

# Supplement To “Non-negative Matrix Factorization via Archetypal Analysis”

Hamid Javadi\* and Andrea Montanari†

March 8, 2019

## A Further details on numerical experiments

The data in Figures 1, 4, and 5 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  $1186 \text{ cm}^{-1}$  and  $1530 \text{ cm}^{-1}$ , 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 generated 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

---

\*Department of Electrical Engineering, Stanford University

†Department of Electrical Engineering and Statistics, Stanford University

$1 \leq i \leq 9$ , a pair of indices  $(\ell_1, \ell_2)$  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  $(\ell_1, \ell_2, \ell_3)$  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,  $E^{r,d} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r\}$  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  $\Delta^m$  the  $(m - 1)$ -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  $Q_{r,n}$  the set of  $r$  by  $n$  row stochastic matrices. Namely,

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

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

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

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

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

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

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

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

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  $\boldsymbol{\Pi}_{\mathcal{C}}(\mathbf{x})$ . Namely,

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

Further, we use  $D(\mathbf{x}, \mathcal{C})$  to denote the distance of  $\mathbf{x}$  from  $\mathcal{C}$ , i.e.

$$D(\mathbf{x}, \mathcal{C}) = \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.9})$$

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 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

$$\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha \mu}{30 r^{3/2}}, \quad (\text{B.10})$$

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

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_*^2 r^5}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (\text{B.11})$$

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

Further, if

$$\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha \mu}{330 \kappa(\mathbf{H}_0) r^{5/2}}, \quad (\text{B.12})$$

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

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_{**}^2 r^4}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (\text{B.13})$$

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

### B.3 Proof

### B.4 Proof strategy

Before providing a detailed proof, let us explain the proof scheme and main intuition:

1. Notice that  $\widehat{\mathbf{H}}$  minimizes the distance  $\mathcal{D}(\mathbf{H}; \mathbf{X})$  from the data, among all the possible sets of archetypes that can explain the data itself (in the sense that  $\mathbf{x}_i$  is close to  $\text{conv}(\mathbf{H})$ ), cf. Eq. (2.5). Using the fact that  $\mathbf{X}$  is the perturbed version of  $\mathbf{X}_0$ , we show in Lemma B.5 below that

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.14})$$

2. In absence of separability, there might be multiple (non-equivalent) sets of archetypes that minimize  $\mathcal{D}(\mathbf{H}; \mathbf{X})$ , and hence reconstruction is fundamentally non-unique. We combine the separability assumption with the bound in step 1 to bound

$$\alpha(\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}) \quad (\text{B.15})$$

in terms of  $\delta, \kappa(\widehat{\mathbf{H}}), \sigma_{\max}(\widehat{\mathbf{H}})$ .

3. The last step gives us an error in terms of distances between convex hulls,  $\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)$ ,  $\mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})$ . Lemma B.2 translates this into an upper bound on  $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})$  in terms of  $\delta, \kappa(\widehat{\mathbf{H}}), \sigma_{\max}(\widehat{\mathbf{H}})$ .
4. Finally, we want to translate this upper bound in terms of geometric properties of the true archetypes (as opposed to properties of the estimated archetypes, namely  $\kappa(\widehat{\mathbf{H}}), \sigma_{\max}(\widehat{\mathbf{H}})$ ). In Lemmas B.3, B.4 we bound the quantities  $\kappa(\widehat{\mathbf{H}}), \sigma_{\max}(\widehat{\mathbf{H}})$  in terms of  $\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2}, \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}$  as well as  $\kappa(\mathbf{H}_0), \sigma_{\max}(\mathbf{H}_0)$ .
5. Finally, we aggregate the last three steps to finish the proof.

### B.4.1 Lemmas

In the following two lemmas we bound the notion we use for estimation error  $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})$ , that we defined in (3.2), in terms of  $\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)$ ,  $\mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})$ .

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

$$\gamma_{\mathcal{C}} = \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{v} \in \mathcal{C}, \|\mathbf{v}\|_2=1} \langle \mathbf{u}, \mathbf{v} \rangle. \quad (\text{B.16})$$

*We have*

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

An illustration of this lemma in the case of  $\mathcal{R} \subset \mathcal{C}$  is given in Figure 1. 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  $\gamma_{\mathcal{C}}$ .

*Proof.* Note that the claim is trivial in the case where  $\gamma_{\mathcal{C}} \leq 0$ , as by Cauchy-Schwarz inequality  $\gamma_{\mathcal{C}} \geq -1$  and hence the left hand side of (B.17) is nonnegative. Therefore, we focus on the case where  $\gamma_{\mathcal{C}} > 0$ , i.e.  $\mathcal{C}^*$ , the dual cone of  $\mathcal{C}$  has a nonempty interior. We write

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

Replacing

$$\mathbf{x} = \Pi_{\mathcal{C}}(\mathbf{x}) + (\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})), \quad (\text{B.19})$$

we get

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \langle \mathbf{u}, \Pi_{\mathcal{C}}(\mathbf{x}) + (\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})) \rangle \quad (\text{B.20})$$

$$\geq \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \langle \mathbf{u}, \Pi_{\mathcal{C}}(\mathbf{x}) \rangle - \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.21})$$

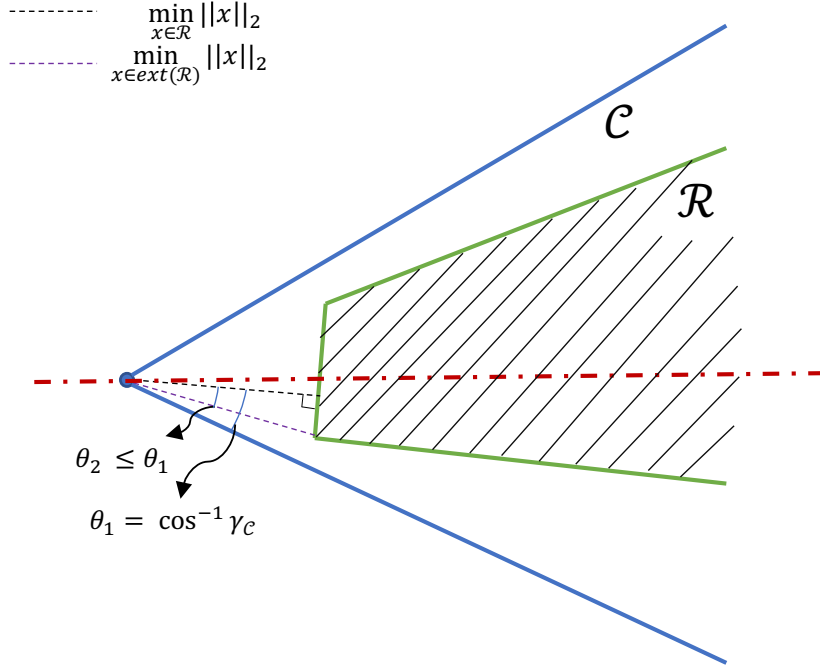


Figure 1: Picture of Lemma B.1, in the case,  $\mathcal{R} \subset \mathcal{C}$ .

Note that  $\tilde{\mathcal{R}}^*$ , the dual cone of the set  $\tilde{\mathcal{R}} \subseteq \mathcal{C}$  where  $\tilde{\mathcal{R}} := \{\Pi_{\mathcal{C}}(\mathbf{x}) ; \mathbf{x} \in \text{ext}(\mathcal{R})\}$ , contains  $\mathcal{C}^*$  and hence it has a nonempty interior. Therefore, the maximizer of the first term in the right hand side of (B.21) is in  $\tilde{\mathcal{R}}^*$  and

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \max_{\mathbf{u} \in \tilde{\mathcal{R}}^*; \|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \langle \mathbf{u}, \Pi_{\mathcal{C}}(\mathbf{x}) \rangle - \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2 \quad (\text{B.22})$$

$$= \max_{\mathbf{u} \in \tilde{\mathcal{R}}^*; \|\mathbf{u}\|_2=1} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \left[ \left\langle \mathbf{u}, \frac{\Pi_{\mathcal{C}}(\mathbf{x})}{\|\Pi_{\mathcal{C}}(\mathbf{x})\|_2} \right\rangle \|\Pi_{\mathcal{C}}(\mathbf{x})\|_2 \right] - \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.23})$$

Since for  $\mathbf{u} \in \tilde{\mathcal{R}}^*$ ,  $\mathbf{x} \in \text{ext}(\mathcal{R})$  we have  $\langle \mathbf{u}, \Pi_{\mathcal{C}}(\mathbf{x}) \rangle \geq 0$ , using the definition of  $\gamma_{\mathcal{C}}$ , we have

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \gamma_{\mathcal{C}} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\Pi_{\mathcal{C}}(\mathbf{x})\|_2 - \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.24})$$

Note that

$$\|\Pi_{\mathcal{C}}(\mathbf{x})\|_2 \geq \|\mathbf{x}\|_2 - \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2. \quad (\text{B.25})$$

Therefore,

$$\min_{\mathbf{x} \in \mathcal{R}} \|\mathbf{x}\|_2 \geq \gamma_{\mathcal{C}} \min_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x}\|_2 - (1 + \gamma_{\mathcal{C}}) \max_{\mathbf{x} \in \text{ext}(\mathcal{R})} \|\mathbf{x} - \Pi_{\mathcal{C}}(\mathbf{x})\|_2, \quad (\text{B.26})$$

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*

$$\mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \sqrt{2}\kappa(\mathbf{H}_0)\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} + (1 + \sqrt{2})\sqrt{r}\kappa(\mathbf{H}_0)\mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2}. \quad (\text{B.27})$$

*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.,

$$\mathcal{C}_1 = \left\{ \mathbf{v} \in \mathbb{R}^d; \mathbf{v} = \sum_{i=2}^r v_i(\mathbf{e}_i - \mathbf{e}_1), v_i \geq 0 \right\}. \quad (\text{B.28})$$

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

$$\mathbf{v} = (-\langle \mathbf{1}, \mathbf{x} \rangle, \mathbf{x}, 0, 0, \dots, 0), \quad (\text{B.29})$$

