): archetypal reconstruction approach in this paper. Interpolating lines are just guides for the eye. The thick horizontal line corresponds to the trivial estimator $\widehat{\mathbf{H}} = 0$.

Gaussian noise with variance σ^2 . We minimize the Lagrangian $\mathcal{R}_\lambda(\mathbf{H})$, with⁴ $\lambda = 4$ (for Figure 3) and $\lambda = 0.8$ (for Figure 4). The reconstructed spectra of the pure analytes appear to be accurate and robust to noise.

In Figure 5 we repeated the same experiment systematically for 10 noise realizations for each noise level σ , and report the resulting average loss. Among various reconstruction methods, the approach described in this paper seem to have good robustness to noise and achieves exact reconstruction as $\sigma \rightarrow 0$.

5 Discussion

We introduced a new optimization formulation of the non-negative matrix factorization problem. In its Lagrangian formulation, our approach consists in minimizing the cost function $\mathcal{R}_\lambda(\mathbf{H})$ defined in Eq. (4.1). This encompasses applications in which only one of the factors is required to be non-negative. A special case of this formulation ($\lambda \rightarrow \infty$) corresponds to the ‘archetypal analysis’ of [CB94]. In this case, the archetype estimates coincide with a subset of the data points, which is appropriate only under the separability assumption of [DS03].

Our main technical result (Theorem 1) is a robustness guarantee for the reconstructed archetypes. This holds under the uniqueness assumption, which appears to hold for generic geometries of the dataset. In particular, while separability implies uniqueness (with optimal constant $\alpha = 1$), uniqueness holds for non-separable data as well. To the best of our knowledge, similar robustness results have been obtained under separability [RRTB12, AGH⁺13, GV14, GV15] (albeit these works ob-

⁴These values were chosen as to approximately minimize the estimation error.

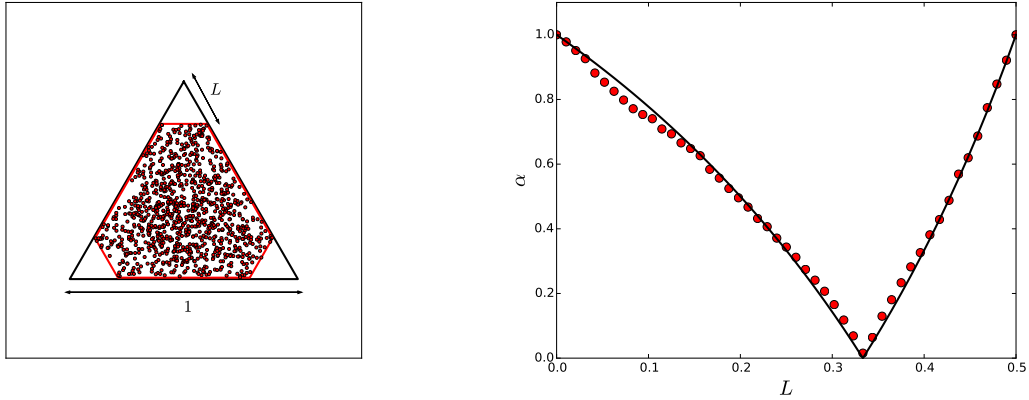


Figure 6: Numerical computation of the uniqueness parameter α . Left: data geometry. The red hexagon corresponds to $\text{conv}(\mathbf{X})$, and the black (equilateral) triangle to the archetypes \mathbf{H}_0 , for $L < 1/3$. For $L = 1/3$ the archetypes are not unique, and for $L \in (1/3, 1/2]$, they are given by an equilateral triangle rotated by $\pi/3$ (pointing down). Right: numerical evaluation of the uniqueness constant (red circles). The continuous line corresponds to an analytical upper bound (triangle rotated by $\pi/3$ with respect to \mathbf{H}_0).

tain a better dependence on r). The only exception is the recent work of Ge and Zou [GZ15] who prove robustness under a ‘subset separability’ condition, which provides a significant relaxation of separability. Under this condition, [GZ15] develops a polynomial-time algorithm to estimate the archetypes by identifying and intersecting the faces of $\text{conv}(\mathbf{H}_0)$. However, the algorithm of [GZ15] exploits collinearities to identify the faces, and this requires additional ‘genericity’ assumptions.

Admittedly, the uniqueness constant α is difficult to evaluate analytically, even for simple geometries of the data. However, by definition it does not vanish except in the case of multiple minimizers, and we expect it typically to be of order one. Figure 6 illustrate this point by computing numerically α for a simple one-parameter family of geometries with $r = 3$, $d = 2$. The parameter α vanishes at a single point, corresponding to a degenerate problem with multiple solutions.

Finally, several earlier works addressed the non-uniqueness problem in classical non-negative matrix factorization. Among others, Miao and Qi [MQ07] penalize a matrix of archetypes \mathbf{H} by the corresponding volume. Closely related to our work is the approach of Mørup and Hansen [MH12] tha also builds on archetypal analysis. To the best of our knowledge, none of these works establishes robustness of the proposed methods.

We conclude by mentioning three important problems that are not addressed by this paper: (1) Are there natural condition under which the risk function $\mathcal{R}_\lambda(\mathbf{H})$ of Eq. (4.1) can be optimized in polynomial time? We only provided an algorithm that is guaranteed to converge to a critical point. (2) We assumed the rank r to be known. In practice it will need to be estimated from the data. (3) Similarly, the regularization parameter λ should be chosen from data.

Acknowledgements

This work was partially supported by the NSF grant CCF-1319979 and a Stanford Graduate Fellowship.

References

- [AGH⁺13] Sanjeev Arora, Rong Ge, Yonatan Halpern, David M Mimno, Ankur Moitra, David Sontag, Yichen Wu, and Michael Zhu, *A practical algorithm for topic modeling with provable guarantees.*, ICML (2), 2013, pp. 280–288.
- [AGKM12] Sanjeev Arora, Rong Ge, Ravindran Kannan, and Ankur Moitra, *Computing a non-negative matrix factorization—provably*, Proceedings of the forty-fourth annual ACM symposium on Theory of computing, ACM, 2012, pp. 145–162.
- [ASG⁺01] Mário César Ugulino Araújo, Teresa Cristina Bezerra Saldanha, Roberto Kawakami Harrop Galvao, Takashi Yoneyama, Henrique Caldas Chame, and Valeria Visani, *The successive projections algorithm for variable selection in spectroscopic multicomponent analysis*, Chemometrics and Intelligent Laboratory Systems **57** (2001), no. 2, 65–73.
- [BST14] Jérôme Bolte, Shoham Sabach, and Marc Teboulle, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Mathematical Programming **146** (2014), no. 1-2, 459–494.
- [CB94] Adele Cutler and Leo Breiman, *Archetypal analysis*, Technometrics **36** (1994), no. 4, 338–347.
- [DS03] David L Donoho and Victoria Stodden, *When does non-negative matrix factorization give a correct decomposition into parts?*
- [GV14] Nicolas Gillis and Stephen A Vavasis, *Fast and robust recursive algorithms for separable nonnegative matrix factorization*, IEEE transactions on pattern analysis and machine intelligence **36** (2014), no. 4, 698–714.
- [GV15] ———, *Semidefinite programming based preconditioning for more robust near-separable nonnegative matrix factorization*, SIAM Journal on Optimization **25** (2015), no. 1, 677–698.
- [GZ15] Rong Ge and James Zou, *Intersecting faces: Non-negative matrix factorization with new guarantees*, Proceedings of the 32nd International Conference on Machine Learning (ICML-15), 2015, pp. 2295–2303.
- [KSD08] Dongmin Kim, Suvrit Sra, and Inderjit S Dhillon, *Fast projection-based methods for the least squares nonnegative matrix approximation problem*, Statistical Analysis and Data Mining **1** (2008), no. 1, 38–51.
- [Lin07] Chih-Jen Lin, *Projected gradient methods for nonnegative matrix factorization*, Neural computation **19** (2007), no. 10, 2756–2779.
- [LM] P.J. Linstrom and W.G. Mallard (eds.), *Nist chemistry webbook, nist standard reference database number 69*, National Institute of Standards and Technology, Gaithersburg MD, 20899, <http://webbook.nist.gov>, (retrieved January 5, 2017).
- [LS99] Daniel D Lee and H Sebastian Seung, *Learning the parts of objects by non-negative matrix factorization*, Nature **401** (1999), no. 6755, 788–791.

- [LS01] ———, *Algorithms for non-negative matrix factorization*, Advances in neural information processing systems, 2001, pp. 556–562.
- [MBDC⁺14] Wing-Kin Ma, José M Bioucas-Dias, Tsung-Han Chan, Nicolas Gillis, Paul Gader, Antonio J Plaza, ArulMurugan Ambikapathi, and Chong-Yung Chi, *A signal processing perspective on hyperspectral unmixing: Insights from remote sensing*, IEEE Signal Processing Magazine **31** (2014), no. 1, 67–81.
- [MH12] Morten Mørup and Lars Kai Hansen, *Archetypal analysis for machine learning and data mining*, Neurocomputing **80** (2012), 54–63.
- [MN13] Boris S Mordukhovich and Nguyen Mau Nam, *An easy path to convex analysis and applications*, Synthesis Lectures on Mathematics and Statistics **6** (2013), no. 2, 1–218.
- [MQ07] Lidan Miao and Hairong Qi, *Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization*, IEEE Transactions on Geoscience and Remote Sensing **45** (2007), no. 3, 765–777.
- [Paa97] Pentti Paatero, *Least squares formulation of robust non-negative factor analysis*, Chemometrics and intelligent laboratory systems **37** (1997), no. 1, 23–35.
- [PT94] Pentti Paatero and Unto Tapper, *Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values*, Environmetrics **5** (1994), no. 2, 111–126.
- [RCP09] Jonathan L Roux, Alain D Cheveigné, and Lucas C Parra, *Adaptive template matching with shift-invariant semi-nmf*, Advances in neural information processing systems, 2009, pp. 921–928.
- [RRTB12] Ben Recht, Christopher Re, Joel Tropp, and Victor Bittorf, *Factoring nonnegative matrices with linear programs*, Advances in Neural Information Processing Systems, 2012, pp. 1214–1222.
- [WL10] Fei Wang and Ping Li, *Efficient nonnegative matrix factorization with random projections*, Proceedings of the 2010 SIAM International Conference on Data Mining, SIAM, 2010, pp. 281–292.
- [XLG03] Wei Xu, Xin Liu, and Yihong Gong, *Document clustering based on non-negative matrix factorization*, Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval, ACM, 2003, pp. 267–273.
- [Zie12] Günter M Ziegler, *Lectures on polytopes*, vol. 152, Springer Science & Business Media, 2012.

A Further details on numerical experiments

The data in Figures 1, 3, and 4 were generated as follows. We retrieved infrared reflection spectra of caffeine, sucrose, lactose and trioctanoin from the NIST Chemistry WebBook dataset [LM]. We restricted these spectra to the wavenumbers between 1186 cm^{-1} and 1530 cm^{-1} , and denote by $\mathbf{h}_{0,1}, \dots, \mathbf{h}_{0,4} \in \mathbb{R}^d$, $d = 87$ the vector representations of these spectra. We then generate data $\mathbf{x}_i \in \mathbb{R}^d$, $i \leq n = 250$ by letting

$$\mathbf{x}_i = \sum_{\ell=1}^4 w_{i,\ell} \mathbf{h}_\ell + \mathbf{z}_i, \quad (\text{A.1})$$

where $\mathbf{z}_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ are i.i.d. Gaussian noise vectors. The weights $\mathbf{w}_i = (w_{i,\ell})_{\ell \leq 4}$ were generated as follows. The weight vectors $\{\mathbf{w}_i\}_{1 \leq i \leq 9}$ are generated such that they have 2 nonzero entries. In other words, 9 data points are on one dimensional facets of the polytope generated by $\mathbf{h}_{0,1}, \dots, \mathbf{h}_{0,4}$. In order to randomly generate these weight vectors, for each $1 \leq i \leq 9$, a pair of indices (ℓ_1, ℓ_2) between 1 and 4 is chosen uniformly at random. Then $\{\tilde{\mathbf{w}}\}_{1 \leq i \leq 9}$, $\tilde{\mathbf{w}} \in \mathbb{R}^2$ are generated as independent Dirichlet random vectors with parameter $(5, 5)$. Then we let $w_{i,\ell_1} = \tilde{w}_{i,1}$ and $w_{i,\ell_2} = \tilde{w}_{i,2}$ for $1 \leq i \leq 9$. The weight vectors $\{\mathbf{w}_i\}_{10 \leq i \leq 20}$ each have 3 nonzero entries. Similar to above, for each of these weight vectors a 3-tuple of indices (ℓ_1, ℓ_2, ℓ_3) between 1 and 4 is chosen uniformly at random. Then we let $w_{i,\ell_1} = \tilde{w}_{i,1}$, $w_{i,\ell_2} = \tilde{w}_{i,2}$, $w_{i,\ell_3} = \tilde{w}_{i,3}$ for $10 \leq i \leq 20$, where $\{\tilde{\mathbf{w}}\}_{10 \leq i \leq 20}$, $\tilde{\mathbf{w}} \in \mathbb{R}^3$ are i.i.d. Dirichlet random vectors with parameter $(5, 5, 5)$. The rest of the weight vectors have cardinality equal to 4. Hence, for $21 \leq i \leq 250$, \mathbf{w}_i are generated as i.i.d. Dirichlet random vectors with parameter $(5, 5, 5, 5)$.

B Proof of Theorem 1

In this appendix we prove Theorem 1. We start by recalling some notations already defined in the main text, and introducing some new ones. We will then state a stronger form of the theorem (with better dependence on the problem geometry in some regimes). Finally, we will present the actual proof.

Throughout this appendix, we assume the square loss $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$.

B.1 Notations and definitions

We use bold capital letters (e.g. $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$) for matrices, bold lower case for vectors (e.g. $\mathbf{x}, \mathbf{y}, \dots$) and plain lower case for scalars (a, b, c and so on). In particular, $\mathbf{e}_i \in \mathbb{R}^d$ denotes the i 'th vector in the canonical basis, $\mathbf{E}^{r,d} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r\}$ and for $r \leq d$, $\mathbf{E}_{r,d} \in \{0, 1\}^{r \times d}$ is the matrix whose i 'th column is \mathbf{e}_i , and whose columns after the r -th one are equal to $\mathbf{0}$. For a matrix \mathbf{X} , $\mathbf{X}_{i,\cdot}$ and $\mathbf{X}_{\cdot,i}$ are its i 'th row and column, respectively.