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

$$\|\mathbf{x}\|_2^2 + \langle \mathbf{1}, \mathbf{x} \rangle^2 = 1. \quad (\text{B.30})$$

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  $j = 1, 2, \dots, r$

$$\begin{aligned} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2 &\geq \gamma \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2 \\ &\quad - (1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \mathbb{R}_{\geq 0}^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{e}_j - \mathbf{E}_{r,d}^\top \mathbf{q} + \mathbf{e}_j \langle \mathbf{1}, \mathbf{q} \rangle\|_2 \end{aligned} \quad (\text{B.31})$$

$$\geq \gamma \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2 - (1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{E}_{r,d}^\top \mathbf{q}\|_2. \quad (\text{B.32})$$

Hence, for  $j = 1, 2, \dots, r$

$$\min_{\mathbf{q} \in \Delta^r} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2 + (1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{E}_{r,d}\|_2 \geq \gamma \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2. \quad (\text{B.33})$$

For simplicity, define  $\mathbf{a}, \mathbf{b} \in \mathbb{R}_{\geq 0}^r$ ,  $c \geq 0$  as

$$a_j \equiv \min_{\mathbf{q} \in \Delta^r} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2, b_j := \gamma \min_{\mathbf{q} \in E^{r,r}} \|\mathbf{e}_j - \mathbf{H}^\top \mathbf{q}\|_2, c := (1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{E}_{r,d}\|_2. \quad (\text{B.34})$$

Using triangle inequality,

$$\|\mathbf{a}\|_2 + c\sqrt{r} \geq \|\mathbf{a} + c\mathbf{1}\|_2. \quad (\text{B.35})$$

Using definition of  $\mathbf{a}, \mathbf{b}, c$ , (B.33) implies that  $a_j + c \geq b_j$ . Therefore, since  $a_j, b_j, c \geq 0$ ,

$$\|\mathbf{a} + c\mathbf{1}\|_2 \geq \|\mathbf{b}\|_2. \quad (\text{B.36})$$

Hence, using (B.35)

$$\|\mathbf{a}\|_2 \geq \|\mathbf{b}\|_2 - c\sqrt{r}. \quad (\text{B.37})$$

In other words,

$$\min_{\mathbf{Q} \in \mathcal{Q}_r} \|\mathbf{E}_{r,d} - \mathbf{Q}\mathbf{H}\|_F \geq \gamma \min_{\mathbf{Q} \in \mathcal{S}_r} \|\mathbf{E}_{r,d} - \mathbf{Q}\mathbf{H}\|_F - (1 + \gamma)\sqrt{r} \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{E}_{r,d}^\top \mathbf{q}\|_2. \quad (\text{B.38})$$

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

$$\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} = \min_{\mathbf{Q} \in \mathcal{Q}_r} \|\mathbf{H}_0 - \mathbf{Q}\mathbf{H}\|_F = \min_{\mathbf{Q} \in \mathcal{Q}_r} \|(\mathbf{E}_{r,d} - \mathbf{Q}\mathbf{Y})\mathbf{M}\|_F \quad (\text{B.39})$$

$$\geq \sigma_{\min}(\mathbf{M}) \min_{\mathbf{Q} \in \mathcal{Q}_r} \|\mathbf{E}_{r,d} - \mathbf{Q}\mathbf{Y}\|_F \quad (\text{B.40})$$

$$\geq \gamma \sigma_{\min}(\mathbf{M}) \min_{\mathbf{Q} \in \mathcal{S}_r} \|\mathbf{E}_{r,d} - \mathbf{Q}\mathbf{Y}\|_F - \sigma_{\min}(\mathbf{M})\sqrt{r}(1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{Y}_{i,\cdot}^\top - \mathbf{E}_{r,d}^\top \mathbf{q}\|_2 \quad (\text{B.41})$$

$$= \gamma \sigma_{\min}(\mathbf{M}) \min_{\mathbf{Q} \in \mathcal{S}_r} \|(\mathbf{H}_0 - \mathbf{Q}\mathbf{H})\mathbf{M}^{-1}\|_F - \sigma_{\min}(\mathbf{M})\sqrt{r}(1 + \gamma) \max_{i \in [r]} \min_{\mathbf{q} \in \Delta^r} \|(\mathbf{M}^{-1})^\top (\mathbf{H}_{i,\cdot}^\top - \mathbf{H}_0^\top \mathbf{q})\|_2. \quad (\text{B.42})$$

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

$$\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \geq \frac{\gamma}{\kappa(\mathbf{H}_0)} \mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} - (1 + \gamma)\sqrt{r} \min_{\mathbf{q} \in \Delta^r} \|\mathbf{H}_{i,\cdot}^\top - \mathbf{H}_0^\top \mathbf{q}\|_2 \quad (\text{B.43})$$

$$\geq \frac{\gamma}{\kappa(\mathbf{H}_0)} \mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} - (1 + \gamma)\sqrt{r}(\mathbf{H}, \mathbf{H}_0)^{1/2}. \quad (\text{B.44})$$

Therefore,

$$\mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \frac{\kappa(\mathbf{H}_0)}{\gamma} \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} + \frac{(1 + \gamma)\kappa(\mathbf{H}_0)\sqrt{r}}{\gamma} \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2}. \quad (\text{B.45})$$

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

$$\mathcal{L}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \sqrt{2}\kappa(\mathbf{H}_0)\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} + (1 + \sqrt{2})\sqrt{r}\kappa(\mathbf{H}_0)\mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} \quad (\text{B.46})$$

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*

$$\sigma_{\max}(\mathbf{H}) \leq \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} + \sqrt{r} \sigma_{\max}(\mathbf{H}_0), \quad (\text{B.47})$$

*In addition, if*

$$\mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \frac{\sigma_{\min}(\mathbf{H}_0)}{2}, \quad (\text{B.48})$$

*then*

$$\kappa(\mathbf{H}) \leq \frac{2r \sigma_{\max}(\mathbf{H}_0) + 2 \mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} \sqrt{r}}{\sigma_{\min}(\mathbf{H}_0)}. \quad (\text{B.49})$$

*Further, if*

$$\mathcal{D}(\mathbf{H}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \mathbf{H})^{1/2} \leq \frac{\sigma_{\min}(\mathbf{H}_0)}{6\sqrt{r}}, \quad (\text{B.50})$$

*then*

$$\sigma_{\max}(\mathbf{H}) \leq 2\sigma_{\max}(\mathbf{H}_0), \quad (\text{B.51})$$

$$\kappa(\mathbf{H}) \leq (7/2)\kappa(\mathbf{H}_0). \quad (\text{B.52})$$

*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

$$\begin{aligned} \mathbf{H}_0 &= \mathbf{P}\mathbf{H} + \mathbf{A}_2; \quad \|\mathbf{A}_2\|_F = \mathcal{D}_2, \\ \mathbf{H} &= \mathbf{R}\mathbf{H}_0 + \mathbf{A}_1; \quad \|\mathbf{A}_1\|_F = \mathcal{D}_1, \end{aligned} \quad (\text{B.53})$$

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,

$$\begin{aligned} \sigma_{\max}(\mathbf{P})\sigma_{\min}(\mathbf{H}) &\geq \sigma_{\min}(\mathbf{P}\mathbf{H}) = \sigma_{\min}(\mathbf{H}_0 - \mathbf{A}_2) \\ &\geq \sigma_{\min}(\mathbf{H}_0) - \sigma_{\max}(\mathbf{A}_2) \geq \sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2. \end{aligned} \quad (\text{B.54})$$

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

$$\sigma_{\max}(\mathbf{P}) \leq \|\mathbf{P}\|_F = \left( \sum_{i=1}^r \|\mathbf{P}_{i,\cdot}\|_2^2 \right)^{1/2} \leq \left( \sum_{i=1}^r \|\mathbf{P}_{i,\cdot}\|_1^2 \right)^{1/2} \leq \sqrt{r}. \quad (\text{B.55})$$

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

$$\sigma_{\min}(\mathbf{H}) \geq \frac{\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2}{\sqrt{r}}. \quad (\text{B.56})$$

In addition, using (B.53)

$$\begin{aligned} \sigma_{\max}(\mathbf{H}) &= \sigma_{\max}(\mathbf{R}\mathbf{H}_0 + \mathbf{A}_1) \leq \sigma_{\max}(\mathbf{R}\mathbf{H}_0) + \sigma_{\max}(\mathbf{A}_1) \\ &\leq \sigma_{\max}(\mathbf{R})\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_1 \leq \sqrt{r}\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_1. \end{aligned} \quad (\text{B.57})$$

Hence, using (B.56), (B.57), for  $\mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)$  we have

$$\kappa(\mathbf{H}) = \frac{\sigma_{\max}(\mathbf{H})}{\sigma_{\min}(\mathbf{H})} \leq \frac{r\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_1\sqrt{r}}{\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2}. \quad (\text{B.58})$$

Thus, for  $\mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)/2$ , by replacing  $\mathcal{D}_2$  with  $\sigma_{\min}(\mathbf{H}_0)/2$  in (B.58) we get Eqs. (B.47), (B.49).