As in the main text, we denote by Δ^m the m -dimensional standard simplex, i.e. $\Delta^m = \{\mathbf{x} \in \mathbb{R}_{\geq 0}^m, \langle \mathbf{x}, \mathbf{1} \rangle = 1\}$, where $\mathbf{1} \in \mathbb{R}^m$ is the all ones vector. For a matrix $\mathbf{H} \in \mathbb{R}^{r \times d}$, we use $\sigma_{\max}(\mathbf{H})$, $\sigma_{\min}(\mathbf{H})$ to denote its largest and smallest nonzero singular values and $\kappa(\mathbf{H}) = \sigma_{\max}(\mathbf{H})/\sigma_{\min}(\mathbf{H})$ to denote its condition number. We denote by $\text{conv}(\mathbf{H})$, $\text{aff}(\mathbf{H})$ the convex hull and the affine hull of the rows of \mathbf{H} , respectively. In other words,

$$\text{conv}(\mathbf{H}) = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x} = \mathbf{H}^\top \boldsymbol{\pi}, \boldsymbol{\pi} \in \Delta^r\}, \quad (\text{B.1})$$

$$\text{aff}(\mathbf{H}) = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x} = \mathbf{H}^\top \boldsymbol{\alpha}, \langle \mathbf{1}, \boldsymbol{\alpha} \rangle = 1\}. \quad (\text{B.2})$$

We denote by $Q_{r,n}$ is the set of r by n row stochastic matrices. Namely,

$$Q_{r,n} = \left\{ \mathbf{\Pi} \in \mathbb{R}_{\geq 0}^{r \times n} : \langle \mathbf{\Pi}_{i,\cdot}, \mathbf{1} \rangle = 1 \right\}. \quad (\text{B.3})$$

with use $Q_r \equiv Q_{r,r}$. Further, S_r is defined as

$$S_r = \{ \mathbf{\Pi} \in Q_r : \Pi_{i,j} \in \{0, 1\} \}. \quad (\text{B.4})$$

As a consequence, given $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{H}_1, \mathbf{H}_2 \in \mathbb{R}^{r \times d}$, the loss functions $\mathcal{D}(\cdot, \cdot)$ and $\mathcal{L}(\cdot, \cdot)$ take the form

$$\mathcal{D}(\mathbf{H}_1, \mathbf{X}) = \min_{\mathbf{\Pi} \in Q_{r,n}} \|\mathbf{H}_1 - \mathbf{\Pi}\mathbf{X}\|_F^2, \quad (\text{B.5})$$

$$\mathcal{L}(\mathbf{H}_1, \mathbf{H}_2) = \min_{\mathbf{\Pi} \in S_r} \|\mathbf{H}_1 - \mathbf{\Pi}\mathbf{H}_2\|_F^2. \quad (\text{B.6})$$

We use $B_m(\rho)$ to denote the closed ball with radius ρ in m dimensions, centered at 0. In addition, for $\mathbf{H} \in \mathbb{R}^{m \times d}$ we define the ρ -neighborhood of $\text{conv}(\mathbf{H})$ as

$$B_r(\rho; \mathbf{H}) := \{ \mathbf{x} \in \mathbb{R}^d : \mathcal{D}(\mathbf{x}, \mathbf{H}) \leq \rho^2 \}. \quad (\text{B.7})$$

For a convex set \mathcal{C} we denote the set of its extremal points by $\text{ext}(\mathcal{C})$ and the projection of a point $\mathbf{x} \in \mathbb{R}^d$ onto \mathcal{C} by $\mathbf{\Pi}_{\mathcal{C}}(\mathbf{x})$. Namely,

$$\mathbf{\Pi}_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\|_2. \quad (\text{B.8})$$

Also, for a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, and a mapping (not necessarily linear) $\mathbf{P} : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $\mathbf{P}(\mathbf{X}) \in \mathbb{R}^{n \times d}$ is the matrix whose i 'th row is $\mathbf{P}(\mathbf{X}_{i,\cdot})$.

B.2 Theorem statement

The statement below provides more detailed result with respect to the one in Theorem 1.

Theorem 2. *Assume $\mathbf{X} = \mathbf{W}_0\mathbf{H}_0 + \mathbf{Z}$ where the factorization $\mathbf{X}_0 = \mathbf{W}_0\mathbf{H}_0$ satisfies the uniqueness assumption with parameter $\alpha > 0$, and that $\text{conv}(\mathbf{X}_0)$ has internal radius $\mu > 0$. Consider the estimator $\widehat{\mathbf{H}}$ defined by Eq. (2.5), with $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ (square loss) and $\delta = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2$. If*

$$\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha\mu}{30r^{3/2}}, \quad (\text{B.9})$$

then, setting $\delta = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2$ in the problem (2.5) we get

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_*^2 r^5}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (\text{B.10})$$

where C_* is a coefficient that depends uniquely on the geometry of $\mathbf{H}_0, \mathbf{X}_0$, namely $C_* = 120(\sigma_{\max}(\mathbf{H}_0)/\mu) \cdot \max(1, \kappa(\mathbf{H}_0)/\sqrt{r})$.

Further, if

$$\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha\mu}{330\kappa(\mathbf{H}_0)r^{5/2}}, \quad (\text{B.11})$$

then, setting $\delta = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2$ in the problem (2.5) we get

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_{**}^2 r^4}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (\text{B.12})$$

where $C_{**} = 120 \max(\kappa(\mathbf{H}_0), (\sigma_{\max}(\mathbf{H}_0)/r + \|\mathbf{z}_0\|_2)/(\mu r^{1/2})) \cdot \max(1, \kappa(\mathbf{H}_0)/\sqrt{r})$.

B.3 Proof

B.3.1 Lemmas

Lemma B.1. *Let \mathcal{R} be a convex set and \mathcal{C} be a convex cone. Define*

$$\gamma_{\mathcal{C}} = \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{v} \in \mathcal{C}, \|\mathbf{v}\|_2=1} \langle \mathbf{u}, \mathbf{v} \rangle. \quad (\text{B.13})$$

We have

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 + (1 + \gamma_{\mathcal{C}}) \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2 \geq \gamma_{\mathcal{C}} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x}\|_2. \quad (\text{B.14})$$

An illustration of this lemma in the case of $\mathcal{R} \subset \mathcal{C}$ is given in Figure 7. Note that, $\gamma_{\mathcal{C}}$ measures the pointedness of the cone \mathcal{C} . Geometrically (for $\mathcal{R} \subseteq \mathcal{C}$) the lemma states that the cosine of the angle between $\arg \min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2$ and $\arg \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x}\|_2$ is smaller than $\gamma_{\mathcal{C}}$.

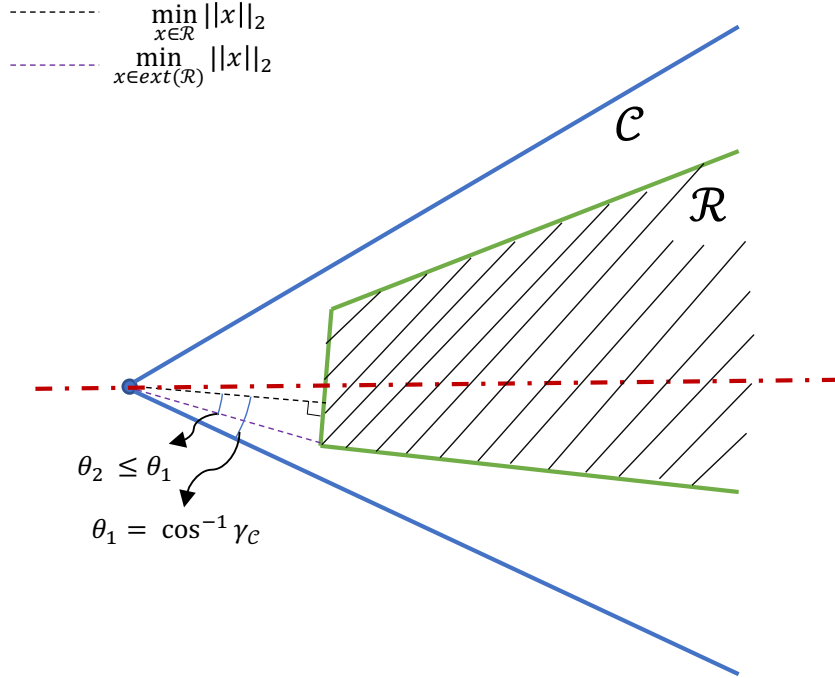


Figure 7: Picture of Lemma B.1, in the case, $\mathcal{R} \subset \mathcal{C}$.

Proof. We write

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 = \min_{\mathbf{x} \in \mathcal{R}} \max_{\|\mathbf{u}\|_2=1} \langle \mathbf{u}, \mathbf{x} \rangle \geq \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \mathcal{R}} \langle \mathbf{u}, \mathbf{x} \rangle = \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \langle \mathbf{u}, \mathbf{x} \rangle. \quad (\text{B.15})$$

Replacing

$$\mathbf{x} = \Pi_{\mathcal{C}}(\mathbf{x}) + (\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})), \quad (\text{B.16})$$

we get

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \langle \mathbf{u}, \Pi_{\mathcal{C}}(\mathbf{x}) + (\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})) \rangle \quad (\text{B.17})$$

$$\geq \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \langle \mathbf{u}, \Pi_{\mathcal{C}}(\mathbf{x}) \rangle - \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.18})$$

Hence, using the definition of $\gamma_{\mathcal{C}}$, we have

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \gamma_{\mathcal{C}} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{\Pi}_{\mathcal{C}}(\mathbf{x})\|_2 - \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \mathbf{\Pi}_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.19})$$

Note that

$$\|\mathbf{\Pi}_{\mathcal{C}}(\mathbf{x})\|_2 \geq \|\mathbf{x}\|_2 - \|\mathbf{x} - \mathbf{\Pi}_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.20})$$

Therefore,

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \gamma_{\mathcal{C}} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x}\|_2 - (1 + \gamma_{\mathcal{C}}) \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \mathbf{\Pi}_{\mathcal{C}}(\mathbf{x})\|_2, \quad (\text{B.21})$$

and this completes the proof. \square

The next lemma is a consequence of Lemma B.1.

Lemma B.2. *Let $\mathbf{H}, \mathbf{H}_0 \in \mathbb{R}^{r \times d}$, $r \leq d$, be matrices with linearly independent rows. We have*

$$\mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \sqrt{2}\kappa(\mathbf{H}_0)\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} + (1 + \sqrt{2})\sqrt{r}\mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2}. \quad (\text{B.22})$$

Proof. Consider the cone $\mathcal{C}_1 \subset \mathbb{R}^d$, generated by vectors $\mathbf{e}_2 - \mathbf{e}_1, \dots, \mathbf{e}_r - \mathbf{e}_1 \in \mathbb{R}^d$, i.e.,

$$\mathcal{C}_1 = \left\{ \mathbf{v} \in \mathbb{R}^d; \mathbf{v} = \sum_{i=2}^r v_i(\mathbf{e}_i - \mathbf{e}_1), v_i \geq 0 \right\}. \quad (\text{B.23})$$

For $\mathbf{v} \in \mathcal{C}_1$, $\|\mathbf{v}\|_2 = 1$ we have

$$\mathbf{v} = (-\langle \mathbf{1}, \mathbf{x} \rangle, \mathbf{x}, 0, 0, \dots, 0), \quad (\text{B.24})$$

where $\mathbf{x} \in \mathbb{R}_{\geq 0}^{r-1}$ and

$$\|\mathbf{x}\|_2^2 + \langle \mathbf{1}, \mathbf{x} \rangle^2 = 1. \quad (\text{B.25})$$

Since, $\langle \mathbf{1}, \mathbf{x} \rangle = \|\mathbf{x}\|_1 \geq \|\mathbf{x}\|_2$, we get $\langle \mathbf{1}, \mathbf{x} \rangle \geq 1/\sqrt{2}$. Thus, for $\mathbf{u} = -\mathbf{e}_1$, we have $\langle \mathbf{u}, \mathbf{v} \rangle \geq 1/\sqrt{2}$. Therefore, for $\gamma_{\mathcal{C}_1}$ defined as in Lemma B.1, we have $\gamma_{\mathcal{C}_1} \geq 1/\sqrt{2}$. In addition, by symmetry, for $i \in \{1, 2, \dots, r\}$, for the cone $\mathcal{C}_i \subset \mathbb{R}^d$, generated by vectors $\mathbf{e}_1 - \mathbf{e}_i, \mathbf{e}_2 - \mathbf{e}_i, \dots, \mathbf{e}_r - \mathbf{e}_i \in \mathbb{R}^d$ we have $\gamma_{\mathcal{C}_i} = \gamma \geq 1/\sqrt{2}$. Hence, using Lemma B.1 for $\mathbf{H} \in \mathbb{R}^{r \times d}$, $\mathcal{R} = \text{conv}(\mathbf{H}) - \mathbf{e}_j$ (the set obtained by translating $\text{conv}(\mathbf{H})$ by $-\mathbf{e}_j$), $\mathcal{C} = \mathcal{C}_j$ we get for $j = 1, 2, \dots, r$

$$\min_{\mathbf{q} \in \Delta^r} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2 \geq \gamma \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2 - (1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \mathbb{R}_{\geq 0}^r} \|\mathbf{H}_{i,\cdot}^{\top} - \mathbf{e}_j - \mathbf{E}_{r,d}^{\top} \mathbf{q} + \mathbf{e}_j \langle \mathbf{1}, \mathbf{q} \rangle\|_2 \quad (\text{B.26})$$

$$\geq \gamma \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2 - (1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^{\top} - \mathbf{E}_{r,d}^{\top} \mathbf{q}\|_2. \quad (\text{B.27})$$