Now assume that  $\mathcal{D}_1 + \mathcal{D}_2 \leq \sigma_{\min}(\mathbf{H}_0)/(6\sqrt{r})$ . In this case, using (B.53) we have

$$\mathbf{H}_0 = \mathbf{P}(\mathbf{R}\mathbf{H}_0 + \mathbf{A}_1) + \mathbf{A}_2. \quad (\text{B.59})$$

Therefore,

$$(\mathbf{I} - \mathbf{P}\mathbf{R})\mathbf{H}_0 = \mathbf{P}\mathbf{A}_1 + \mathbf{A}_2, \quad (\text{B.60})$$

hence,

$$\mathbf{I} - \mathbf{P}\mathbf{R} = (\mathbf{P}\mathbf{A}_1 + \mathbf{A}_2)\mathbf{H}_0^\dagger \quad (\text{B.61})$$

and

$$\mathbf{P}\mathbf{R} = \mathbf{I} - \mathbf{P}\mathbf{A}_1\mathbf{H}_0^\dagger - \mathbf{A}_2\mathbf{H}_0^\dagger. \quad (\text{B.62})$$

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

$$\sigma_{\max}(\mathbf{H}_0^\dagger) = \sigma_{\min}(\mathbf{H}_0)^{-1}. \quad (\text{B.63})$$

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

$$R_{ii} = \|\mathbf{R}_{:,i}\|_\infty \geq \langle \mathbf{P}_{i,:}, \mathbf{R}_{:,i} \rangle = (\mathbf{P}\mathbf{R})_{ii} = 1 - (\mathbf{P}\mathbf{A}_1\mathbf{H}_0^\dagger)_{ii} - (\mathbf{A}_2\mathbf{H}_0^\dagger)_{ii} \quad (\text{B.64})$$

$$\geq 1 - \|(\mathbf{P}\mathbf{A}_1\mathbf{H}_0^\dagger)_{i,:}\|_2 - \|(\mathbf{A}_2\mathbf{H}_0^\dagger)_{i,:}\|_2 \quad (\text{B.65})$$

$$\geq 1 - \max_{\mathbf{u} \in \Delta^r} \|\mathbf{A}_1^\top \mathbf{u}\|_2 \sigma_{\max}(\mathbf{H}_0^\dagger) - \|(\mathbf{A}_2)_{i,:}\|_2 \sigma_{\max}(\mathbf{H}_0^\dagger) \quad (\text{B.66})$$

$$\geq 1 - \max_{\mathbf{u} \in \Delta^r} \|\mathbf{u}\|_2 \sigma_{\max}(\mathbf{A}_1) \sigma_{\max}(\mathbf{H}_0^\dagger) - \|\mathbf{A}_2\|_F \sigma_{\max}(\mathbf{H}_0^\dagger) \quad (\text{B.67})$$

$$\geq 1 - \frac{\mathcal{D}_1 + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0)}. \quad (\text{B.68})$$

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

$$R_{ji} \leq \frac{\mathcal{D}_1 + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0)}. \quad (\text{B.69})$$

Thus,

$$\langle \mathbf{P}_{i,:}, \mathbf{R}_{:,i} \rangle = R_{ii}P_{ii} + \sum_{j \neq i} P_{ij}R_{ji} \leq R_{ii}P_{ii} + \left( \max_{j \neq i} R_{ji} \right) \sum_{j \neq i} P_{ij} \quad (\text{B.70})$$

$$\leq P_{ii} + \frac{\mathcal{D}_1 + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0)} (1 - P_{ii}). \quad (\text{B.71})$$

Therefore, using (B.68),

$$P_{ii} \geq \frac{\sigma_{\min}(\mathbf{H}_0) - 2(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.72})$$

Thus, we can write

$$\mathbf{P} = \mathbf{I} + \Delta; \quad \|\Delta_{i,\cdot}\|_1 \leq \frac{2(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.73})$$

Therefore,

$$\sigma_{\max}(\Delta) \leq \|\Delta\|_F = \left( \sum_{i=1}^r \|\Delta_{i,\cdot}\|_2^2 \right)^{1/2} \leq \left( \sum_{i=1}^r \|\Delta_{i,\cdot}\|_1^2 \right)^{1/2} \leq \frac{2(\mathcal{D}_1 + \mathcal{D}_2)\sqrt{r}}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.74})$$

Hence,

$$\sigma_{\max}(\mathbf{P}) \leq 1 + \frac{2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}, \quad \sigma_{\min}(\mathbf{P}) \geq 1 - \frac{2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.75})$$

From (B.53) we have  $\sigma_{\min}(\mathbf{P}\mathbf{H}) \geq \sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2$ . Using  $\sigma_{\min}(\mathbf{P}\mathbf{H}) \leq \sigma_{\max}(\mathbf{P})\sigma_{\min}(\mathbf{H})$ , we get

$$\sigma_{\min}(\mathbf{H}) \geq \frac{(\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2)(\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2))}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2) + 2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.76})$$

Further, from (B.53) we have  $\sigma_{\max}(\mathbf{P}\mathbf{H}) \leq \sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_2$ . Therefore, using the fact that  $\sigma_{\max}(\mathbf{P}\mathbf{H}) \geq \sigma_{\min}(\mathbf{P})\sigma_{\max}(\mathbf{H})$ , we get

$$\sigma_{\max}(\mathbf{H}) \leq \frac{(\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_2)(\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2))}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2) - 2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}. \quad (\text{B.77})$$

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,

$$\kappa(\mathbf{H}) \leq \left( \frac{\sigma_{\max}(\mathbf{H}_0) + \mathcal{D}_2}{\sigma_{\min}(\mathbf{H}_0) - \mathcal{D}_2} \right) \left( 1 + \frac{4\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)}{\sigma_{\min}(\mathbf{H}_0) - (\mathcal{D}_1 + \mathcal{D}_2) - 2\sqrt{r}(\mathcal{D}_1 + \mathcal{D}_2)} \right) \quad (\text{B.78})$$

$$\leq \frac{6\kappa(\mathbf{H}_0) + 1}{5} \left( 1 + \frac{4}{3} \right) \leq \frac{42\kappa(\mathbf{H}_0) + 7}{15} < \frac{7\kappa(\mathbf{H}_0)}{2}, \quad (\text{B.79})$$

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*

$$\sigma_{\min}(\mathbf{H}) \geq \sqrt{2}(\mu - 2\delta). \quad (\text{B.80})$$

*Proof.* Let

$$\mathbf{X}'_{i,\cdot} = \Pi_{\text{conv}(\mathbf{H})}(\mathbf{X}_{i,\cdot}) \equiv \arg \min_{\mathbf{x} \in \text{conv}(\mathbf{H})} \|\mathbf{X}_{i,\cdot} - \mathbf{x}\|_2. \quad (\text{B.81})$$

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

$$\|(\mathbf{X}_0)_{i,\cdot} - \mathbf{X}'_{i,\cdot}\|_2 \leq \|(\mathbf{X}_0)_{i,\cdot} - \mathbf{X}_{i,\cdot}\|_2 + \|\mathbf{X}_{i,\cdot} - \mathbf{X}'_{i,\cdot}\|_2 \leq 2\delta. \quad (\text{B.82})$$

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

$$\mathcal{D}(\mathbf{x}_0, \mathbf{X}')^{1/2} = \min_{\mathbf{a} \in \Delta^n} \|\mathbf{X}_0^\top \mathbf{a}_0 - \mathbf{X}'^\top \mathbf{a}\|_2 \leq \|\mathbf{X}_0^\top \mathbf{a}_0 - \mathbf{X}'^\top \mathbf{a}_0\|_2 \quad (\text{B.83})$$

$$\leq \left( \sum_{i=1}^n (a_0)_i \right) \|(\mathbf{X}_0)_{i,\cdot} - \mathbf{X}'_{i,\cdot}\|_2 \leq 2\delta. \quad (\text{B.84})$$

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

$$\mu \mathbf{U} \mathbf{z} + \mathbf{z}_0 = \mathbf{X}_0^\top \mathbf{a}. \quad (\text{B.85})$$

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

$$l_{\mathbf{u},\mu} = \{\mathbf{z} : \mathbf{z} = \mathbf{z}_0 + \alpha \mathbf{u}, |\alpha| \leq \mu\} \subseteq \text{conv}(\mathbf{X}_0). \quad (\text{B.86})$$

Thus,

$$l_{\mathbf{u},\mu} \subseteq \mathbf{P}_{\mathbf{u}}(\text{conv}(\mathbf{X}_0)) \quad (\text{B.87})$$

where  $\mathbf{P}_u$  is the orthogonal projection onto the line containing  $l_{u,\mu}$ . Note that since  $\mathbf{P}_u(\cdot)$  is a nonexpansive mapping, using (B.84), for any  $\mathbf{x}_0 \in \text{conv}(\mathbf{X}_0)$  we have

$$D(\mathbf{P}_u(\mathbf{x}_0), \mathbf{P}_u(\text{conv}(\mathbf{X}'))^{1/2} \leq \mathcal{D}(\mathbf{x}_0, \mathbf{X}')^{1/2} \leq 2\delta. \quad (\text{B.88})$$

In other words, for any  $\mathbf{x}_0 \in \mathbf{P}_u(\text{conv}(\mathbf{X}_0))$ ,  $D(\mathbf{x}_0, \mathbf{P}_u(\text{conv}(\mathbf{X}')) \leq 2\delta$ . Therefore, using (B.86) for any  $\mathbf{u}$  in column space of  $\mathbf{U}$ , we have

$$l_{u,\mu-2\delta} \subseteq \mathbf{P}_u(\text{conv}(\mathbf{X}')). \quad (\text{B.89})$$

This implies that

$$\mathbf{z}_0 + \mathbf{U}\mathbf{B}_{r'}(\mu - 2\delta) \subseteq \text{conv}(\mathbf{X}') \subseteq \text{conv}(\mathbf{H}). \quad (\text{B.90})$$

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

$$(\mu - 2\delta)\mathbf{U}\mathbf{z} + \mathbf{z}_0 = \mathbf{H}^\top \mathbf{a}. \quad (\text{B.91})$$

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

$$(\mu - 2\delta)(\mathbf{H}^\top)^\dagger \mathbf{U}\mathbf{z} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0 = \mathbf{a}. \quad (\text{B.92})$$

Let

$$\mathbf{a}_1 = (\mu - 2\delta)(\mathbf{H}^\top)^\dagger \mathbf{U}\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0, \quad (\text{B.93})$$

$$\mathbf{a}_2 = -(\mu - 2\delta)(\mathbf{H}^\top)^\dagger \mathbf{U}\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0, \quad (\text{B.94})$$

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

Therefore, we have

$$\mathbf{a}_1 = (\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U})\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0, \quad (\text{B.95})$$

$$\mathbf{a}_2 = -(\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U})\mathbf{v} + (\mathbf{H}^\top)^\dagger \mathbf{z}_0. \quad (\text{B.96})$$

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

$$\|\mathbf{a}_1 - \mathbf{a}_2\|_2 = 2(\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U}). \quad (\text{B.97})$$