Hence,

$$\begin{aligned} \sum_{j=1}^r \min_{\mathbf{q} \in \Delta^r} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2^2 &\geq \gamma^2 \sum_{j=1}^r \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2^2 + (1 + \gamma)^2 r \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^{\top} - \mathbf{E}_{r,d}^{\top} \mathbf{q}\|_2^2 \\ &\quad - 2\gamma(1 + \gamma) \left(\max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^{\top} - \mathbf{E}_{r,d}^{\top} \mathbf{q}\|_2 \right) \sum_{j=1}^r \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2 \quad (\text{B.28}) \end{aligned}$$

$$\begin{aligned} &\geq \left[\gamma \left(\sum_{j=1}^r \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^{\top} \mathbf{q}\|_2^2 \right)^{1/2} \right. \\ &\quad \left. - (1 + \gamma)\sqrt{r} \left(\max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^{\top} - \mathbf{E}_{r,d}^{\top} \mathbf{q}\|_2 \right) \right]^2. \quad (\text{B.29}) \end{aligned}$$

Therefore,

$$\min_{\mathbf{Q} \in \mathcal{Q}_r} \|\mathbf{E}_{r,d} - \mathbf{QH}\|_F \geq \gamma \min_{\mathbf{Q} \in \mathcal{S}_r} \|\mathbf{E}_{r,d} - \mathbf{QH}\|_F - (1 + \gamma)\sqrt{r} \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{E}_{r,d}^\top \mathbf{q}\|_2. \quad (\text{B.30})$$

Now consider $\mathbf{H}_0 \in \mathbb{R}^{r \times d}$ where $\mathbf{H}_0 = \mathbf{E}_{r,d} \mathbf{M}$, $\mathbf{H} = \mathbf{Y} \mathbf{M}$, where $\mathbf{M} \in \mathbb{R}^{d \times d}$ is invertible. We have

$$\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} = \min_{\mathbf{Q} \in \mathcal{Q}_r} \|\mathbf{H}_0 - \mathbf{QH}\|_F = \min_{\mathbf{Q} \in \mathcal{Q}_r} \|(\mathbf{E}_{r,d} - \mathbf{QY})\mathbf{M}\|_F \quad (\text{B.31})$$

$$\geq \sigma_{\min}(\mathbf{M}) \min_{\mathbf{Q} \in \mathcal{Q}_r} \|\mathbf{E}_{r,d} - \mathbf{QY}\|_F \quad (\text{B.32})$$

$$\geq \gamma \sigma_{\min}(\mathbf{M}) \min_{\mathbf{Q} \in \mathcal{S}_r} \|\mathbf{E}_{r,d} - \mathbf{QY}\|_F - \sigma_{\min}(\mathbf{M})\sqrt{r}(1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{Y}_{i,\cdot}^\top - \mathbf{E}_{r,d}^\top \mathbf{q}\|_2 \quad (\text{B.33})$$

$$\begin{aligned} &= \gamma \sigma_{\min}(\mathbf{M}) \min_{\mathbf{Q} \in \mathcal{S}_r} \|(\mathbf{H}_0 - \mathbf{QH})\mathbf{M}^{-1}\|_F \\ &\quad - \sigma_{\min}(\mathbf{M})\sqrt{r}(1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|(\mathbf{M}^{-1})^\top (\mathbf{H}_{i,\cdot}^\top - \mathbf{H}_0^\top \mathbf{q})\|_2. \end{aligned} \quad (\text{B.34})$$

Thus, using the fact that $\sigma_{\max}(\mathbf{M})/\sigma_{\min}(\mathbf{M}) = \kappa(\mathbf{M}) = \kappa(\mathbf{H}_0)$,

$$\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \geq \frac{\gamma}{\kappa(\mathbf{H}_0)} \mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} - \frac{(1 + \gamma)\sqrt{r}}{\kappa(\mathbf{H}_0)} \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{H}_0^\top \mathbf{q}\|_2 \quad (\text{B.35})$$

$$\geq \frac{\gamma}{\kappa(\mathbf{H}_0)} \mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} - \frac{(1 + \gamma)\sqrt{r}}{\kappa(\mathbf{H}_0)} \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2}. \quad (\text{B.36})$$

Therefore,

$$\mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \frac{\kappa(\mathbf{H}_0)}{\gamma} \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} + \frac{(1 + \gamma)\sqrt{r}}{\gamma} \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2}. \quad (\text{B.37})$$

Finally, note that the function $f(x) = (1 + x)/x$ is monotone decreasing over $\mathbb{R}_{>0}$. Hence, for $\gamma \geq 1/\sqrt{2}$, $(1 + \gamma)/\gamma \leq 1 + \sqrt{2}$. Therefore, we get

$$\mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \sqrt{2}\kappa(\mathbf{H}_0) \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} + (1 + \sqrt{2})\sqrt{r} \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} \quad (\text{B.38})$$

and this completes the proof. \square

We continue with the following lemmas on the condition number of the matrix \mathbf{H} .

Lemma B.3. *Let $\mathbf{H}_0, \mathbf{H} \in \mathbb{R}^{r \times d}$, $r \leq d$, with \mathbf{H} having full row rank. We have*

$$\sigma_{\max}(\mathbf{H}) \leq \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} + \sqrt{r} \sigma_{\max}(\mathbf{H}_0), \quad (\text{B.39})$$

In addition, if

$$\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \frac{\sigma_{\min}(\mathbf{H}_0)}{2}, \quad (\text{B.40})$$

then

$$\kappa(\mathbf{H}) \leq \frac{2r \sigma_{\max}(\mathbf{H}_0) + 2 \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} \sqrt{r}}{\sigma_{\min}(\mathbf{H}_0)}. \quad (\text{B.41})$$

Further, if

$$\mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \frac{\sigma_{\min}(\mathbf{H}_0)}{6\sqrt{r}}, \quad (\text{B.42})$$

then

$$\sigma_{\max}(\mathbf{H}) \leq 2\sigma_{\max}(\mathbf{H}_0), \quad (\text{B.43})$$

$$\kappa(\mathbf{H}) \leq (7/2)\kappa(\mathbf{H}_0). \quad (\text{B.44})$$

Proof. For the sake of simplicity, we will write $\mathcal{D}_1 = \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2}$, $\mathcal{D}_2 = \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2}$. Note that using the assumptions of Lemma B.3 we have

$$\begin{aligned} \mathbf{H}_0 &= \mathbf{P}\mathbf{H} + \mathbf{A}_2; & \|\mathbf{A}_2\|_F &= \mathcal{D}_2, \\ \mathbf{H} &= \mathbf{R}\mathbf{H}_0 + \mathbf{A}_1; & \|\mathbf{A}_1\|_F &= \mathcal{D}_1, \end{aligned} \quad (\text{B.45})$$

where $\mathbf{P}, \mathbf{R} \in \mathbb{R}_{\geq 0}^{r \times r}$ are row-stochastic matrices and $\mathbf{A}_1, \mathbf{A}_2 \in \mathbb{R}^{r \times d}$. Also, $\sigma_{\max}(\mathbf{A}_1) \leq \|\mathbf{A}_1\|_F = \mathcal{D}_1$, $\sigma_{\max}(\mathbf{A}_2) \leq \|\mathbf{A}_2\|_F = \mathcal{D}_2$. Therefore,

$$\sigma_{\max}(\mathbf{P})\sigma_{\min}(\mathbf{H}) \geq \sigma_{\min}(\mathbf{P}\mathbf{H}) \geq \sigma_{\min}(\mathbf{H}_0) - \sigma_{\max}(\mathbf{A}_2) \geq \sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2. \quad (\text{B.46})$$

In addition, note that for a row stochastic matrix $\mathbf{P} \in Q_r$, we have

$$\sigma_{\max}(\mathbf{P}) \leq \|\mathbf{P}\|_F = \left(\sum_{i=1}^r \|\mathbf{P}_{i,\cdot}\|_2^2 \right)^{1/2} \leq \left(\sum_{i=1}^r \|\mathbf{P}_{i,\cdot}\|_1^2 \right)^{1/2} \leq \sqrt{r}. \quad (\text{B.47})$$

Hence, for $\mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)$ we get

$$\sigma_{\min}(\mathbf{H}) \geq \frac{\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2}{\sqrt{r}}. \quad (\text{B.48})$$

In addition,

$$\sigma_{\max}(\mathbf{H}) \leq \sigma_{\max}(\mathbf{R}\mathbf{H}_0) + \sigma_{\max}(\mathbf{A}_1) \leq \sigma_{\max}(\mathbf{R})\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_1 \leq \sqrt{r}\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_1. \quad (\text{B.49})$$

Hence, using (B.48), (B.49), for $\mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)$ we have

$$\kappa(\mathbf{H}) \leq \frac{r\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_1\sqrt{r}}{\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2}. \quad (\text{B.50})$$

Thus, for $\mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)/2$, we get Eqs. (B.39), (B.41).

Now assume that $\mathcal{D}_1 + \mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)/(6\sqrt{r})$. In this case, using (B.45) we have

$$\mathbf{H}_0 = \mathbf{P}(\mathbf{R}\mathbf{H}_0 + \mathbf{A}_1) + \mathbf{A}_2. \quad (\text{B.51})$$

Therefore,

$$(\mathbf{I} - \mathbf{P}\mathbf{R})\mathbf{H}_0 = \mathbf{P}\mathbf{A}_1 + \mathbf{A}_2, \quad (\text{B.52})$$

hence,

$$\mathbf{I} - \mathbf{P}\mathbf{R} = (\mathbf{P}\mathbf{A}_1 + \mathbf{A}_2)\mathbf{H}_0^\dagger \quad (\text{B.53})$$

and

$$\mathbf{P}\mathbf{R} = \mathbf{I} - \mathbf{P}\mathbf{A}_1\mathbf{H}_0^\dagger - \mathbf{A}_2\mathbf{H}_0^\dagger. \quad (\text{B.54})$$

where \mathbf{H}_0^\dagger is the right inverse of matrix \mathbf{H}_0 . Note that

$$\sigma_{\max}(\mathbf{H}_0^\dagger) = \sigma_{\min}(\mathbf{H}_0)^{-1}. \quad (\text{B.55})$$

By permuting the rows and columns of \mathbf{H}_0 , without loss of generality, we can assume that $R_{ii} = \|\mathbf{R}_{\cdot,i}\|_\infty$. We can write

$$R_{ii} \geq \langle \mathbf{P}_{i,\cdot}, \mathbf{R}_{\cdot,i} \rangle = 1 - (\mathbf{P}\mathbf{A}_1\mathbf{H}_0^\dagger)_{ii} - (\mathbf{A}_2\mathbf{H}_0^\dagger)_{ii} \quad (\text{B.56})$$

$$\geq 1 - \|(\mathbf{P}\mathbf{A}_1\mathbf{H}_0^\dagger)_{i,\cdot}\|_2 - \|(\mathbf{A}_2\mathbf{H}_0^\dagger)_{i,\cdot}\|_2 \quad (\text{B.57})$$

$$\geq 1 - \max_{\mathbf{u} \in \Delta^r} \|\mathbf{A}_1^\top \mathbf{u}\|_2 \sigma_{\max}(\mathbf{H}_0^\dagger) - \|(\mathbf{A}_2)_{i,\cdot}\|_2 \sigma_{\max}(\mathbf{H}_0^\dagger) \quad (\text{B.58})$$

$$\geq 1 - \max_{\mathbf{u} \in \Delta^r} \|\mathbf{u}\|_2 \sigma_{\max}(\mathbf{A}_1) \sigma_{\max}(\mathbf{H}_0^\dagger) - \|\mathbf{A}_2\|_F \sigma_{\max}(\mathbf{H}_0^\dagger) \quad (\text{B.59})$$

$$\geq 1 - \frac{\mathcal{D}_1 + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0)}. \quad (\text{B.60})$$

Hence, for all $i, j \in [r], i \neq j$, since \mathbf{R} is row-stochastic,

$$R_{ji} \leq \frac{\mathcal{D}_1 + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0)}. \quad (\text{B.61})$$

Thus,

$$\langle \mathbf{P}_{i,\cdot}, \mathbf{R}_{\cdot,i} \rangle = R_{ii}P_{ii} + \sum_{j \neq i} P_{ij}R_{ji} \leq R_{ii}P_{ii} + \left(\max_{j \neq i} R_{ji} \right) \sum_{j \neq i} P_{ij} \quad (\text{B.62})$$

$$\leq P_{ii} + \frac{\mathcal{D}_1 + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0)} (1 - P_{ii}). \quad (\text{B.63})$$

Therefore, using (B.60),

$$P_{ii} \geq \frac{\sigma_{\min}(\mathbf{H}_0) - 2(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.64})$$

Thus, we can write

$$\mathbf{P} = \mathbf{I} + \Delta; \quad \|\Delta_{i,\cdot}\|_1 \leq \frac{2(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.65})$$

Therefore,

$$\sigma_{\max}(\Delta) \leq \|\Delta\|_F = \left(\sum_{i=1}^r \|\Delta_{i,\cdot}\|_2^2 \right)^{1/2} \leq \left(\sum_{i=1}^r \|\Delta_{i,\cdot}\|_1^2 \right)^{1/2} \leq \frac{2(\mathcal{D}_1 + \mathcal{D}_2)\sqrt{r}}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.66})$$

Hence,

$$\sigma_{\max}(\mathbf{P}) \leq 1 + \frac{2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}, \quad \sigma_{\min}(\mathbf{P}) \geq 1 - \frac{2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.67})$$

From (B.45) we have $\sigma_{\min}(\mathbf{PH}) \geq \sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2$. Using $\sigma_{\min}(\mathbf{PH}) \leq \sigma_{\max}(\mathbf{P})\sigma_{\min}(\mathbf{H})$, we get

$$\sigma_{\min}(\mathbf{H}) \geq \frac{(\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2)(\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2))}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2) + 2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.68})$$

Further, from (B.45) we have $\sigma_{\max}(\mathbf{PH}) \leq \sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_2$. Using $\sigma_{\max}(\mathbf{PH}) \geq \sigma_{\min}(\mathbf{P})\sigma_{\max}(\mathbf{H})$, we get

$$\sigma_{\max}(\mathbf{H}) \leq \frac{(\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_2)(\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2))}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2) - 2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.69})$$

Hence, for $\mathcal{D}_1 + \mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)/(6\sqrt{r})$, we have $\sigma_{\max}(\mathbf{H}) \leq 35\sigma_{\max}(\mathbf{H}_0)/18 < 2\sigma_{\max}(\mathbf{H}_0)$. In addition,

$$\kappa(\mathbf{H}) \leq \left(\frac{\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2} \right) \left(1 + \frac{4\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2) - 2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)} \right) \quad (\text{B.70})$$

$$\leq \frac{6\kappa(\mathbf{H}_0) + 1}{5} \left(1 + \frac{4}{3} \right) \leq \frac{42\kappa(\mathbf{H}_0) + 7}{15} < \frac{7\kappa(\mathbf{H}_0)}{2}, \quad (\text{B.71})$$

and this completes the proof. \square

Lemma B.4. *Let $\mathbf{X}_0 = \mathbf{W}_0\mathbf{H}_0 \in \mathbb{R}^{n \times d}$ be such that $\text{conv}(\mathbf{X}_0)$ has internal radius at least $\mu > 0$, and $\mathbf{X} = \mathbf{X}_0 + \mathbf{Z}$ with $\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta$. If $\mathbf{H} \in \mathbb{R}^{r \times d}$, $\mathbf{H}_{i,\cdot} \in \text{aff}(\mathbf{H}_0)$ is feasible for problem (2.5) and has linearly independent rows, then we have*

$$\sigma_{\min}(\mathbf{H}) \geq \sqrt{2}(\mu - 2\delta). \quad (\text{B.72})$$

Proof. Let

$$\mathbf{X}'_{i,\cdot} = \Pi_{\text{conv}(\mathbf{H})}(\mathbf{X}_{i,\cdot}) \equiv \arg \min_{\mathbf{x} \in \text{conv}(\mathbf{H})} \|\mathbf{X}_{i,\cdot} - \mathbf{x}\|_2. \quad (\text{B.73})$$

Note that since \mathbf{H} is feasible for problem (2.5) and $\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta$

$$\|(\mathbf{X}_0)_{i,\cdot} - \mathbf{X}'_{i,\cdot}\|_2 \leq \|(\mathbf{X}_0)_{i,\cdot} - \mathbf{X}_{i,\cdot}\|_2 + \|\mathbf{X}_{i,\cdot} - \mathbf{X}'_{i,\cdot}\|_2 \leq 2\delta. \quad (\text{B.74})$$

Therefore, for any $\mathbf{x}_0 \in \text{conv}(\mathbf{X}_0)$, writing $\mathbf{x}_0 = \mathbf{X}_0^\top \mathbf{a}_0$, $\mathbf{a}_0 \in \Delta^n$, we have

$$\mathcal{D}(\mathbf{x}_0, \mathbf{X}')^{1/2} = \min_{\mathbf{a} \in \Delta^n} \left\| \mathbf{X}_0^\top \mathbf{a}_0 - \mathbf{X}'^\top \mathbf{a} \right\|_2 \leq \left\| \mathbf{X}_0^\top \mathbf{a}_0 - \mathbf{X}'^\top \mathbf{a}_0 \right\|_2 \quad (\text{B.75})$$

$$\leq \left(\sum_{i=1}^n (a_0)_i \right) \|(\mathbf{X}_0)_{i,\cdot} - \mathbf{X}'_{i,\cdot}\|_2 \leq 2\delta. \quad (\text{B.76})$$

Since $\text{conv}(\mathbf{X}_0)$ has internal radius at least μ , there exists $\mathbf{z}_0 \in \mathbb{R}^d$, and an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{d \times r'}$, $r' = r - 1$, such that $\mathbf{z}_0 + \mathbf{U}\mathbf{B}_{r'}(\mu) \subseteq \text{conv}(\mathbf{X}_0)$. Hence, for every $\mathbf{z} \in \mathbb{R}^{r'}$, $\|\mathbf{z}\|_2 = 1$ there exists $\mathbf{a} \in \Delta^n$ such that

$$\mu\mathbf{U}\mathbf{z} + \mathbf{z}_0 = \mathbf{X}_0^\top \mathbf{a}. \quad (\text{B.77})$$

Therefore, for any unit vector \mathbf{u} in column space of \mathbf{U} , for the line segment

$$l_{\mathbf{u},\mu} = \{\mathbf{z} : \mathbf{z} = \mathbf{z}_0 + \alpha\mathbf{u}, |\alpha| \leq \mu\} \subseteq \text{conv}(\mathbf{X}_0). \quad (\text{B.78})$$

Thus,

$$l_{\mathbf{u},\mu} \subseteq \mathbf{P}_{\mathbf{u}}(\text{conv}(\mathbf{X}_0)) \quad (\text{B.79})$$

where $\mathbf{P}_{\mathbf{u}}$ is the orthogonal projection onto the line containing $l_{\mathbf{u},\mu}$. Note that using (B.76), for any $\mathbf{x}_0 \in \text{conv}(\mathbf{X}_0)$ we have

$$\mathcal{D}(\mathbf{P}_{\mathbf{u}}(\mathbf{x}_0), \mathbf{P}_{\mathbf{u}}(\text{conv}(\mathbf{X}'))^{1/2}) \leq \mathcal{D}(\mathbf{x}_0, \mathbf{X}')^{1/2} \leq 2\delta. \quad (\text{B.80})$$

In other words, for any $\mathbf{x}_0 \in \mathbf{P}_{\mathbf{u}}(\text{conv}(\mathbf{X}_0))$, $D(\mathbf{x}_0, \mathbf{P}_{\mathbf{u}}(\text{conv}(\mathbf{X}'))) \leq 2\delta$. Therefore, using (B.78) for any \mathbf{u} in column space of \mathbf{U} , we have

$$l_{\mathbf{u},\mu-2\delta} \subseteq \mathbf{P}_{\mathbf{u}}(\text{conv}(\mathbf{X}')). \quad (\text{B.81})$$

This implies that

$$\mathbf{z}_0 + \mathbf{U}\mathbf{B}_{r'}(\mu - 2\delta) \subseteq \text{conv}(\mathbf{X}') \subseteq \text{conv}(\mathbf{H}). \quad (\text{B.82})$$

Hence, for every $\mathbf{z} \in \mathbb{R}^{r'}$, $\|\mathbf{z}\|_2 = 1$ there exists $\mathbf{a} \in \Delta^r$ such that

$$(\mu - 2\delta)\mathbf{U}\mathbf{z} + \mathbf{z}_0 = \mathbf{H}^\top \mathbf{a}. \quad (\text{B.83})$$

Note that \mathbf{H}^\top has linearly independent columns. Multiplying the previous equation by $(\mathbf{H}^\top)^\dagger$ the left inverse of \mathbf{H}^\top , we get

$$(\mu - 2\delta)(\mathbf{H}^\top)^\dagger \mathbf{U}\mathbf{z} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0 = \mathbf{a}. \quad (\text{B.84})$$

Let

$$\mathbf{a}_1 = (\mu - 2\delta)(\mathbf{H}^\top)^\dagger \mathbf{U}\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0, \quad (\text{B.85})$$

$$\mathbf{a}_2 = -(\mu - 2\delta)(\mathbf{H}^\top)^\dagger \mathbf{U}\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0, \quad (\text{B.86})$$

where \mathbf{v} is the right singular vector corresponding to the largest singular value of $(\mathbf{H}^\top)^\dagger \mathbf{U}$. Therefore, we have

$$\mathbf{a}_1 = (\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U})\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0, \quad (\text{B.87})$$

$$\mathbf{a}_2 = -(\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U})\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0. \quad (\text{B.88})$$

Thus, for $\mathbf{a}_1, \mathbf{a}_2 \in \Delta^r$

$$\|\mathbf{a}_1 - \mathbf{a}_2\|_2 = 2(\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U}). \quad (\text{B.89})$$

Note that

$$\|\mathbf{a}_1 - \mathbf{a}_2\|_2 \leq \sqrt{2}. \quad (\text{B.90})$$

Thus,

$$2(\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U}) = \frac{2(\mu - 2\delta)}{\sigma_{\min}(\mathbf{H})} \leq \sqrt{2}. \quad (\text{B.91})$$

Hence,

$$\sigma_{\min}(\mathbf{H}) \geq \sqrt{2}(\mu - 2\delta). \quad (\text{B.92})$$

□

The following lemma states an important property of $\widehat{\mathbf{H}}$ the optimal solution of problem (2.5).

Lemma B.5. *If $\max_i \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta$ and $\widehat{\mathbf{H}}$ is the optimal solution of problem (2.5), then we have*

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.93})$$

Proof. First note that since $\delta \geq \max_i \|\mathbf{Z}_{i,\cdot}\|_2$, we have

$$\max_{i \leq n} \mathcal{D}(\mathbf{X}_{i,\cdot}, \text{conv}(\mathbf{H}_0))^{1/2} \leq \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta. \quad (\text{B.94})$$

Hence, \mathbf{H}_0 is a feasible solution for the problem (2.5). Therefore, we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}). \quad (\text{B.95})$$

Letting $\tilde{\alpha}_i = \arg \min_{\alpha \in \Delta^n} \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}^\top \alpha\|_2$, we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) = \sum_{i=1}^r \min_{\alpha_i \in \Delta^r} \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \alpha_i - \mathbf{Z}^\top \alpha_i\|_2^2 \quad (\text{B.96})$$

$$= \sum_{i=1}^r \min_{\alpha_i \in \Delta^r} \left(\|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \alpha_i\|_2^2 - 2 \langle \mathbf{Z}^\top \alpha_i, \widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \alpha_i \rangle + \|\mathbf{Z}^\top \alpha_i\|_2^2 \right) \quad (\text{B.97})$$

$$= \sum_{i=1}^r \left(\|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\alpha}_i\|_2^2 - 2 \langle \mathbf{Z}^\top \tilde{\alpha}_i, \widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\alpha}_i \rangle + \|\mathbf{Z}^\top \tilde{\alpha}_i\|_2^2 \right). \quad (\text{B.98})$$

Using the fact that (by triangle inequality) $\|\mathbf{Z}^\top \tilde{\alpha}_i\|_2 \leq \delta$, we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq \sum_{i=1}^r \left(\|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\alpha}_i\|_2^2 - 2\delta \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\alpha}_i\|_2 \right) \quad (\text{B.99})$$

$$\geq U^2 - 2\delta\sqrt{r}U \quad (\text{B.100})$$

where $U^2 = \sum_{i=1}^r \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\alpha}_i\|_2^2$. Note that $\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq 0$ and for $U \geq 2\delta\sqrt{r}$, the function $U^2 - 2\delta\sqrt{r}U$ is increasing. Hence, since

$$U \geq \left(\sum_{i=1}^r \min_{\alpha_i} \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \alpha_i\|_2^2 \right)^{1/2} = \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2}, \quad (\text{B.101})$$

we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq (U^2 - 2\delta\sqrt{r}U) \mathbb{1}_{U \geq 2\delta\sqrt{r}} \geq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0) - 2\delta\sqrt{r} \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2}. \quad (\text{B.102})$$

Therefore,

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X})^{1/2} \geq \left(\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0) - 2\delta\sqrt{r} \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \right)_+^{1/2} \geq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} - 2\delta\sqrt{r}. \quad (\text{B.103})$$

In addition,

$$\mathcal{D}(\mathbf{H}_0, \mathbf{X}) = \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^n} \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i - \mathbf{Z}^\top \boldsymbol{\alpha}_i\|_2^2 \quad (\text{B.104})$$

$$\leq \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^n} \left\{ \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2 + \|\mathbf{Z}^\top \boldsymbol{\alpha}_i\|_2 \right\}^2 \quad (\text{B.105})$$

$$\leq \sum_{i=1}^r \left\{ \min_{\boldsymbol{\alpha}_i \in \Delta^n} \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2 + \max_{\boldsymbol{\alpha}_i \in \Delta^n} \|\mathbf{Z}^\top \boldsymbol{\alpha}_i\|_2 \right\}^2 \quad (\text{B.106})$$

$$\leq \left\{ \left(\sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^n} \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2^2 \right)^{1/2} + \delta\sqrt{r} \right\}^2 \quad (\text{B.107})$$

$$\leq \left(\mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \delta\sqrt{r} \right)^2. \quad (\text{B.108})$$

Hence,

$$\mathcal{D}(\mathbf{H}_0, \mathbf{X})^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \delta\sqrt{r}. \quad (\text{B.109})$$

Combining equations (B.103), (B.109), and (B.95), we get

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.110})$$

This completes the proof of lemma. \square

Lemma B.6. *Let \mathbf{X}_0 be such that the uniqueness assumption holds with parameter $\alpha > 0$, and $\text{conv}(\mathbf{X}_0)$ has internal radius at least $\mu > 0$. In particular, we have $\mathbf{z}_0 + \mathbf{U}\mathbf{B}_{r-1}(\mu) \subseteq \text{conv}(\mathbf{X}_0)$ for $\mathbf{z}_0 \in \mathbb{R}^d$, and an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{d \times (r-1)}$. Finally assume $\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta$. Then for $\widehat{\mathbf{H}}$ the optimal solution of problem (2.5), we have*

$$\alpha(\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}) \leq 2(1 + 2\alpha) \left[r^{3/2} \delta \kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{\delta\sqrt{r}}{\mu} \sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top) \right] + 3\delta\sqrt{r} \quad (\text{B.111})$$

where $\mathbf{P}_0 : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the orthogonal projector onto $\text{aff}(\mathbf{H}_0)$ (in particular, \mathbf{P}_0 is an affine map).