Note that

$$\|\mathbf{a}_1 - \mathbf{a}_2\|_2 \leq \sqrt{2}. \quad (\text{B.98})$$

Thus,

$$2(\mu - 2\delta)\sigma_{\max}((\mathbf{H}^\top)^\dagger \mathbf{U}) = \frac{2(\mu - 2\delta)}{\sigma_{\min}(\mathbf{H})} \leq \sqrt{2}. \quad (\text{B.99})$$

Hence,

$$\sigma_{\min}(\mathbf{H}) \geq \sqrt{2}(\mu - 2\delta). \quad (\text{B.100})$$

□

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*

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.101})$$

*Proof.* First note that  $\mathbf{X} = \mathbf{W}_0 \mathbf{H}_0 + \mathbf{Z}$  and the rows of  $\mathbf{W}_0 \mathbf{H}_0$  are in  $\text{conv}(\mathbf{H}_0)$ . Thus, since  $\max_i \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta$ , we have

$$\max_{i \leq n} \mathcal{D}(\mathbf{X}_{i,\cdot}, \mathbf{H}_0)^{1/2} \leq \max_{i \leq n} \left\| (\mathbf{X}_0 - \mathbf{W}_0 \mathbf{H}_0)_{i,\cdot} \right\|_2 = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta. \quad (\text{B.102})$$

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

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}). \quad (\text{B.103})$$

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

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

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

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

$$= \sum_{i=1}^r \left( \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\boldsymbol{\alpha}}_i\|_2^2 - 2 \langle \mathbf{Z}^\top \tilde{\boldsymbol{\alpha}}_i, \widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\boldsymbol{\alpha}}_i \rangle + \|\mathbf{Z}^\top \tilde{\boldsymbol{\alpha}}_i\|_2^2 \right). \quad (\text{B.107})$$

Using the fact that (by triangle inequality)  $\|\mathbf{Z}^\top \tilde{\boldsymbol{\alpha}}_i\|_2 \leq (\max_i \|\mathbf{Z}_{i,\cdot}\|_2) \|\tilde{\boldsymbol{\alpha}}_i\|_1 \leq \delta$ , we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq \sum_{i=1}^r \left( \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\boldsymbol{\alpha}}_i\|_2^2 - 2\delta \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\boldsymbol{\alpha}}_i\|_2 \right) \quad (\text{B.108})$$

$$\geq U^2 - 2\delta\sqrt{r}U \quad (\text{B.109})$$

where  $U^2 = \sum_{i=1}^r \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \tilde{\boldsymbol{\alpha}}_i\|_2^2$ . In addition, since  $U \geq 0$ ,  $\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq 0$ ,

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq (U^2 - 2\delta\sqrt{r}U) \mathbb{I}_{U \geq 2\delta\sqrt{r}}. \quad (\text{B.110})$$

Note that for  $U \geq 2\delta\sqrt{r}$ , the function  $U^2 - 2\delta\sqrt{r}U$  is increasing. Hence, since

$$U \geq \left( \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i} \|\widehat{\mathbf{H}}_{i,\cdot}^\top - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2^2 \right)^{1/2} = \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2}, \quad (\text{B.111})$$

we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}) \geq (U^2 - 2\delta\sqrt{r}U) \mathbb{I}_{U \geq 2\delta\sqrt{r}} \geq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0) - 2\delta\sqrt{r} \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2}. \quad (\text{B.112})$$

Therefore, letting  $(x)_+ := \max\{x, 0\}$  be the positive part of  $x$ , using the fact that for all  $x, a \geq 0$  we have  $(x^2 - ax)_+^{1/2} \geq x - a$ , we have

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X})^{1/2} \geq \left( \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0) - 2\delta\sqrt{r}\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \right)_+^{1/2} \geq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} - 2\delta\sqrt{r}. \quad (\text{B.113})$$

In addition,

$$\mathcal{D}(\mathbf{H}_0, \mathbf{X}) = \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^n} \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i - \mathbf{Z}^\top \boldsymbol{\alpha}_i\|_2^2 \quad (\text{B.114})$$

$$\leq \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^n} \{ \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2 + \|\mathbf{Z}^\top \boldsymbol{\alpha}_i\|_2 \}^2 \quad (\text{B.115})$$

$$\leq \sum_{i=1}^r \left\{ \min_{\boldsymbol{\alpha}_i \in \Delta^n} \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2 + \max_{\boldsymbol{\alpha}_i \in \Delta^n} \|\mathbf{Z}^\top \boldsymbol{\alpha}_i\|_2 \right\}^2 \quad (\text{B.116})$$

$$\leq \left\{ \left( \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^n} \|(\mathbf{H}_0)_{i,\cdot} - \mathbf{X}_0^\top \boldsymbol{\alpha}_i\|_2^2 \right)^{1/2} + \delta\sqrt{r} \right\}^2 \quad (\text{B.117})$$

$$\leq \left( \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \delta\sqrt{r} \right)^2. \quad (\text{B.118})$$

Hence,

$$\mathcal{D}(\mathbf{H}_0, \mathbf{X})^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \delta\sqrt{r}. \quad (\text{B.119})$$

Combining equations (B.113), (B.119), and (B.103), we get

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta\sqrt{r}. \quad (\text{B.120})$$

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 that  $\max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \delta$ . Then for  $\widehat{\mathbf{H}}$  the optimal solution of problem (2.5), we have

$$\alpha(\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}) \leq 2(1 + 2\alpha) \left[ r^{3/2} \delta \kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{\delta \sqrt{r}}{\mu} \sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1} \mathbf{z}_0^\top) \right] + 3\delta \sqrt{r} \quad (\text{B.121})$$

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

$$\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} \geq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + \alpha(\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2}). \quad (\text{B.122})$$

Note that Lemma B.5 implies

$$\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} \leq \mathcal{D}(\mathbf{H}_0, \mathbf{X}_0)^{1/2} + 3\delta \sqrt{r}. \quad (\text{B.123})$$

Therefore,

$$\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} \geq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} - 3\delta \sqrt{r} + \alpha(\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2}). \quad (\text{B.124})$$

Hence,

$$\alpha(\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2}) \leq \mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2} + 3\delta \sqrt{r}. \quad (\text{B.125})$$

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

$$\left[ \sum_{i=1}^n D(\widehat{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} - \left[ \sum_{i=1}^n D(\widetilde{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} \leq \left[ \sum_{i=1}^n \left( D(\widehat{\mathbf{H}}_{i,\cdot}, S)^{1/2} - D(\widetilde{\mathbf{H}}_{i,\cdot}, S)^{1/2} \right)^2 \right]^{1/2}. \quad (\text{B.126})$$

In addition, note that for  $\mathbf{x} \in \mathbb{R}^d$ ,  $S \subseteq \mathbb{R}^d$ ,  $D(\mathbf{x}, S)^{1/2} = \inf_{\mathbf{z} \in S} \|\mathbf{x} - \mathbf{z}\|_2$ . Hence, Using triangle inequality, for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,  $\mathbf{z} \in S$ ,

$$\|\mathbf{x} - \mathbf{z}\|_2 \leq \|\mathbf{y} - \mathbf{z}\|_2 + \|\mathbf{x} - \mathbf{y}\|_2. \quad (\text{B.127})$$

Taking infimum over  $\mathbf{z} \in S$  above and rearranging terms, we get

$$D(\mathbf{x}, S)^{1/2} - D(\mathbf{y}, S)^{1/2} \leq \|\mathbf{x} - \mathbf{y}\|_2. \quad (\text{B.128})$$

Similarly, by replacing the roles of  $\mathbf{x}, \mathbf{y}$  above we conclude that for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,  $S \subseteq \mathbb{R}^d$

$$|D(\mathbf{x}, S)^{1/2} - D(\mathbf{y}, S)^{1/2}| \leq \|\mathbf{x} - \mathbf{y}\|_2. \quad (\text{B.129})$$

In particular, we have

$$|D(\widetilde{\mathbf{H}}_{i,\cdot}, S)^{1/2} - D(\widehat{\mathbf{H}}_{i,\cdot}, S)^{1/2}| \leq \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2. \quad (\text{B.130})$$

Therefore, using (B.126) we have

$$\left[ \sum_{i=1}^n D(\widehat{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} - \left[ \sum_{i=1}^n D(\widetilde{\mathbf{H}}_{i,\cdot}, S) \right]^{1/2} \leq \left[ \sum_{i=1}^n \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2^2 \right]^{1/2} = \|\widehat{\mathbf{H}} - \widetilde{\mathbf{H}}\|_F. \quad (\text{B.131})$$

Hence,

$$|\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{X}_0)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{X}_0)^{1/2}| \leq \|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F \quad (\text{B.132})$$

and

$$|\mathcal{D}(\widetilde{\mathbf{H}}, \mathbf{H}_0)^{1/2} - \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2}| \leq \|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F. \quad (\text{B.133})$$

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

$$\mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}}) = \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widetilde{\mathbf{H}}^\top \boldsymbol{\alpha}_i\|_2^2 \quad (\text{B.134})$$

$$= \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \boldsymbol{\alpha}_i - (\widehat{\mathbf{H}} - \widetilde{\mathbf{H}})^\top \boldsymbol{\alpha}_i\|_2^2 \quad (\text{B.135})$$

$$\leq \sum_{i=1}^r \min_{\boldsymbol{\alpha}_i \in \Delta^r} \left\{ \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \boldsymbol{\alpha}_i\|_2 + \|(\widehat{\mathbf{H}} - \widetilde{\mathbf{H}})^\top \boldsymbol{\alpha}_i\|_2 \right\}^2 \quad (\text{B.136})$$

$$\leq \sum_{i=1}^r \left\{ \min_{\boldsymbol{\alpha} \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \boldsymbol{\alpha}\|_2 + \max_{\boldsymbol{\alpha} \in \Delta^r} \|(\widehat{\mathbf{H}} - \widetilde{\mathbf{H}})^\top \boldsymbol{\alpha}\|_2 \right\}^2 \quad (\text{B.137})$$

$$\leq \left\{ \left( \sum_{i=1}^r \min_{\boldsymbol{\alpha} \in \Delta^r} \|(\mathbf{H}_0)_{i,\cdot} - \widehat{\mathbf{H}}^\top \boldsymbol{\alpha}_i\|_2^2 \right)^{1/2} + \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \right\}^2 \quad (\text{B.138})$$

$$\leq \left( \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} + \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \right)^2. \quad (\text{B.139})$$