Proof. Let $\widetilde{\mathbf{H}}$ be such that $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\widetilde{\mathbf{H}})$. The uniqueness assumption implies

$$\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} \geq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \alpha(\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2}). \quad (\text{B.112})$$

Note that Lemma B.5 implies

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.113})$$

Therefore,

$$\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} \geq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} - 3\delta\sqrt{r} + \alpha(\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2}). \quad (\text{B.114})$$

Hence,

$$\alpha(\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2}) \leq \mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.115})$$

In addition, for a convex set S , by triangle inequality we have

$$\left[\sum_{i=1}^n \mathcal{D}(\widehat{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} - \left[\sum_{i=1}^n \mathcal{D}(\widetilde{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} \leq \left[\sum_{i=1}^n \left(\mathcal{D}(\widehat{\mathbf{H}}_{i,\cdot}, S)^{1/2} - \mathcal{D}(\widetilde{\mathbf{H}}_{i,\cdot}, S)^{1/2} \right)^2 \right]^{1/2} \quad (\text{B.116})$$

Therefore, using

$$|\mathcal{D}(\widetilde{\mathbf{H}}_{i,\cdot}, S)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}_{i,\cdot}, S)^{1/2}| \leq \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \quad (\text{B.117})$$

we have

$$\left[\sum_{i=1}^n \mathcal{D}(\widehat{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} - \left[\sum_{i=1}^n \mathcal{D}(\widetilde{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} \leq \left[\sum_{i=1}^n \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2^2 \right]^{1/2} = \|\widehat{\mathbf{H}} - \widetilde{\mathbf{H}}\|_F. \quad (\text{B.118})$$

Hence,

$$|\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2}| \leq \|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F \quad (\text{B.119})$$

and

$$|\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2}| \leq \|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F. \quad (\text{B.120})$$

In addition, similarly to the proof of Lemma B.5, we can write

$$\mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}}) = \sum_{i=1}^r \min_{\alpha_i \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widetilde{\mathbf{H}}^\top \alpha_i\|_2^2 \quad (\text{B.121})$$

$$= \sum_{i=1}^r \min_{\alpha_i \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \alpha_i - (\widehat{\mathbf{H}} - \widetilde{\mathbf{H}})^\top \alpha_i\|_2^2 \quad (\text{B.122})$$

$$\leq \sum_{i=1}^r \min_{\alpha_i \in \Delta^r} \left\{ \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \alpha_i\|_2 + \|(\widehat{\mathbf{H}} - \widetilde{\mathbf{H}})^\top \alpha_i\|_2 \right\}^2 \quad (\text{B.123})$$

$$\leq \sum_{i=1}^r \left\{ \min_{\alpha \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \alpha\|_2 + \max_{\alpha \in \Delta^r} \|(\widehat{\mathbf{H}} - \widetilde{\mathbf{H}})^\top \alpha\|_2 \right\}^2 \quad (\text{B.124})$$

$$\leq \left\{ \left(\sum_{i=1}^r \min_{\alpha \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \alpha_i\|_2^2 \right)^{1/2} + \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \right\}^2 \quad (\text{B.125})$$

$$\leq \left(\mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} + \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \right)^2. \quad (\text{B.126})$$

Thus,

$$|\mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2} - \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}| \leq \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \quad (\text{B.127})$$

Therefore, combining (B.115), (B.119), (B.120), (B.127), we get

$$\alpha \left(\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \right) \leq (1 + \alpha) \|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F + \alpha \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 + 3\delta \sqrt{r}. \quad (\text{B.128})$$

Now, we would like to bound the terms $\|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F, \max_{i \in [r]} \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2$. Note that using the fact that $\widehat{\mathbf{H}}$ is feasible for Problem (2.5), we have

$$\mathcal{D}(\mathbf{X}_{i,\cdot}, \widehat{\mathbf{H}}) \leq \delta^2. \quad (\text{B.129})$$

Thus,

$$\mathcal{D}((\mathbf{X}_0)_{i,\cdot}, \widehat{\mathbf{H}})^{1/2} \leq \|\mathbf{X}_{i,\cdot} - (\mathbf{X}_0)_{i,\cdot}\|_2 + \mathcal{D}(\mathbf{X}_{i,\cdot}, \widehat{\mathbf{H}})^{1/2} \leq 2\delta. \quad (\text{B.130})$$

In addition, we know that $(\mathbf{X}_0)_{i,\cdot} \in \text{aff}(\mathbf{H}_0)$, where $\text{aff}(\mathbf{H}_0)$ is a $r-1$ dimensional affine subspace. Therefore, $\text{conv}(\mathbf{X}_0) \subseteq \text{aff}(\mathbf{H}_0)$ and, by convexity of $\text{B}_d(2\delta, \widehat{\mathbf{H}})$, we get

$$\text{conv}(\mathbf{X}_0) \subseteq \text{B}_d(2\delta, \widehat{\mathbf{H}}) \cap \text{aff}(\mathbf{H}_0). \quad (\text{B.131})$$

First consider the case in which $\widehat{\mathbf{H}}_{i,\cdot} \in \text{aff}(\mathbf{H}_0)$, for all $i \in \{1, 2, \dots, r\}$. By a perturbation argument, we can assume that the rows of $\widehat{\mathbf{H}}$ are linearly independent, and hence $\text{aff}(\widehat{\mathbf{H}}) = \text{aff}(\mathbf{H}_0)$. Consider $\widetilde{\mathbf{Q}} \in \mathbb{R}^{r \times d}$ defined by

$$\widetilde{Q}_{ii} = 1 + \xi, \quad \text{if } i = j \in \{1, 2, \dots, r\}, \quad (\text{B.132})$$

$$\widetilde{Q}_{ij} = -\frac{\xi}{r-1}, \quad \text{if } i \neq j \in \{1, 2, \dots, r\}, \quad (\text{B.133})$$

$$\widetilde{Q}_{ij} = 0, \quad \text{if } j \in \{r+1, r+2, \dots, d\} \quad (\text{B.134})$$

where $\xi = 2r\delta_0$. Note that for every $\mathbf{y} \in \text{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d})$, we have $\mathcal{D}(\mathbf{y}, \mathbf{E}_{r,d})^{1/2} \leq 2\delta_0$. In addition, since $\mathbf{y} \in \text{aff}(\mathbf{E}_{r,d})$, $\langle \mathbf{y}, \mathbf{1} \rangle = 1$. Hence, for $\mathbf{y} \in \text{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d})$, we can write

$$\mathbf{y} = \boldsymbol{\pi} + \mathbf{x} \quad (\text{B.135})$$

where $\boldsymbol{\pi} \in \text{conv}(\mathbf{E}_{r,d})$, $\mathbf{x} \in \mathbb{R}^d$, $\langle \mathbf{1}, \mathbf{x} \rangle = 0$, $\|\mathbf{x}\|_2 \leq 2\delta_0$. It is easy to check that for this \mathbf{y} we have

$$\mathbf{y} = \sum_{i=1}^r \beta_i \widetilde{\mathbf{Q}}_{i,\cdot} \quad (\text{B.136})$$

where $\boldsymbol{\beta} \in \mathbb{R}^r$ is such that for $i = 1, 2, \dots, r$,

$$\beta_i = \frac{r-1}{r-1+\xi r} (\pi_i + x_i) + \frac{\xi}{r-1+\xi r}. \quad (\text{B.137})$$

Further, note that since $\boldsymbol{\pi} \in \text{conv}(\mathbf{E}_{r,d})$, $\pi_i \geq 0$ and $x_i \geq -\|\mathbf{x}\|_2 \geq -2\delta_0$, we have $\pi_i + x_i \geq -2\delta_0$. Hence, for $i \in \{1, 2, \dots, r\}$,

$$\beta_i \geq \frac{-2\delta_0(r-1) + \xi}{r-1+\xi r} = \frac{2\delta_0}{r-1+\xi r} \geq 0. \quad (\text{B.138})$$

In addition,

$$\sum_{i=1}^r \beta_i = \frac{r\xi}{r-1+\xi r} + \frac{r-1}{r-1+\xi r} \left(\sum_{i=1}^r (\pi_i + x_i) \right) = 1. \quad (\text{B.139})$$

Therefore, every $\mathbf{y} \in \mathbb{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d})$ can be written as a convex combination of the rows of $\tilde{\mathbf{Q}}$. Hence,

$$\mathbb{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d}) \subseteq \text{conv}(\tilde{\mathbf{Q}}). \quad (\text{B.140})$$

Let $\widehat{\mathbf{H}} = \mathbf{E}_{r,d}\mathbf{M}$, $\mathbf{M} \in \mathbb{R}^{d \times d}$. Since $\text{aff}(\widehat{\mathbf{H}}) = \text{aff}(\mathbf{H}_0)$, by taking $\widetilde{\mathbf{H}} = \tilde{\mathbf{Q}}\mathbf{M}$, we have

$$\text{conv}(\widetilde{\mathbf{H}}) \supseteq \left[\bigcup_{\mathbf{x} \in \text{conv}(\mathbf{E}_{r,d})} \mathbf{M}^\top \mathbb{B}_d(2\delta_0; \mathbf{x}) \right] \cap \text{aff}(\widetilde{\mathbf{H}}) \quad (\text{B.141})$$

$$\supseteq \left[\bigcup_{\mathbf{x} \in \text{conv}(\widehat{\mathbf{H}})} \mathbb{B}_d(2\delta_0 \sigma_{\min}(\mathbf{M}); \mathbf{x}) \right] \cap \text{aff}(\widetilde{\mathbf{H}}) \quad (\text{B.142})$$

$$\supseteq \mathbb{B}_d(2\delta; \widehat{\mathbf{H}}) \cap \text{aff}(\mathbf{H}_0), \quad (\text{B.143})$$

provided that $\delta_0 = \delta / \sigma_{\min}(\mathbf{M}) = \delta / \sigma_{\min}(\widehat{\mathbf{H}})$. Hence, using (B.131) for this δ_0 , $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\widetilde{\mathbf{H}})$. Note that for $\tilde{\mathbf{Q}}$, we have $\|\tilde{\mathbf{Q}}_{i,\cdot} - \mathbf{e}_i\|_2 \leq 2r\delta_0$. Thus,

$$\|\tilde{\mathbf{Q}} - \mathbf{E}_{r,d}\|_F \leq 2r^{3/2}\delta_0. \quad (\text{B.144})$$

Therefore, there exists $\widetilde{\mathbf{H}} \in \mathbb{R}^{r \times d}$ such that $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\widetilde{\mathbf{H}})$ and

$$\begin{aligned} \|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F &= \|(\tilde{\mathbf{Q}} - \mathbf{E}_{r,d})\mathbf{M}\|_F \leq 2r^{3/2}\delta_0\sigma_{\max}(\mathbf{M}) = 2r^{3/2}\delta_0\sigma_{\max}(\widehat{\mathbf{H}}) = 2r^{3/2}\delta\kappa(\widehat{\mathbf{H}}), \\ \max_{i \in [r]} \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 &= \max_{i \in [r]} \|(\tilde{\mathbf{Q}}_{i,\cdot} - \mathbf{e}_i)\mathbf{M}\|_2 \leq 2r\delta_0\sigma_{\max}(\mathbf{M}) = 2r\delta_0\sigma_{\max}(\widehat{\mathbf{H}}) = 2r\delta\kappa(\widehat{\mathbf{H}}). \end{aligned}$$

Now consider the general case in which $\text{aff}(\widehat{\mathbf{H}}) \neq \text{aff}(\mathbf{H}_0)$. Let $\mathbf{H}' \in \mathbb{R}^{r \times d}$ be such that $\mathbf{H}'_{i,\cdot}$ is the projection of $\widehat{\mathbf{H}}_{i,\cdot}$ onto $\text{aff}(\mathbf{H}_0)$. Assuming that the rows of \mathbf{H}' are linearly independent, $\text{aff}(\mathbf{H}') = \text{aff}(\mathbf{H}_0)$. Note that since $\text{conv}(\mathbf{X}_0) \in \text{aff}(\mathbf{H}_0)$, for every point $\mathbf{x} \in \text{conv}(\mathbf{X}_0)$, $\mathcal{D}(\mathbf{x}, \mathbf{H}')^{1/2} \leq \mathcal{D}(\mathbf{x}, \widehat{\mathbf{H}})^{1/2} \leq 2\delta$. Thus,

$$(\mathbf{X}_0)_{i,\cdot} \in \mathbb{B}_d(2\delta, \mathbf{H}') \cap \text{aff}(\mathbf{H}'). \quad (\text{B.145})$$

Therefore, using the above argument for the case where $\text{aff}(\widehat{\mathbf{H}}) = \text{aff}(\mathbf{H}_0)$, we can find $\widetilde{\mathbf{H}}$ such that $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\widetilde{\mathbf{H}})$ and

$$\|\widetilde{\mathbf{H}} - \mathbf{H}'\|_F \leq 2r^{3/2}\delta\kappa(\mathbf{H}'), \quad (\text{B.146})$$

$$\max_{i \in [r]} \|\widetilde{\mathbf{H}}_{i,\cdot} - \mathbf{H}'_{i,\cdot}\|_2 \leq 2r\delta\kappa(\mathbf{H}'). \quad (\text{B.147})$$

Hence, for every $i = 1, 2, \dots, r$,

$$\begin{aligned} \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 &\leq \|\widetilde{\mathbf{H}}_{i,\cdot} - \mathbf{H}'_{i,\cdot}\|_2 + \|\mathbf{H}'_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \\ &\leq 2r\delta\kappa(\mathbf{H}') + \|\mathbf{P}_0(\widehat{\mathbf{H}}_{i,\cdot}) - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \end{aligned} \quad (\text{B.148})$$

where \mathbf{P}_0 orthogonal projection onto $\text{aff}(\mathbf{H}_0)$. We next use the assumption on the internal radius of $\text{conv}(\mathbf{X}_0)$ to upper bound the term $\|\mathbf{P}_0(\widehat{\mathbf{H}}_{i,\cdot}) - \widehat{\mathbf{H}}_{i,\cdot}\|_2$. Note that since $\text{conv}(\mathbf{X}_0) \subseteq \mathbb{B}_d(2\delta, \widehat{\mathbf{H}})$, letting $\bar{\mathbf{H}} = \widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top$, for some orthogonal matrix $\mathbf{U} \in \mathbb{R}^{d \times r'}$, $r' = r - 1$, we have

$$\max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a}\|_2^2 = \max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} - (\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)^\top \mathbf{a}\|_2^2 \quad (\text{B.149})$$

$$\leq \max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} + \mathbf{z}_0 - \widehat{\mathbf{H}}^\top \mathbf{a}\|_2^2 \leq 4\delta^2. \quad (\text{B.150})$$

Now, using Cauchy-Schwarz inequality we can write

$$\max_{\|z\|_2 \leq \mu} \min_{\|a\|_2 \leq 1} \|Uz - \bar{H}^T a\|_2^2 \leq \max_{\|z\|_2 \leq \mu} \min_{\langle a, \mathbf{1} \rangle = 1, a \geq 0} \|Uz - \bar{H}^T a\|_2^2 \leq 4\delta^2. \quad (\text{B.151})$$