Thus,

$$|\mathcal{D}(\mathbf{H}_0, \widetilde{\mathbf{H}})^{1/2} - \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}| \leq \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \quad (\text{B.140})$$

Therefore, combining (B.125), (B.132), (B.133), (B.140), we get

$$\begin{aligned} \alpha \left( \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \right) &\leq (1 + \alpha) \|\widehat{\mathbf{H}} - \widetilde{\mathbf{H}}\|_F \\ &\quad + \alpha \sqrt{r} \max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 + 3\delta \sqrt{r}. \end{aligned} \quad (\text{B.141})$$

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

$$\mathcal{D}(\mathbf{X}_{i,\cdot}, \widehat{\mathbf{H}}) \leq \delta^2. \quad (\text{B.142})$$

Thus by triangle inequality,

$$\mathcal{D}((\mathbf{X}_0)_{i,\cdot}, \widehat{\mathbf{H}})^{1/2} \leq \|\mathbf{X}_{i,\cdot} - (\mathbf{X}_0)_{i,\cdot}\|_2 + \mathcal{D}(\mathbf{X}_{i,\cdot}, \widehat{\mathbf{H}})^{1/2} \leq 2\delta. \quad (\text{B.143})$$

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  $\mathbf{B}_d(2\delta, \widehat{\mathbf{H}})$ , we get

$$\text{conv}(\mathbf{X}_0) \subseteq \mathbf{B}_d(2\delta, \widehat{\mathbf{H}}) \cap \text{aff}(\mathbf{H}_0). \quad (\text{B.144})$$

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  $\tilde{\mathbf{Q}} \in \mathbb{R}^{r \times d}$  defined by

$$\tilde{Q}_{ii} = 1 + \xi, \quad \text{if } i = j \in \{1, 2, \dots, r\}, \quad (\text{B.145})$$

$$\tilde{Q}_{ij} = -\frac{\xi}{r-1}, \quad \text{if } i \neq j \in \{1, 2, \dots, r\}, \quad (\text{B.146})$$

$$\tilde{Q}_{ij} = 0, \quad \text{if } j \in \{r+1, r+2, \dots, d\} \quad (\text{B.147})$$

where  $\xi = 2r\delta_0$ . Note that for every  $\mathbf{y} \in \mathbf{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d})$ , 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  $\mathbf{y} \in \mathbf{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d})$ , we can write

$$\mathbf{y} = \boldsymbol{\pi} + \mathbf{x} \quad (\text{B.148})$$

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

$$\mathbf{y} = \sum_{i=1}^r \beta_i \tilde{\mathbf{Q}}_{i,\cdot} \quad (\text{B.149})$$

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

$$\beta_i = \frac{r-1}{r-1+\xi r}(\pi_i + x_i) + \frac{\xi}{r-1+\xi r}. \quad (\text{B.150})$$

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  $i \in \{1, 2, \dots, r\}$ ,

$$\beta_i \geq \frac{-2\delta_0(r-1) + \xi}{r-1+\xi r} = \frac{2\delta_0}{r-1+\xi r} \geq 0. \quad (\text{B.151})$$

In addition,

$$\sum_{i=1}^r \beta_i = \frac{r\xi}{r-1+\xi r} + \frac{r-1}{r-1+\xi r} \left( \sum_{i=1}^r (\pi_i + x_i) \right) = 1. \quad (\text{B.152})$$

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

$$\mathcal{B}_d(2\delta_0; \mathbf{E}_{r,d}) \cap \text{aff}(\mathbf{E}_{r,d}) \subseteq \text{conv}(\tilde{\mathbf{Q}}). \quad (\text{B.153})$$

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

$$\text{conv}(\widetilde{\mathbf{H}}) \supseteq \left[ \bigcup_{\mathbf{x} \in \text{conv}(\mathbf{E}_{r,d})} \mathbf{M}^\top \mathcal{B}_d(2\delta_0; \mathbf{x}) \right] \cap \text{aff}(\widehat{\mathbf{H}}) \quad (\text{B.154})$$

$$\supseteq \left[ \bigcup_{\mathbf{x} \in \text{conv}(\widehat{\mathbf{H}})} \mathcal{B}_d(2\delta_0 \sigma_{\min}(\mathbf{M}); \mathbf{x}) \right] \cap \text{aff}(\widehat{\mathbf{H}}) \quad (\text{B.155})$$

$$\supseteq \mathcal{B}_d(2\delta; \widehat{\mathbf{H}}) \cap \text{aff}(\mathbf{H}_0), \quad (\text{B.156})$$

provided that  $\delta_0 = \delta / \sigma_{\min}(\mathbf{M}) = \delta / \sigma_{\min}(\widehat{\mathbf{H}})$ . Hence, using (B.144) for this  $\delta_0$ ,  $\text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\widetilde{\mathbf{H}})$ . Note that for  $\tilde{\mathbf{Q}}$ , we have  $\|\tilde{\mathbf{Q}}_{i,\cdot} - \mathbf{e}_i\|_2 \leq 2r\delta_0$ . Thus,

$$\|\tilde{\mathbf{Q}} - \mathbf{E}_{r,d}\|_F \leq 2r^{3/2}\delta_0. \quad (\text{B.157})$$

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

$$\|\widetilde{\mathbf{H}} - \widehat{\mathbf{H}}\|_F = \|(\tilde{\mathbf{Q}} - \mathbf{E}_{r,d}) \mathbf{M}\|_F \leq 2r^{3/2}\delta_0 \sigma_{\max}(\mathbf{M}) = 2r^{3/2}\delta_0 \sigma_{\max}(\widehat{\mathbf{H}}) = 2r^{3/2}\delta \kappa(\widehat{\mathbf{H}}),$$

$$\max_{i \in [r]} \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 = \max_{i \in [r]} \|(\tilde{\mathbf{Q}}_{i,\cdot} - \mathbf{e}_i) \mathbf{M}\|_2 \leq 2r\delta_0 \sigma_{\max}(\mathbf{M}) = 2r\delta_0 \sigma_{\max}(\widehat{\mathbf{H}}) = 2r\delta \kappa(\widehat{\mathbf{H}}).$$

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,

$$(\mathbf{X}_0)_{i,\cdot} \in \mathcal{B}_d(2\delta, \mathbf{H}') \cap \text{aff}(\mathbf{H}'). \quad (\text{B.158})$$

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

$$\|\widetilde{\mathbf{H}} - \mathbf{H}'\|_F \leq 2r^{3/2}\delta\kappa(\mathbf{H}'), \quad (\text{B.159})$$

$$\max_{i \in [r]} \|\widetilde{\mathbf{H}}_{i,\cdot} - \mathbf{H}'_{i,\cdot}\|_2 \leq 2r\delta\kappa(\mathbf{H}'). \quad (\text{B.160})$$

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

$$\begin{aligned} \|\widetilde{\mathbf{H}}_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 &\leq \|\widetilde{\mathbf{H}}_{i,\cdot} - \mathbf{H}'_{i,\cdot}\|_2 + \|\mathbf{H}'_{i,\cdot} - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \\ &\leq 2r\delta\kappa(\mathbf{H}') + \|\mathbf{P}_0(\widehat{\mathbf{H}}_{i,\cdot}) - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \end{aligned} \quad (\text{B.161})$$

where  $\mathbf{P}_0$  is the 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  $\text{conv}(\mathbf{X}_0) \subseteq \text{B}_d(2\delta, \widehat{\mathbf{H}})$ , 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

$$\max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a}\|_2^2 = \max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} - (\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)^\top \mathbf{a}\|_2^2 \quad (\text{B.162})$$

$$\leq \max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} + \mathbf{z}_0 - \widehat{\mathbf{H}}^\top \mathbf{a}\|_2^2 \leq 4\delta^2. \quad (\text{B.163})$$

Now, using Cauchy-Schwarz inequality we can write

$$\max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\|\mathbf{a}\|_2 \leq 1} \|\mathbf{U}\mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a}\|_2^2 \leq \max_{\|\mathbf{z}\|_2 \leq \mu} \min_{\langle \mathbf{a}, \mathbf{1} \rangle = 1, \mathbf{a} \geq 0} \|\mathbf{U}\mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a}\|_2^2 \leq 4\delta^2. \quad (\text{B.164})$$

Note that,

$$\min_{\|\mathbf{a}\|_2 \leq 1} \|\mathbf{U}\mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a}\|_2^2 = \max_{\rho \geq 0} \min_{\mathbf{a}} \left\{ \|\mathbf{z}\|_2^2 - 2 \left\langle \mathbf{z}, \mathbf{U}^\top \bar{\mathbf{H}}^\top \mathbf{a} \right\rangle + \left\langle \mathbf{a}, (\bar{\mathbf{H}}\bar{\mathbf{H}}^\top + \rho\mathbf{I})\mathbf{a} \right\rangle - \rho \right\} \quad (\text{B.165})$$

$$= \max_{\rho \geq 0} \left\{ \|\mathbf{z}\|_2^2 - \left\langle \bar{\mathbf{H}}\mathbf{U}\mathbf{z}, (\bar{\mathbf{H}}\bar{\mathbf{H}}^\top + \rho\mathbf{I})^{-1} \bar{\mathbf{H}}\mathbf{U}\mathbf{z} \right\rangle - \rho \right\} \quad (\text{B.166})$$

Hence, using (B.164)

$$\mu^2 \max_{\rho \geq 0} \left\{ \lambda_{\max}(\mathbf{I} - \mathbf{U}^\top \bar{\mathbf{H}}^\top (\bar{\mathbf{H}} \bar{\mathbf{H}}^\top + \rho \mathbf{I})^{-1} \bar{\mathbf{H}} \mathbf{U}) - \rho \right\} \leq 4\delta^2. \quad (\text{B.167})$$

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

$$\mu^2 \lambda_{\max}(\mathbf{I} - \mathbf{U}^\top \bar{\mathbf{H}}^\top (\bar{\mathbf{H}} \bar{\mathbf{H}}^\top)^{-1} \bar{\mathbf{H}} \mathbf{U}) \leq 4\delta^2. \quad (\text{B.168})$$