Note that,

$$\min_{\|a\|_2 \leq 1} \|Uz - \bar{H}^T a\|_2^2 = \max_{\rho \geq 0} \min_a \left\{ \|z\|_2^2 - 2 \langle z, U^T \bar{H}^T a \rangle + \langle a, (\bar{H}\bar{H}^T + \rho \mathbf{I}) a \rangle - \rho \right\} \quad (\text{B.152})$$

$$= \max_{\rho \geq 0} \left\{ \|z\|_2^2 - \langle \bar{H}Uz, (\bar{H}\bar{H}^T + \rho \mathbf{I})^{-1} \bar{H}Uz \rangle - \rho \right\} \quad (\text{B.153})$$

Hence, using (B.151)

$$\mu^2 \max_{\rho \geq 0} \left\{ \lambda_{\max}(\mathbf{I} - U^T \bar{H}^T (\bar{H}\bar{H}^T + \rho \mathbf{I})^{-1} \bar{H}U) - \rho \leq 4\delta^2 \right\}. \quad (\text{B.154})$$

In particular, for $\rho = 0$ we get

$$\mu^2 \lambda_{\max}(\mathbf{I} - U^T \bar{H}^T (\bar{H}\bar{H}^T)^{-1} \bar{H}U) \leq 4\delta^2. \quad (\text{B.155})$$

Taking $\bar{H} = \tilde{U}\Sigma\tilde{V}^T$, the singular value decomposition of \bar{H} , we have $\sigma_{\max}(\bar{H}) = \sigma_{\max}(\widehat{H} - \mathbf{1}z_0^T) = \max_i \Sigma_{ii}$. Letting $U^T \tilde{V} = Q$, we get

$$\max_{\rho \geq 0} \lambda_{\max}(\mathbf{I} - QQ^T) \leq \frac{4\delta^2}{\mu^2}. \quad (\text{B.156})$$

Letting $q = \sigma_{\min}(Q)$, this results in

$$1 - q^2 \leq \frac{4\delta^2}{\mu^2}. \quad (\text{B.157})$$

In addition, note that, by the internal radius assumption, for any $z \in \mathbb{R}^r$, $z_0 + Uz \in \text{aff}(\mathbf{H}_0)$. Further, since $z_0 \in \text{aff}(\mathbf{H}_0)$,

$$\max_{i \in [r]} \|P_0(\widehat{H}_{i,\cdot}) - \widehat{H}_{i,\cdot}\|_2 = \max_{i \in [r]} \|P_U(\bar{H}_{i,\cdot}) - \bar{H}_{i,\cdot}\|_2 \quad (\text{B.158})$$

$$\leq \max_{\|a\|_2 \leq 1} \|P_U(\bar{H}^T a) - \bar{H}^T a\|_2 \quad (\text{B.159})$$

$$\leq \max_{\|a\|_2 \leq 1} \|P_U(\bar{H}^T a) - \bar{H}^T a\|_2 \quad (\text{B.160})$$

$$\leq \max_{\|a\|_2 \leq 1} \min_z \|Uz - \bar{H}^T a\|_2^2 \quad (\text{B.161})$$

where P_U is the projector onto the column space of U . Note that,

$$\max_{\|a\|_2 \leq 1} \min_z \|Uz - \bar{H}^T a\|_2^2 = \max_{\|a\|_2 \leq 1} \left\{ - \langle a, \bar{H}U U^T \bar{H}^T a \rangle + \langle a, \bar{H}\bar{H}^T a \rangle \right\} \quad (\text{B.162})$$

$$= \lambda_{\max}(\bar{H}\bar{H}^T - \bar{H}U U^T \bar{H}^T) \quad (\text{B.163})$$

$$= \lambda_{\max}(\Sigma(\mathbf{I} - Q^T Q)\Sigma) \quad (\text{B.164})$$

$$\leq \sigma_{\max}(\bar{H})^2 \lambda_{\max}(\mathbf{I} - Q^T Q) \quad (\text{B.165})$$

$$\leq \sigma_{\max}(\bar{H})^2 (1 - q^2) \leq \frac{4\sigma_{\max}(\bar{H})^2 \delta^2}{\mu^2} \quad (\text{B.166})$$

where the last inequality follows from (B.157). This results in

$$\max_{i \in [r]} \|\mathbf{P}_0(\widehat{\mathbf{H}}_{i,\cdot}) - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \leq \frac{2\sigma_{\max}(\bar{\mathbf{H}})\delta}{\mu} = \frac{2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta}{\mu} \quad (\text{B.167})$$

Therefore, $\|\mathbf{P}_0(\widehat{\mathbf{H}}) - \widehat{\mathbf{H}}\|_F \leq 2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta\sqrt{r}/\mu$. Hence, using (B.148) we get

$$\max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \leq 2r\delta\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta}{\mu}, \quad (\text{B.168})$$

$$\|\widehat{\mathbf{H}} - \widetilde{\mathbf{H}}\|_F \leq 2r^{3/2}\delta\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta\sqrt{r}}{\mu}. \quad (\text{B.169})$$

Replacing this in (B.128) completes the proof. \square

B.3.2 Proof of Theorem 2

For simplicity, let $\mathcal{D} = \alpha(\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2})$. First note that under the assumption of Theorem 2 we have

$$\mathbf{z}_0 + \mathbf{UB}_{r'}(\mu) \subseteq \text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\mathbf{H}_0). \quad (\text{B.170})$$

Therefore, using Lemma B.4 with $\mathbf{H} = \mathbf{H}_0$ and $\delta = 0$, we have

$$\mu\sqrt{2} \leq \sigma_{\min}(\mathbf{H}_0) \leq \sigma_{\max}(\mathbf{H}_0). \quad (\text{B.171})$$

In addition, since $\mathbf{z}_0 \in \text{conv}(\mathbf{H}_0)$ we have $\mathbf{z}_0 = \mathbf{H}_0^\top \boldsymbol{\alpha}_0$ for some $\boldsymbol{\alpha}_0 \in \Delta^r$. Therefore,

$$\|\mathbf{z}_0\|_2 \leq \sigma_{\max}(\mathbf{H}_0)\|\boldsymbol{\alpha}_0\|_2 \leq \sigma_{\max}(\mathbf{H}_0). \quad (\text{B.172})$$

Note that

$$\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top) \leq \sigma_{\max}(\widehat{\mathbf{H}}) + \sigma_{\max}(\mathbf{1}\mathbf{z}_0^\top) = \sigma_{\max}(\widehat{\mathbf{H}}) + \sqrt{r}\|\mathbf{z}_0\|_2. \quad (\text{B.173})$$

Therefore, using Lemma B.6 we have

$$\mathcal{D} \leq 2(1 + 2\alpha) \left(r^{3/2}\delta\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{\sigma_{\max}(\widehat{\mathbf{H}})\delta r^{1/2}}{\mu} + \frac{r\delta\|\mathbf{z}_0\|_2}{\mu} \right) + 3\delta r^{1/2}. \quad (\text{B.174})$$

In addition, Lemma B.2 implies that

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{1}{\alpha} \max \left\{ (1 + \sqrt{2})\sqrt{r}, \sqrt{2}\kappa(\mathbf{H}_0) \right\} \mathcal{D}. \quad (\text{B.175})$$

Further, let \mathbf{P}_0 denote the orthogonal projector on $\text{aff}(\mathbf{H}_0)$. Hence, \mathbf{P}_0 is a non-expansive mapping: for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $D(\mathbf{P}_0(\mathbf{x}), \mathbf{P}_0(\mathbf{y})) \leq D(\mathbf{x}, \mathbf{y})$. Therefore, since $\text{conv}(\mathbf{H}_0) \subset \text{aff}(\mathbf{H}_0)$, for any $\mathbf{h} \in \mathbb{R}^d$

$$\mathcal{D}(\mathbf{P}_0(\mathbf{h}), \mathbf{H}_0) \leq D(\mathbf{P}_0(\mathbf{h}), \mathbf{P}_0(\boldsymbol{\Pi}_{\text{conv}(\mathbf{H}_0)}(\mathbf{h}))) \leq D(\mathbf{h}, \boldsymbol{\Pi}_{\text{conv}(\mathbf{H}_0)}(\mathbf{h})) = \mathcal{D}(\mathbf{h}, \mathbf{H}_0). \quad (\text{B.176})$$

Therefore,

$$\mathcal{D}(\mathbf{P}_0(\widehat{\mathbf{H}}), \mathbf{H}_0) \leq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0). \quad (\text{B.177})$$

First consider the case in which

$$\delta \leq \frac{\alpha\mu}{30r^{3/2}}. \quad (\text{B.178})$$

Note that in this case $\delta \leq \mu/2$. Hence, using Lemma B.3 to upper bound $\sigma_{\max}(\widehat{\mathbf{H}})$, $\sigma_{\max}(\mathbf{P}_0(\widehat{\mathbf{H}}))$ and Lemma B.4 to lower bound $\sigma_{\min}(\mathbf{P}_0(\widehat{\mathbf{H}}))$, by (B.177), we get

$$\sigma_{\max}(\widehat{\mathbf{H}}) \leq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + r^{1/2}\sigma_{\max}(\mathbf{H}_0) \leq \frac{\mathcal{D}}{\alpha} + r^{1/2}\sigma_{\max}(\mathbf{H}_0), \quad (\text{B.179})$$

$$\begin{aligned} \kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) &= \frac{\sigma_{\max}(\mathbf{P}_0(\widehat{\mathbf{H}}))}{\sigma_{\min}(\mathbf{P}_0(\widehat{\mathbf{H}}))} \leq \frac{\mathcal{D}(\mathbf{P}_0(\widehat{\mathbf{H}}), \mathbf{H}_0)^{1/2} + r^{1/2}\sigma_{\max}(\mathbf{H}_0)}{\sqrt{2}(\mu - 2\delta)} \\ &\leq \frac{\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + r^{1/2}\sigma_{\max}(\mathbf{H}_0)}{\sqrt{2}(\mu - 2\delta)} \leq \frac{\mathcal{D}}{\alpha(\mu - 2\delta)\sqrt{2}} + \frac{r^{1/2}\sigma_{\max}(\mathbf{H}_0)}{(\mu - 2\delta)\sqrt{2}}. \end{aligned} \quad (\text{B.180})$$

Replacing these in (B.174) we have

$$\mathcal{D} \leq 2(1 + 2\alpha) \left[\frac{r^{3/2}\mathcal{D}\delta}{\alpha(\mu - 2\delta)\sqrt{2}} + \frac{r^2\sigma_{\max}(\mathbf{H}_0)\delta}{(\mu - 2\delta)\sqrt{2}} + \frac{\mathcal{D}r^{1/2}\delta}{\alpha\mu} + \frac{r\sigma_{\max}(\mathbf{H}_0)\delta}{\mu} + \frac{r\|\mathbf{z}_0\|_2\delta}{\mu} \right] + 3\delta\sqrt{r}. \quad (\text{B.181})$$

Therefore,

$$\begin{aligned} &\mathcal{D} \left[1 - \frac{\sqrt{2}(1 + 2\alpha)r^{3/2}\delta}{\alpha(\mu - 2\delta)} - \frac{2(1 + 2\alpha)r^{1/2}\delta}{\alpha\mu} \right] \\ &\leq 2(1 + 2\alpha) \left[\frac{r^2\sigma_{\max}(\mathbf{H}_0)\delta}{(\mu - 2\delta)\sqrt{2}} + \frac{r\sigma_{\max}(\mathbf{H}_0)\delta}{\mu} + \frac{r\|\mathbf{z}_0\|_2\delta}{\mu} \right] + 3\delta\sqrt{r} \end{aligned} \quad (\text{B.182})$$

Notice that condition (B.178) implies that $\mu - 2\delta \geq \mu/2$ and

$$\frac{\sqrt{2}(1 + 2\alpha)r^{3/2}\delta}{\alpha(\mu - 2\delta)} + \frac{2(1 + 2\alpha)r^{1/2}\delta}{\alpha\mu} \leq \frac{1}{2}. \quad (\text{B.183})$$

Using the previous two equations, under condition (B.178) we have

$$\begin{aligned} \mathcal{D} &\leq \frac{4(1 + 2\alpha)r\delta}{\mu} \left[\frac{5r\sigma_{\max}(\mathbf{H}_0)}{2} + \|\mathbf{z}_0\|_2 \right] + 3\delta\sqrt{r} \\ &\leq \frac{4(1 + 2\alpha)r^2}{\mu} \left[\frac{5\sigma_{\max}(\mathbf{H}_0)}{2} + \frac{\|\mathbf{z}_0\|_2}{r} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}} \right] \delta. \end{aligned} \quad (\text{B.184})$$

Combining this with (B.175), and using the fact that $1 + 2\alpha \leq 3$, we have under condition (B.178)

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{12r^2}{\mu\alpha} \left(\frac{5\sigma_{\max}(\mathbf{H}_0)}{2} + \frac{\|\mathbf{z}_0\|_2}{r} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}} \right) \max \left\{ (1 + \sqrt{2})\sqrt{r}, \sqrt{2}\kappa(\mathbf{H}_0) \right\} \delta \quad (\text{B.185})$$

$$\leq \frac{29\sigma_{\max}(\mathbf{H}_0)r^{5/2}}{\alpha\mu} \max \left\{ 1, \frac{\kappa(\mathbf{H}_0)}{\sqrt{r}} \right\} \left(\frac{5}{2} + \frac{\|\mathbf{z}_0\|_2}{r\sigma_{\max}(\mathbf{H}_0)} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}\sigma_{\max}(\mathbf{H}_0)} \right) \delta. \quad (\text{B.186})$$

Note that using (B.171), (B.172) and since $\alpha \geq 0$

$$\frac{\|\mathbf{z}_0\|_2}{r\sigma_{\max}(\mathbf{H}_0)} \leq 1, \quad \frac{3\mu}{4(1+2\alpha)r^{3/2}\sigma_{\max}(\mathbf{H}_0)} \leq \frac{3}{4\sqrt{2}}. \quad (\text{B.187})$$

Therefore,

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{120\sigma_{\max}(\mathbf{H}_0)r^{5/2}}{\alpha\mu} \max\left\{1, \frac{\kappa(\mathbf{H}_0)}{\sqrt{r}}\right\} \delta. \quad (\text{B.188})$$

Thus,

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_*^2 r^5}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (\text{B.189})$$

where C_* is defined in Theorem 2.

Next, consider the case in which

$$\delta = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha\mu}{330\kappa(\mathbf{H}_0)r^{5/2}}, \quad (\text{B.190})$$

Note that using (B.171),(B.172) and since $1 + 2\alpha \leq 3$, this condition on δ implies that

$$\delta \leq \frac{\alpha\mu\sigma_{\min}(\mathbf{H}_0)}{12r(1+2\alpha)(5r^{3/2}\sigma_{\max}(\mathbf{H}_0) + 2\|\mathbf{z}_0\|_2r^{1/2} + 3\mu)}. \quad (\text{B.191})$$

In particular, condition (B.178) holds. Hence, using equation (B.184) we get

$$\mathcal{D} \leq \frac{4(1+2\alpha)r^2}{\mu} \left[\frac{5\sigma_{\max}(\mathbf{H}_0)}{2} + \frac{\|\mathbf{z}_0\|_2}{r} + \frac{3\mu}{4(1+2\alpha)r^{3/2}} \right] \delta \leq \frac{\alpha\sigma_{\min}(\mathbf{H}_0)}{6\sqrt{r}}. \quad (\text{B.192})$$

Further, note that since \mathbf{P}_0 is a projection onto an affine subspace, for $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{P}_0(\mathbf{x}) = \widetilde{\mathbf{P}}_0\mathbf{x} + \mathbf{x}_0$ for some $\widetilde{\mathbf{P}}_0 \in \mathbb{R}^{d \times d}$, $\mathbf{x}_0 \in \mathbb{R}^d$. Hence, for any $\boldsymbol{\pi} \in \Delta^r$, $\mathbf{h} = \widehat{\mathbf{H}}^\top \boldsymbol{\pi} \in \text{conv}(\widehat{\mathbf{H}})$, we have

$$\mathbf{P}_0(\mathbf{h}) = \widetilde{\mathbf{P}}_0\mathbf{h} + \mathbf{x}_0 = \widetilde{\mathbf{P}}_0\widehat{\mathbf{H}}^\top \boldsymbol{\pi} + \mathbf{x}_0 = \sum_{i=1}^r \pi_i \left(\widetilde{\mathbf{P}}_0\widehat{\mathbf{H}}^\top \mathbf{e}_i + \mathbf{x}_0 \right) = \sum_{i=1}^r \pi_i \mathbf{P}_0(\widehat{\mathbf{h}}_i) \in \text{conv}(\mathbf{P}_0(\widehat{\mathbf{H}})) \quad (\text{B.193})$$

where \mathbf{e}_i is the i 'th standard unit vector. Hence,

$$\mathbf{P}_0(\text{conv}(\widehat{\mathbf{H}})) \subseteq \text{conv}(\mathbf{P}_0(\widehat{\mathbf{H}})). \quad (\text{B.194})$$

Thus, for $\mathbf{h}_0 \in \mathbb{R}^d$ an arbitrary row of \mathbf{H}_0 , we have

$$\mathcal{D}(\mathbf{h}_0, \mathbf{P}_0(\widehat{\mathbf{H}})) = D(\mathbf{h}_0, \text{conv}(\mathbf{P}_0(\widehat{\mathbf{H}}))) \leq D(\mathbf{h}_0, \mathbf{P}_0(\text{conv}(\widehat{\mathbf{H}}))) \leq D(\mathbf{h}_0, \mathbf{P}_0(\boldsymbol{\Pi}_{\text{conv}(\widehat{\mathbf{H}})}(\mathbf{h}_0))). \quad (\text{B.195})$$

In addition, using non-expansivity of \mathbf{P}_0 , we have

$$D(\mathbf{h}_0, \mathbf{P}_0(\boldsymbol{\Pi}_{\text{conv}(\widehat{\mathbf{H}})}(\mathbf{h}_0))) \leq D(\mathbf{h}_0, \boldsymbol{\Pi}_{\text{conv}(\widehat{\mathbf{H}})}(\mathbf{h}_0)) = D(\mathbf{h}_0, \text{conv}(\widehat{\mathbf{H}})) = \mathcal{D}(\mathbf{h}_0, \widehat{\mathbf{H}}). \quad (\text{B.196})$$

This implies that

$$\mathcal{D}(\mathbf{H}_0, \mathbf{P}_0(\widehat{\mathbf{H}})) \leq \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}}). \quad (\text{B.197})$$

Therefore, using (B.177), (B.197) and (B.192) we get

$$\mathcal{D}(\mathbf{H}_0, \mathbf{P}_0(\widehat{\mathbf{H}}))^{1/2} + \mathcal{D}(\mathbf{P}_0(\widehat{\mathbf{H}}), \mathbf{H}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} + \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} \leq \frac{\mathcal{D}}{\alpha} \leq \frac{\sigma_{\min}(\mathbf{H}_0)}{6\sqrt{r}}. \quad (\text{B.198})$$

Hence, in this case Lemma B.3 implies that

$$\sigma_{\max}(\widehat{\mathbf{H}}) \leq 2\sigma_{\max}(\mathbf{H}_0), \quad (\text{B.199})$$

$$\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) \leq \frac{7\kappa(\mathbf{H}_0)}{2}. \quad (\text{B.200})$$

Replacing this in (B.174), we have

$$\mathcal{D} \leq (1 + 2\alpha)r^{1/2} \left(7r\delta\kappa(\mathbf{H}_0) + \frac{4\sigma_{\max}(\mathbf{H}_0)\delta + 2\sqrt{r}\|\mathbf{z}_0\|_2\delta}{\mu} \right) + 3\delta r^{1/2} \quad (\text{B.201})$$

$$\leq 3\delta\sqrt{r} \left(8r\kappa(\mathbf{H}_0) + \frac{4\sigma_{\max}(\mathbf{H}_0) + 2\sqrt{r}\|\mathbf{z}_0\|_2}{\mu} \right) \quad (\text{B.202})$$

Hence, using (B.175) under assumption (B.190), we have

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq 3\sqrt{r} \max \left\{ (1 + \sqrt{2})\sqrt{r}, \sqrt{2}\kappa(\mathbf{H}_0) \right\} \left(8r\kappa(\mathbf{H}_0) + \frac{4\sigma_{\max}(\mathbf{H}_0) + 2\sqrt{r}\|\mathbf{z}_0\|_2}{\mu} \right) \frac{\delta}{\alpha} \quad (\text{B.203})$$

$$\leq 120 \max \left\{ 1, \frac{\kappa(\mathbf{H}_0)}{\sqrt{r}} \right\} \max \left\{ r\kappa(\mathbf{H}_0), \frac{\sigma_{\max}(\mathbf{H}_0) + \sqrt{r}\|\mathbf{z}_0\|_2}{\mu} \right\} \frac{r\delta}{\alpha}. \quad (\text{B.204})$$

Hence, for C_*'' as defined in the statement of the theorem, we get

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{C_*'' r}{\alpha} \max_{i \leq n} \|\mathbf{Z}_{i \cdot}\|_2 \quad (\text{B.205})$$

This completes the proof.

C Proof of Proposition 4.2

The proof follows immediately from the following two propositions.

Proposition C.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$. Then the gradient of the function $\mathbf{u} \mapsto \mathcal{D}(\mathbf{u}, \mathbf{X})$ is given by*

$$\nabla_{\mathbf{u}} \mathcal{D}(\mathbf{u}, \mathbf{X}) = 2(\mathbf{u} - \mathbf{\Pi}_{\text{conv}(\mathbf{X})}(\mathbf{u})). \quad (\text{C.1})$$

Proof. Note that $\mathcal{D}(\mathbf{u}, \mathbf{X})$ is the solution of the following convex optimization problem.

$$\begin{aligned} & \text{minimize} && \|\mathbf{u} - \mathbf{y}\|_2^2, \\ & \text{subject to} && \mathbf{y} = \mathbf{X}^T \boldsymbol{\pi}, \\ & && \boldsymbol{\pi} \geq 0, \\ & && \langle \boldsymbol{\pi}, \mathbf{1} \rangle = 1. \end{aligned} \quad (\text{C.2})$$

The Lagrangian for this problem is

$$\mathcal{L}(\mathbf{y}, \boldsymbol{\pi}, \boldsymbol{\rho}, \tilde{\rho}, \boldsymbol{\lambda}) = \|\mathbf{u} - \mathbf{y}\|_2^2 + \langle \boldsymbol{\rho}, (\mathbf{y} - \mathbf{X}^\top \boldsymbol{\pi}) \rangle - \langle \boldsymbol{\lambda}, \boldsymbol{\pi} \rangle + \tilde{\rho}(1 - \langle \boldsymbol{\pi}, \mathbf{1} \rangle). \quad (\text{C.3})$$

The KKT condition implies that at the minimizer $(\mathbf{y}^*, \boldsymbol{\pi}^*, \boldsymbol{\rho}^*, \tilde{\rho}^*, \boldsymbol{\lambda}^*)$, we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = 0, \quad (\text{C.4})$$

and therefore

$$\boldsymbol{\rho}^* = 2(\mathbf{u} - \mathbf{y}^*) \quad (\text{C.5})$$

and the dual of the above optimization problem is

$$\begin{aligned} & \text{maximize} && -\frac{1}{4}\|\boldsymbol{\rho}\|_2^2 + \langle \boldsymbol{\rho}, \mathbf{u} \rangle + \tilde{\rho}, \\ & \text{subject to} && \boldsymbol{\lambda} \geq 0, \\ & && \mathbf{X}\boldsymbol{\rho} + \tilde{\rho}\mathbf{1} + \boldsymbol{\lambda} = 0. \end{aligned} \quad (\text{C.6})$$

Note that since (C.2) is strictly feasible, Slater condition holds and by strong duality the optimal value of (C.6) is equal to $f(\mathbf{u})$. Hence, we have written $f(\mathbf{u})$ as pointwise supremum of functions. Therefore, subgradient of $f(\mathbf{u})$ can be achieved by taking the derivative of the objective function in (C.6) at the optimal solution (see Section 2.10 in [MN13]). Note that the derivative of this objective function at the optimal solution is equal to $\boldsymbol{\rho}^* = 2(\mathbf{u} - \mathbf{y}^*) = 2(\mathbf{u} - \mathbf{\Pi}_{\text{conv}(\mathbf{X})}(\mathbf{u}))$ (where we used Eq. (C.5)). Since the dual optimum is unique (by strong convexity in $\boldsymbol{\rho}$), the function $\mathbf{u} \mapsto \mathcal{D}(\mathbf{u}, \mathbf{X})$ is differentiable with gradient given by Eq. (C.1). \square

Proposition C.2. *Let $\mathbf{u} \in \mathbb{R}^d$ and $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$, and assume that the rows of $\mathbf{H}_0 \in \mathbb{R}^{r \times d}$ are affine independent. Then the function $\mathbf{H} \mapsto \mathcal{D}(\mathbf{u}, \mathbf{H})$ is differentiable at \mathbf{H}_0 with gradient*

$$\nabla_{\mathbf{H}} \mathcal{D}(\mathbf{u}, \mathbf{H}_0) = 2\boldsymbol{\pi}_0(\mathbf{\Pi}_{\text{conv}(\mathbf{H}_0)}(\mathbf{u}) - \mathbf{u})^\top, \quad \boldsymbol{\pi}_0 = \arg \min_{\boldsymbol{\pi} \in \Delta^r} \left\| \mathbf{H}_0^\top \boldsymbol{\pi} - \mathbf{u} \right\|_2^2. \quad (\text{C.7})$$

Proof. We will denote by \mathbf{G} the right hand side of Eq. (C.7). For $\mathbf{V} \in \mathbb{R}^{r \times d}$, we have

$$\mathcal{D}(\mathbf{u}, \mathbf{H}_0 + \mathbf{V}) = \min_{\boldsymbol{\pi} \in \Delta^r} \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi} - \mathbf{u} \right\|_2^2. \quad (\text{C.8})$$

Note that $(\mathbf{H}_0 + \mathbf{V})$ has affinely independent rows for \mathbf{V} in a neighborhood of $\mathbf{0}$, and hence has a unique minimizer there, that we will denote by $\boldsymbol{\pi}_{\mathbf{V}}$. By optimality of $\boldsymbol{\pi}_{\mathbf{V}}$, we have

$$\mathcal{D}(\mathbf{u}, \mathbf{H}_0 + \mathbf{V}) - \mathcal{D}(\mathbf{u}, \mathbf{H}_0) = \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi}_{\mathbf{V}} - \mathbf{u} \right\|_2^2 - \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi}_0 - \mathbf{u} \right\|_2^2 \quad (\text{C.9})$$

$$\leq \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi}_0 - \mathbf{u} \right\|_2^2 - \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi}_0 - \mathbf{u} \right\|_2^2 \quad (\text{C.10})$$

$$= \langle \mathbf{G}, \mathbf{V} \rangle + \|\mathbf{V}\boldsymbol{\pi}_0\|_2^2. \quad (\text{C.11})$$

On the other hand, by optimality of $\boldsymbol{\pi}_0$,

$$\mathcal{D}(\mathbf{u}, \mathbf{H}_0 + \mathbf{V}) - \mathcal{D}(\mathbf{u}, \mathbf{H}_0) \geq \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi}_{\mathbf{V}} - \mathbf{u} \right\|_2^2 - \left\| (\mathbf{H}_0 + \mathbf{V})^\top \boldsymbol{\pi}_{\mathbf{V}} - \mathbf{u} \right\|_2^2 \quad (\text{C.12})$$

$$= \langle 2\boldsymbol{\pi}_{\mathbf{V}}(\mathbf{\Pi}_{\text{conv}(\mathbf{H}_0)}(\mathbf{u}) - \mathbf{u})^\top, \mathbf{V} \rangle + \|\mathbf{V}\boldsymbol{\pi}_{\mathbf{V}}\|_2^2 \quad (\text{C.13})$$

$$= \langle \mathbf{G}, \mathbf{V} \rangle + 2\langle (\boldsymbol{\pi}_{\mathbf{V}} - \boldsymbol{\pi}_0)(\mathbf{\Pi}_{\text{conv}(\mathbf{H}_0)}(\mathbf{u}) - \mathbf{u})^\top, \mathbf{V} \rangle + \|\mathbf{V}\boldsymbol{\pi}_{\mathbf{V}}\|_2^2. \quad (\text{C.14})$$