Taking  $\bar{\mathbf{H}} = \tilde{\mathbf{U}} \Sigma \tilde{\mathbf{V}}^\top$ , the singular value decomposition of  $\bar{\mathbf{H}}$ , we have  $\sigma_{\max}(\bar{\mathbf{H}}) = \sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1} \mathbf{z}_0^\top) = \max_i \Sigma_{ii}$ . Letting  $\mathbf{U}^\top \tilde{\mathbf{V}} = \mathbf{Q}$ , we get

$$\max_{\rho \geq 0} \lambda_{\max}(\mathbf{I} - \mathbf{Q} \mathbf{Q}^\top) \leq \frac{4\delta^2}{\mu^2}. \quad (\text{B.169})$$

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

$$1 - q^2 \leq \frac{4\delta^2}{\mu^2}. \quad (\text{B.170})$$

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

$$\max_{i \in [r]} \| \mathbf{P}_0(\widehat{\mathbf{H}}_{i,\cdot}) - \widehat{\mathbf{H}}_{i,\cdot} \|_2 = \max_{i \in [r]} \| \mathbf{P}_U(\bar{\mathbf{H}}_{i,\cdot}) - \bar{\mathbf{H}}_{i,\cdot} \|_2 \quad (\text{B.171})$$

$$\leq \max_{\|\mathbf{a}\|_2 \leq 1} \| \mathbf{P}_U(\bar{\mathbf{H}}^\top \mathbf{a}) - \bar{\mathbf{H}}^\top \mathbf{a} \|_2 \quad (\text{B.172})$$

$$\leq \max_{\|\mathbf{a}\|_2 \leq 1} \| \mathbf{P}_U(\bar{\mathbf{H}}^\top \mathbf{a}) - \bar{\mathbf{H}}^\top \mathbf{a} \|_2 \quad (\text{B.173})$$

$$\leq \max_{\|\mathbf{a}\|_2 \leq 1} \min_{\mathbf{z}} \| \mathbf{U} \mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a} \|_2^2 \quad (\text{B.174})$$

where  $\mathbf{P}_U$  is the projector onto the column space of  $\mathbf{U}$ . Note that,

$$\max_{\|\mathbf{a}\|_2 \leq 1} \min_{\mathbf{z}} \|\mathbf{U}\mathbf{z} - \bar{\mathbf{H}}^\top \mathbf{a}\|_2^2 = \max_{\|\mathbf{a}\|_2 \leq 1} \left\{ -\left\langle \mathbf{a}, \bar{\mathbf{H}}\mathbf{U}\mathbf{U}^\top \bar{\mathbf{H}}^\top \mathbf{a} \right\rangle + \left\langle \mathbf{a}, \bar{\mathbf{H}}\bar{\mathbf{H}}^\top \mathbf{a} \right\rangle \right\} \quad (\text{B.175})$$

$$= \lambda_{\max}(\bar{\mathbf{H}}\bar{\mathbf{H}}^\top - \bar{\mathbf{H}}\mathbf{U}\mathbf{U}^\top \bar{\mathbf{H}}^\top) \quad (\text{B.176})$$

$$= \lambda_{\max}(\Sigma(\mathbf{I} - \mathbf{Q}^\top \mathbf{Q})\Sigma) \quad (\text{B.177})$$

$$\leq \sigma_{\max}(\bar{\mathbf{H}})^2 \lambda_{\max}(\mathbf{I} - \mathbf{Q}^\top \mathbf{Q}) \quad (\text{B.178})$$

$$\leq \sigma_{\max}(\bar{\mathbf{H}})^2 (1 - q^2) \leq \frac{4\sigma_{\max}(\bar{\mathbf{H}})^2 \delta^2}{\mu^2} \quad (\text{B.179})$$

where the last inequality follows from (B.170). This results in

$$\max_{i \in [r]} \|\mathbf{P}_0(\widehat{\mathbf{H}}_{i,\cdot}) - \widehat{\mathbf{H}}_{i,\cdot}\|_2 \leq \frac{2\sigma_{\max}(\bar{\mathbf{H}})\delta}{\mu} = \frac{2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta}{\mu}. \quad (\text{B.180})$$

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.161) we get

$$\max_{i \in [r]} \|\widehat{\mathbf{H}}_{i,\cdot} - \widetilde{\mathbf{H}}_{i,\cdot}\|_2 \leq 2r\delta\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta}{\mu}, \quad (\text{B.181})$$

$$\|\widehat{\mathbf{H}} - \widetilde{\mathbf{H}}\|_F \leq 2r^{3/2}\delta\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{2\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top)\delta\sqrt{r}}{\mu}. \quad (\text{B.182})$$

Replacing this in (B.141) completes the proof.  $\square$

#### B.4.2 Proof of Theorem 1

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 1 we have

$$\mathbf{z}_0 + \mathbf{U}\mathbf{B}_{r'}(\mu) \subseteq \text{conv}(\mathbf{X}_0) \subseteq \text{conv}(\mathbf{H}_0). \quad (\text{B.183})$$

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

$$\mu\sqrt{2} \leq \sigma_{\min}(\mathbf{H}_0) \leq \sigma_{\max}(\mathbf{H}_0). \quad (\text{B.184})$$

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,

$$\|\mathbf{z}_0\|_2 \leq \sigma_{\max}(\mathbf{H}_0) \|\boldsymbol{\alpha}_0\|_2 \leq \sigma_{\max}(\mathbf{H}_0). \quad (\text{B.185})$$

Note that

$$\sigma_{\max}(\widehat{\mathbf{H}} - \mathbf{1}\mathbf{z}_0^\top) \leq \sigma_{\max}(\widehat{\mathbf{H}}) + \sigma_{\max}(\mathbf{1}\mathbf{z}_0^\top) = \sigma_{\max}(\widehat{\mathbf{H}}) + \sqrt{r}\|\mathbf{z}_0\|_2. \quad (\text{B.186})$$

Therefore, using Lemma B.6 we have

$$\mathcal{D} \leq 2(1 + 2\alpha) \left( r^{3/2} \delta \kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) + \frac{\sigma_{\max}(\widehat{\mathbf{H}}) \delta r^{1/2}}{\mu} + \frac{r \delta \|\mathbf{z}_0\|_2}{\mu} \right) + 3\delta r^{1/2}. \quad (\text{B.187})$$

In addition, Lemma B.2 implies that

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{\kappa(\mathbf{H}_0)}{\alpha} (1 + \sqrt{2}) \sqrt{r} \mathcal{D}. \quad (\text{B.188})$$

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  $\mathbf{h} \in \mathbb{R}^d$

$$\mathcal{D}(\mathbf{P}_0(\mathbf{h}), \mathbf{H}_0) \leq D(\mathbf{P}_0(\mathbf{h}), \mathbf{P}_0(\Pi_{\text{conv}(\mathbf{H}_0)}(\mathbf{h}))) \leq D(\mathbf{h}, \Pi_{\text{conv}(\mathbf{H}_0)}(\mathbf{h})) = \mathcal{D}(\mathbf{h}, \mathbf{H}_0). \quad (\text{B.189})$$

Therefore,

$$\mathcal{D}(\mathbf{P}_0(\widehat{\mathbf{H}}), \mathbf{H}_0) \leq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0). \quad (\text{B.190})$$

First consider the case in which

$$\delta \leq \frac{\alpha \mu}{30 r^{3/2}}. \quad (\text{B.191})$$

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.190), we get

$$\sigma_{\max}(\widehat{\mathbf{H}}) \leq \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + r^{1/2}\sigma_{\max}(\mathbf{H}_0) \leq \frac{\mathcal{D}}{\alpha} + r^{1/2}\sigma_{\max}(\mathbf{H}_0), \quad (\text{B.192})$$

$$\begin{aligned} \kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) &= \frac{\sigma_{\max}(\mathbf{P}_0(\widehat{\mathbf{H}}))}{\sigma_{\min}(\mathbf{P}_0(\widehat{\mathbf{H}}))} \leq \frac{\mathcal{D}(\mathbf{P}_0(\widehat{\mathbf{H}}), \mathbf{H}_0)^{1/2} + r^{1/2}\sigma_{\max}(\mathbf{H}_0)}{\sqrt{2}(\mu - 2\delta)} \\ &\leq \frac{\mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} + r^{1/2}\sigma_{\max}(\mathbf{H}_0)}{\sqrt{2}(\mu - 2\delta)} \leq \frac{\mathcal{D}}{\alpha(\mu - 2\delta)\sqrt{2}} + \frac{r^{1/2}\sigma_{\max}(\mathbf{H}_0)}{(\mu - 2\delta)\sqrt{2}}. \end{aligned} \quad (\text{B.193})$$

Replacing these in (B.187) we have

$$\begin{aligned} \mathcal{D} &\leq 2(1 + 2\alpha) \left[ \frac{r^{3/2}\mathcal{D}\delta}{\alpha(\mu - 2\delta)\sqrt{2}} + \frac{r^2\sigma_{\max}(\mathbf{H}_0)\delta}{(\mu - 2\delta)\sqrt{2}} + \frac{\mathcal{D}r^{1/2}\delta}{\alpha\mu} \right. \\ &\quad \left. + \frac{r\sigma_{\max}(\mathbf{H}_0)\delta}{\mu} + \frac{r\|\mathbf{z}_0\|_2\delta}{\mu} \right] + 3\delta\sqrt{r}. \end{aligned} \quad (\text{B.194})$$

Therefore,

$$\begin{aligned} \mathcal{D} &\left[ 1 - \frac{\sqrt{2}(1 + 2\alpha)r^{3/2}\delta}{\alpha(\mu - 2\delta)} - \frac{2(1 + 2\alpha)r^{1/2}\delta}{\alpha\mu} \right] \\ &\leq 2(1 + 2\alpha) \left[ \frac{r^2\sigma_{\max}(\mathbf{H}_0)\delta}{(\mu - 2\delta)\sqrt{2}} + \frac{r\sigma_{\max}(\mathbf{H}_0)\delta}{\mu} + \frac{r\|\mathbf{z}_0\|_2\delta}{\mu} \right] + 3\delta\sqrt{r} \end{aligned} \quad (\text{B.195})$$

Notice that condition (B.191) implies that  $\mu - 2\delta \geq \mu/2$  and

$$\frac{\sqrt{2}(1 + 2\alpha)r^{3/2}\delta}{\alpha(\mu - 2\delta)} + \frac{2(1 + 2\alpha)r^{1/2}\delta}{\alpha\mu} \leq \frac{1}{2}. \quad (\text{B.196})$$

Using the previous two equations, under condition (B.191) we have

$$\begin{aligned} \mathcal{D} &\leq \frac{4(1 + 2\alpha)r\delta}{\mu} \left[ \frac{5r\sigma_{\max}(\mathbf{H}_0)}{2} + \|\mathbf{z}_0\|_2 \right] + 3\delta\sqrt{r} \\ &\leq \frac{4(1 + 2\alpha)r^2}{\mu} \left[ \frac{5\sigma_{\max}(\mathbf{H}_0)}{2} + \frac{\|\mathbf{z}_0\|_2}{r} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}} \right] \delta. \end{aligned} \quad (\text{B.197})$$

Combining this with (B.188), and using the fact that  $1 + 2\alpha \leq 3$ , we have under condition (B.191)

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{12(1 + \sqrt{2})r^{5/2}\kappa(\mathbf{H}_0)}{\mu\alpha} \left( \frac{5\sigma_{\max}(\mathbf{H}_0)}{2} + \frac{\|\mathbf{z}_0\|_2}{r} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}} \right) \delta \quad (\text{B.198})$$

$$\leq \frac{29\sigma_{\max}(\mathbf{H}_0)\kappa(\mathbf{H}_0)r^{5/2}}{\alpha\mu} \left( \frac{5}{2} + \frac{\|\mathbf{z}_0\|_2}{r\sigma_{\max}(\mathbf{H}_0)} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}\sigma_{\max}(\mathbf{H}_0)} \right) \delta. \quad (\text{B.199})$$

Note that using (B.184), (B.185) and since  $\alpha \geq 0$

$$\frac{\|\mathbf{z}_0\|_2}{r\sigma_{\max}(\mathbf{H}_0)} \leq 1, \quad \frac{3\mu}{4(1 + 2\alpha)r^{3/2}\sigma_{\max}(\mathbf{H}_0)} \leq \frac{3}{4\sqrt{2}}. \quad (\text{B.200})$$

Therefore,

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{120\sigma_{\max}(\mathbf{H}_0)\kappa(\mathbf{H}_0)r^{5/2}}{\alpha\mu} \delta. \quad (\text{B.201})$$

Thus,

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}}) \leq \frac{C_*^2 r^5}{\alpha^2} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2^2, \quad (\text{B.202})$$

where  $C_*$  is defined in Theorem 1.

Next, consider the case in which

$$\delta = \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \leq \frac{\alpha\mu}{330\kappa(\mathbf{H}_0)r^{5/2}}, \quad (\text{B.203})$$

Note that using (B.184), (B.185) and since  $1 + 2\alpha \leq 3$ , this condition on  $\delta$  implies that

$$\delta \leq \frac{\alpha\mu\sigma_{\min}(\mathbf{H}_0)}{12r(1 + 2\alpha)(5r^{3/2}\sigma_{\max}(\mathbf{H}_0) + 2\|\mathbf{z}_0\|_2r^{1/2} + 3\mu)}. \quad (\text{B.204})$$

In particular, condition (B.191) holds. Hence, using equation (B.197) we get

$$\mathcal{D} \leq \frac{4(1 + 2\alpha)r^2}{\mu} \left[ \frac{5\sigma_{\max}(\mathbf{H}_0)}{2} + \frac{\|\mathbf{z}_0\|_2}{r} + \frac{3\mu}{4(1 + 2\alpha)r^{3/2}} \right] \delta \leq \frac{\alpha\sigma_{\min}(\mathbf{H}_0)}{6\sqrt{r}}. \quad (\text{B.205})$$

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

$$\begin{aligned} \mathbf{P}_0(\mathbf{h}) &= \widetilde{\mathbf{P}}_0 \mathbf{h} + \mathbf{x}_0 = \widetilde{\mathbf{P}}_0 \widehat{\mathbf{H}}^\top \boldsymbol{\pi} + \mathbf{x}_0 = \sum_{i=1}^r \pi_i \left( \widetilde{\mathbf{P}}_0 \widehat{\mathbf{H}}^\top \mathbf{e}_i + \mathbf{x}_0 \right) \\ &= \sum_{i=1}^r \pi_i \mathbf{P}_0(\widehat{\mathbf{h}}_i) \in \text{conv}(\mathbf{P}_0(\widehat{\mathbf{H}})) \end{aligned} \quad (\text{B.206})$$

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

$$\mathbf{P}_0(\text{conv}(\widehat{\mathbf{H}})) \subseteq \text{conv}(\mathbf{P}_0(\widehat{\mathbf{H}})). \quad (\text{B.207})$$

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

$$\begin{aligned} \mathcal{D}(\mathbf{h}_0, \mathbf{P}_0(\widehat{\mathbf{H}})) &= D(\mathbf{h}_0, \text{conv}(\mathbf{P}_0(\widehat{\mathbf{H}}))) \leq D(\mathbf{h}_0, \mathbf{P}_0(\text{conv}(\widehat{\mathbf{H}}))) \\ &\leq D(\mathbf{h}_0, \mathbf{P}_0(\boldsymbol{\Pi}_{\text{conv}(\widehat{\mathbf{H}})}(\mathbf{h}_0))). \end{aligned} \quad (\text{B.208})$$

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

$$D(\mathbf{h}_0, \mathbf{P}_0(\boldsymbol{\Pi}_{\text{conv}(\widehat{\mathbf{H}})}(\mathbf{h}_0))) \leq D(\mathbf{h}_0, \boldsymbol{\Pi}_{\text{conv}(\widehat{\mathbf{H}})}(\mathbf{h}_0)) = D(\mathbf{h}_0, \text{conv}(\widehat{\mathbf{H}})) = \mathcal{D}(\mathbf{h}_0, \widehat{\mathbf{H}}). \quad (\text{B.209})$$

This implies that

$$\mathcal{D}(\mathbf{H}_0, \mathbf{P}_0(\widehat{\mathbf{H}})) \leq \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}}). \quad (\text{B.210})$$

Therefore, using (B.190), (B.210) and (B.205) we get

$$\begin{aligned} \mathcal{D}(\mathbf{H}_0, \mathbf{P}_0(\widehat{\mathbf{H}}))^{1/2} + \mathcal{D}(\mathbf{P}_0(\widehat{\mathbf{H}}), \mathbf{H}_0)^{1/2} &\leq \mathcal{D}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} + \mathcal{D}(\widehat{\mathbf{H}}, \mathbf{H}_0)^{1/2} \\ &\leq \frac{\mathcal{D}}{\alpha} \leq \frac{\sigma_{\min}(\mathbf{H}_0)}{6\sqrt{r}}. \end{aligned} \quad (\text{B.211})$$

Hence, in this case Lemma B.3 implies that

$$\sigma_{\max}(\widehat{\mathbf{H}}) \leq 2\sigma_{\max}(\mathbf{H}_0), \quad (\text{B.212})$$

$$\kappa(\mathbf{P}_0(\widehat{\mathbf{H}})) \leq \frac{7\kappa(\mathbf{H}_0)}{2}. \quad (\text{B.213})$$

Replacing this in (B.187), we have

$$\mathcal{D} \leq (1 + 2\alpha)r^{1/2} \left( 7r\delta\kappa(\mathbf{H}_0) + \frac{4\sigma_{\max}(\mathbf{H}_0)\delta + 2\sqrt{r}\|\mathbf{z}_0\|_2\delta}{\mu} \right) + 3\delta r^{1/2} \quad (\text{B.214})$$

$$\leq 3\delta\sqrt{r} \left( 8r\kappa(\mathbf{H}_0) + \frac{4\sigma_{\max}(\mathbf{H}_0) + 2\sqrt{r}\|\mathbf{z}_0\|_2}{\mu} \right) \quad (\text{B.215})$$

Hence, using (B.188) under assumption (B.203), we have

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq 3(1 + \sqrt{2})\kappa(\mathbf{H}_0)r \left( 8r\kappa(\mathbf{H}_0) + \frac{4\sigma_{\max}(\mathbf{H}_0) + 2\sqrt{r}\|\mathbf{z}_0\|_2}{\mu} \right) \frac{\delta}{\alpha} \quad (\text{B.216})$$

$$\leq 120\kappa(\mathbf{H}_0) \max \left\{ r\kappa(\mathbf{H}_0), \frac{\sigma_{\max}(\mathbf{H}_0) + \sqrt{r}\|\mathbf{z}_0\|_2}{\mu} \right\} \frac{r\delta}{\alpha}. \quad (\text{B.217})$$

Hence, for  $C''_*$  as defined in the statement of the theorem, we get

$$\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} \leq \frac{C''_* r}{\alpha} \max_{i \leq n} \|\mathbf{Z}_{i,\cdot}\|_2 \quad (\text{B.218})$$

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*

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

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

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

The Lagrangian for this problem is

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

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

$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = 0, \tag{C.4}$$

and therefore

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

and the dual of the above optimization problem is

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

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  $\boldsymbol{\rho}^* = 2(\mathbf{u} - \mathbf{y}^*) = 2(\mathbf{u} - \Pi_{\text{conv}(\mathbf{X})}(\mathbf{u}))$  (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*

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

*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

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

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

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

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

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

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

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

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

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

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

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

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

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

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  $\ell_2$  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  $\text{conv}(\mathbf{X})$ ,  $\Delta^r$  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  $\Delta^r$  is closed,  $\Psi$  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  $\Delta^r$  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  $\Delta^r$ . Hence,

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

for some positive constants  $\lambda_1^-$ ,  $\lambda_1^+$ . 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  $\Psi$ .