Letting $R(\mathbf{V}) = |\mathcal{D}(\mathbf{u}, \mathbf{H}_0 + \mathbf{V}) - \mathcal{D}(\mathbf{u}, \mathbf{H}_0) - \langle \mathbf{G}, \mathbf{V} \rangle|$ denote the residual, we get

$$\frac{R(\mathbf{V})}{\|\mathbf{V}\|_F} \leq \|\mathbf{\Pi}_{\text{conv}(\mathbf{H}_0)}(\mathbf{u}) - \mathbf{u}\|_2 \|\boldsymbol{\pi}_{\mathbf{V}} - \boldsymbol{\pi}_0\|_2 + \|\mathbf{V}\|_F (\|\boldsymbol{\pi}_{\mathbf{V}}\|_2 + \|\boldsymbol{\pi}_0\|_2). \quad (\text{C.15})$$

Note that $\boldsymbol{\pi}_{\mathbf{V}}$ must converge to $\boldsymbol{\pi}_0$ as $\mathbf{V} \rightarrow 0$ because $\boldsymbol{\pi}_0$ is the unique minimizer for $\mathbf{V} = \mathbf{0}$. Hence we get $R(\mathbf{V})/\|\mathbf{V}\|_F \rightarrow 0$ as $\|\mathbf{V}\|_F \rightarrow 0$, which proves our claim. \square

D Proof of Proposition 4.1

We use the results of [BST14] to prove Proposition 4.1. We refer the reader to [BST14] for the definitions of the technical terms in this section. First, consider the function

$$f(\mathbf{H}) = \lambda \mathcal{D}(\mathbf{H}, \mathbf{X}). \quad (\text{D.1})$$

Note that using the main theorem of polytope theory (Theorem 1.1 in [Zie12]), we can write

$$\text{conv}(\mathbf{X}) = \left\{ \mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{a}_i, \mathbf{x} \rangle \leq b_i \text{ for } 1 \leq i \leq m \right\} \quad (\text{D.2})$$

for some $\mathbf{a}_i \in \mathbb{R}^d$, $b_i \in \mathbb{R}$ and a finite m . Hence, using the definition of the semi-algebraic sets (see Definition 5 in [BST14]), the set $\text{conv}(\mathbf{X})$ is semi-algebraic. Therefore, the function $f(\mathbf{H})$ which is proportional to the sum of squared ℓ_2 distances of the rows of \mathbf{H} from a semi-algebraic set, is a semi-algebraic function (See Appendix in [BST14]). Further, the function

$$g(\mathbf{W}) = \sum_{i=1}^n \mathbf{I}(\mathbf{w}_i \in \Delta^r) \quad (\text{D.3})$$

is the sum of indicator functions of semi-algebraic sets (Note that using the same argument used for $\text{conv}(\mathbf{X})$, Δ^r is semi-algebraic). Therefore, the function g is semi-algebraic (See Appendix in [BST14]). In addition, the function

$$h(\mathbf{H}, \mathbf{W}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 \quad (\text{D.4})$$

is a polynomial. Hence, it is semi-algebraic. Therefore, we deduce that the function

$$\Psi(\mathbf{H}, \mathbf{W}) = f(\mathbf{H}) + g(\mathbf{W}) + h(\mathbf{H}, \mathbf{W}) \quad (\text{D.5})$$

is semi-algebraic. In addition, since Δ^r is closed, Ψ is proper and lower semi-continuous. Therefore, $\Psi(\mathbf{H}, \mathbf{W})$ is a KL function (See Theorem 3 in [BST14]).

Now, we will show that the Assumptions 1,2 in [BST14] hold for our algorithm. First, note that since Δ^r is closed, the functions $f(\mathbf{H})$ and $g(\mathbf{W})$ are proper and lower semi-continuous. Further, $f(\mathbf{H}) \geq 0$, $g(\mathbf{W}) \geq 0$, $h(\mathbf{H}, \mathbf{W}) \geq 0$ for all $\mathbf{H} \in \mathbb{R}^{r \times d}$, $\mathbf{W} \in \mathbb{R}^{n \times r}$. In addition, the function $h(\mathbf{H}, \mathbf{W})$ is C^2 . Therefore, it is Lipschitz continuous over the bounded subsets of $\mathbb{R}^{r \times d} \times \mathbb{R}^{n \times r}$. Also, the partial derivatives of $h(\mathbf{H}, \mathbf{W})$ are

$$\nabla_{\mathbf{H}} h(\mathbf{H}, \mathbf{W}) = 2\mathbf{W}^\top (\mathbf{W}\mathbf{H} - \mathbf{X}), \quad (\text{D.6})$$

$$\nabla_{\mathbf{W}} h(\mathbf{H}, \mathbf{W}) = 2(\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^\top. \quad (\text{D.7})$$

It can be seen that for any fixed \mathbf{W} , the function $\mathbf{H} \mapsto \nabla_{\mathbf{H}} h(\mathbf{H}, \mathbf{W})$ is Lipschitz continuous with moduli $L_1(\mathbf{W}) = 2\|\mathbf{W}^\top \mathbf{W}\|_F$. Similarly, for any fixed \mathbf{H} , the function $\mathbf{W} \mapsto \nabla_{\mathbf{W}} h(\mathbf{H}, \mathbf{W})$ is

Lipschitz continuous with moduli $L_2(\mathbf{H}) = 2\|\mathbf{H}\mathbf{H}^\top\|_F$. Note that since in each iteration of the algorithm the rows of \mathbf{W}^k are in Δ^r . Hence,

$$\inf \left\{ L_1(\mathbf{W}^k) : k \in \mathbb{N} \right\} \geq \lambda_1^- \quad \sup \left\{ L_1(\mathbf{W}^k) : k \in \mathbb{N} \right\} \leq \lambda_1^+ \quad (\text{D.8})$$

for some some positive constants λ_1^-, λ_1^+ . In addition, note that because the PALM algorithm is a descent algorithm, i.e., $\Psi(\mathbf{H}^k, \mathbf{W}^k) \leq \Psi(\mathbf{H}^{k-1}, \mathbf{W}^{k-1})$ for $k \in \mathbb{N}$, and since $f(\mathbf{H}) \rightarrow \infty$ as $\|\mathbf{H}\|_F \rightarrow \infty$, the value of $L_2(\mathbf{H}^k) = \|\mathbf{H}^k \mathbf{H}^{k\top}\|_F$ remains bounded in every iteration. Finally, note that by taking $\gamma_2^k > \max \left\{ \left\| \mathbf{H}^{k+1} \mathbf{H}^{k+1\top} \right\|_F, \varepsilon \right\}$ for some constant $\varepsilon > 0$, we make sure that the steps in the PALM algorithm remain well defined (See Remark 3(iii) in [BST14]). Hence, we have shown that the assumptions of Theorem 1 in [BST14] hold. Therefore, using this theorem, the sequence $\{\mathbf{H}^k, \mathbf{W}^k\}_{k \in \mathbb{N}}$ generated by the iterations in (4.7) - (4.9) has a finite length and it converges to a stationary point $(\mathbf{H}^*, \mathbf{W}^*)$ of Ψ .

E Other optimization algorithms

Apart from the proximal alternating linearized minimization discussed in Section 4.2, we experimented with two other algorithms, obtaining comparable results. For the sake of completeness, we describe these algorithms here.

E.1 Stochastic gradient descent

Using any of the initializations discussed in Section 4.1 we iterate

$$\mathbf{H}^{(t+1)} = \mathbf{H}^{(0)} - \gamma_t \mathbf{G}^{(t)}. \quad (\text{E.1})$$

The step size γ_t is selected by backtracking line search. Ideally, the direction $\mathbf{G}^{(t)}$ can be taken to be equal to $\nabla \mathcal{R}_\lambda(\mathbf{H}^{(t)})$. However, for large datasets this is computationally impractical, since it requires to compute the projection of each data point onto the set $\text{conv}(\mathbf{H}^{(t)})$. In order to reduce the complexity of the direction calculation, we estimate this sum by subsampling. Namely, we draw a uniformly random set $S_t \subseteq [n]$ of fixed size $|S_t| = s \leq n$, and compute

$$\mathbf{G}^{(t)} = \frac{2n}{|S_t|} \sum_{i \in S_t} \alpha_i^* (\Pi_{\text{conv}(\mathbf{H})}(\mathbf{x}_i) - \mathbf{x}_i) + 2\lambda (\mathbf{H} - \Pi_{\text{conv}(\mathbf{X})}(\mathbf{H})), \quad (\text{E.2})$$

$$\alpha_i^* = \arg \min_{\alpha \in \Delta^r} \left\| \mathbf{H}^\top \alpha - \mathbf{x}_i^\top \right\|_2. \quad (\text{E.3})$$

E.2 Alternating minimization

This approach generalizes the original algorithm of [CB94]. We rewrite the objective as a function of $\mathbf{W} = (w_{ij})_{i \leq n, j \leq r}$, $\mathbf{h}_\ell \in \mathbb{R}^d$ and $\mathbf{A} = (\alpha_{\ell, i})_{\ell \leq r, i \leq n}$

$$\mathcal{R}_\lambda(\mathbf{H}) = \min_{\mathbf{W}, \mathbf{A}} F(\mathbf{H}, \mathbf{W}, \mathbf{A}), \quad (\text{E.4})$$

$$F(\mathbf{H}, \mathbf{W}, \mathbf{A}) = \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{\ell=1}^r w_{i\ell} \mathbf{h}_\ell \right\|_2^3 + \lambda \sum_{\ell=1}^r \left\| \mathbf{h}_\ell - \sum_{i=1}^n \alpha_{\ell, i} \mathbf{x}_i \right\|_2^2. \quad (\text{E.5})$$

The algorithm alternates between minimizing with respect to the weights $(w_{ij})_{i \leq n, j \leq r}$ (this can be done independently across $i \in \{1, \dots, n\}$) and minimizing over $(\mathbf{h}_\ell, \mathbf{A})$, which is done sequentially by

cycling over $\ell \in \{1, \dots, r\}$. Minimization over \mathbf{w}_i can be performed by solving a non-negative least squares problem. As shown in [CB94], minimization over $(\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell)$ is also equivalent to non-negative least squares. Indeed, by a simple calculation

$$F(\mathbf{H}, \mathbf{W}, \mathbf{A}) = w_\ell^{\text{tot}} \|\mathbf{h}_\ell - \mathbf{v}_\ell\|_2^2 + \lambda \left\| \mathbf{h}_\ell - \sum_{i=1}^n \alpha_{\ell,i} \mathbf{x}_i \right\|_2^2 + \tilde{F}(\mathbf{H}, \mathbf{W}, \mathbf{A}) \quad (\text{E.6})$$

$$= f_\ell(\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell; \mathbf{H}_{\neq \ell}, \mathbf{W}, \mathbf{A}) + \tilde{F}(\mathbf{H}, \mathbf{W}, \mathbf{A}). \quad (\text{E.7})$$

where $\mathbf{H}_{\neq \ell} = (\mathbf{h}_i)_{i \neq \ell, i \leq r}$, $\tilde{F}(\mathbf{H}, \mathbf{W}, \mathbf{A})$ does not depend on $(\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell)$, and we defined

$$w_\ell^{\text{tot}} \equiv \sum_{i=1}^n w_{i\ell}^2, \quad (\text{E.8})$$

$$\mathbf{v}_\ell \equiv \frac{1}{w_\ell^{\text{tot}}} \sum_{i=1}^n w_{i,\ell} \left\{ \mathbf{x}_i - \sum_{j \neq \ell, j \leq r} w_{ij} \mathbf{h}_j \right\}. \quad (\text{E.9})$$

It is therefore sufficient to minimize $f_\ell(\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell; \mathbf{H}_{\neq \ell}, \mathbf{W}, \mathbf{A})$ with respect to its first two arguments, which is equivalent to a non-negative least squares problem. This can be seen by minimizing $f_\ell(\dots)$ explicitly with respect to \mathbf{h}_ℓ and writing the resulting objective function.

The pseudocode for this algorithm is given below.

ALTERNATING MINIMIZATION

Input : Data $\{\mathbf{x}_i\}_{i \leq n}$, $\mathbf{x}_i \in \mathbb{R}^d$; integer r ; initial archetypes $\{\mathbf{h}_\ell^{(0)}\}_{1 \leq \ell \leq r}$; number of iterations T ;

Output : Archetype estimates $\{\hat{\mathbf{h}}_\ell^{(T)}\}_{1 \leq \ell \leq r}$;

- 1: For $\ell \in \{1, \dots, r\}$:
 - 2: Set $\boldsymbol{\alpha}_\ell^{(0)} = \arg \min_{\boldsymbol{\alpha} \in \Delta^n} \|\mathbf{h}_\ell^{(0)} - \mathbf{X} \boldsymbol{\alpha}\|_2$;
 - 3: For $t \in \{1, \dots, T\}$:
 - 4: Set $\mathbf{W}^t = \arg \min_{\mathbf{W}} F(\mathbf{H}^{t-1}, \mathbf{W}, \mathbf{A}^{t-1})$
 - 5: For $\ell \in \{1, \dots, r\}$:
 - 6: Set $\mathbf{h}_\ell^{(t)}, \boldsymbol{\alpha}_\ell^{(t)} = \arg \min_{\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell} f_\ell(\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell; \mathbf{H}_{< \ell}^t, \mathbf{H}_{> \ell}^{t-1}, \mathbf{W}^t, \mathbf{A}_{< \ell}^t, \mathbf{A}_{> \ell}^{t-1})$;
 - 7: End For;
 - 8: Return $\{\hat{\mathbf{h}}_\ell^{(T)}\}_{1 \leq \ell \leq r}$
-

Here $\mathbf{H}_{< \ell} = (\mathbf{h}_i)_{i < \ell}$, $\mathbf{H}_{> \ell} = (\mathbf{h}_i)_{\ell < i \leq r}$, and similarly for \mathbf{A} .