## 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  $\gamma_t$  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

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

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

## E.2 Alternating minimization

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

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

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

The algorithm alternates between minimizing with respect to the weights  $(\mathbf{w}_i)_{i \leq n}$  (this can be done independently across  $i \in \{1, \dots, n\}$ ) and minimizing over  $(\mathbf{h}_\ell, \boldsymbol{\alpha}_\ell)$ , 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

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

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

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  $\mathbf{h}_\ell$  and writing the resulting objective function.

The pseudocode for this algorithm is given below.

**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  $\{\mathbf{h}_\ell^{(T)}\}_{1 \leq \ell \leq r}$ ;

- 1: For  $\ell \in \{1, \dots, r\}$ :
- 2:     Set  $\boldsymbol{\alpha}_\ell^{(0)} = \arg \min_{\boldsymbol{\alpha} \in \Delta^n} \|\mathbf{h}_\ell^{(0)} - \mathbf{X}\boldsymbol{\alpha}_\ell\|_2$ ;
- 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  $\mathbf{A}$ .

## F Further simulation results

In this section, we evaluate the performances of the proposed method by comparing it with a number of algorithms for non-negative matrix factorization under different settings.

### F.1 Further comparisons with algorithms from the literature

Figures 2 to 7 extend the comparison of Figures 1, 4 and 5 to seven alternative reconstruction methods in the literature. Namely:

- *No noise.* Figures 2 and 3 repeat the experiment of Figure 1 for seven new methods.

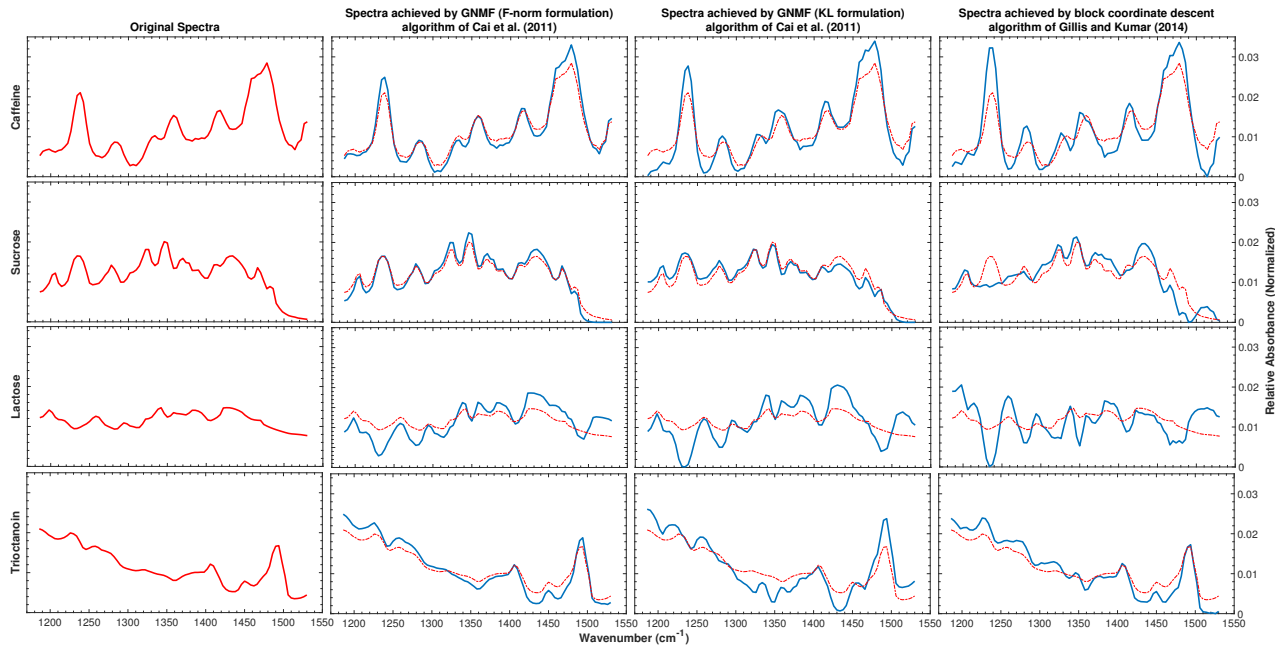


Figure 2: As in Figure 1 for three other methods.

- *Low noise.* Figures 4, 5 repeat the experiment of Figure 4 for the same seven new methods.
- *High noise.* Figures 6, 7 repeat the experiment of Figure 5 for the same seven new methods.

## F.2 Simulations with non-Gaussian correlated noise

In order to show the robustness of the proposed method to the noise model, we have repeated the experiments of Table 1 with dependent, non-Gaussian noise (see caption for a definition of the noise model). In Table 1, we report the average reconstruction error achieved by the same nine algorithms for non-negative matrix factorization. For each noise

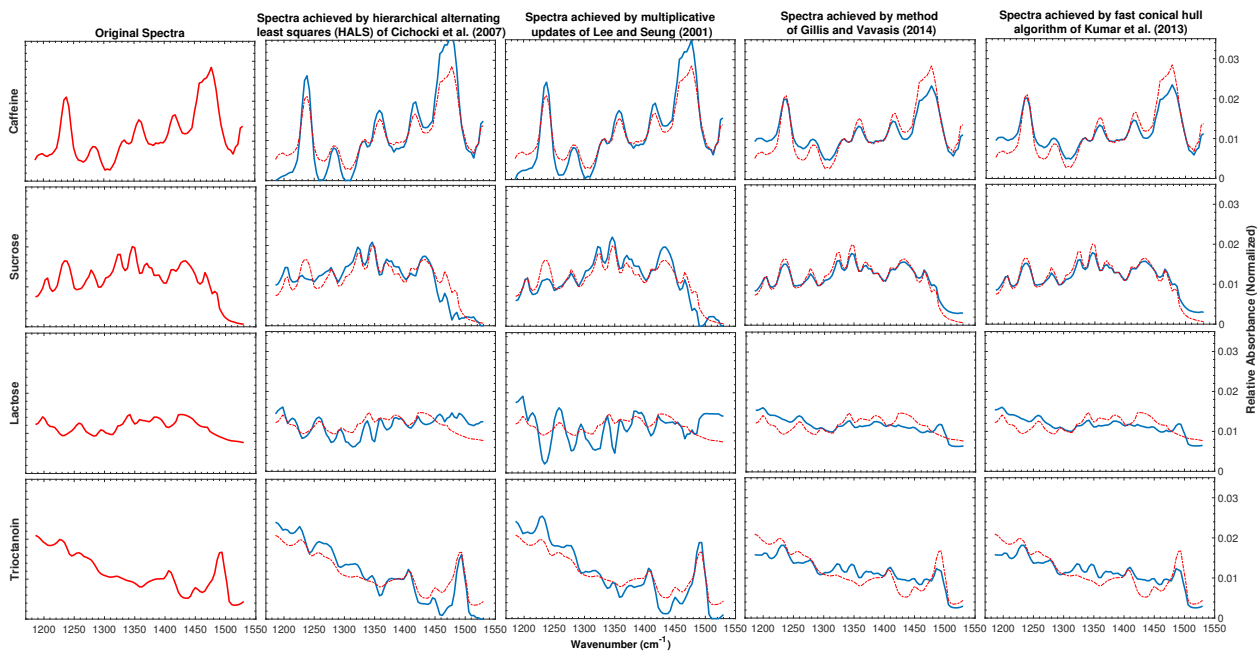


Figure 3: As in Figure 1 for four other methods.

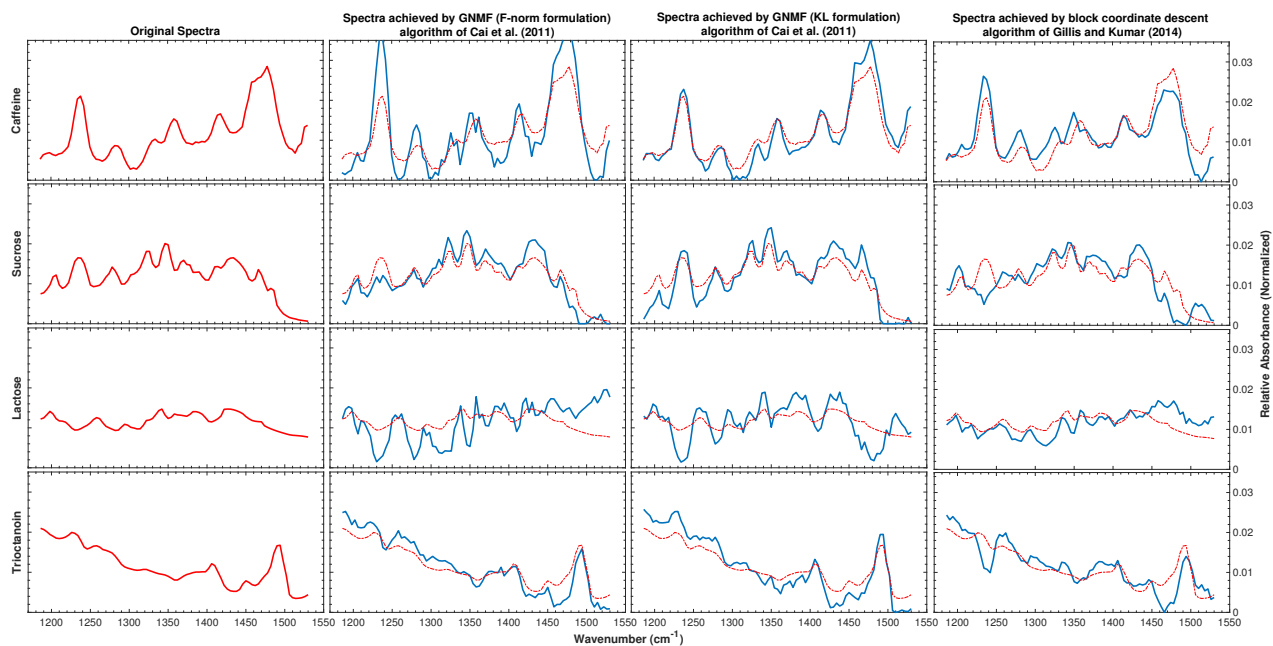


Figure 4: As in Figure 4 for three other methods.

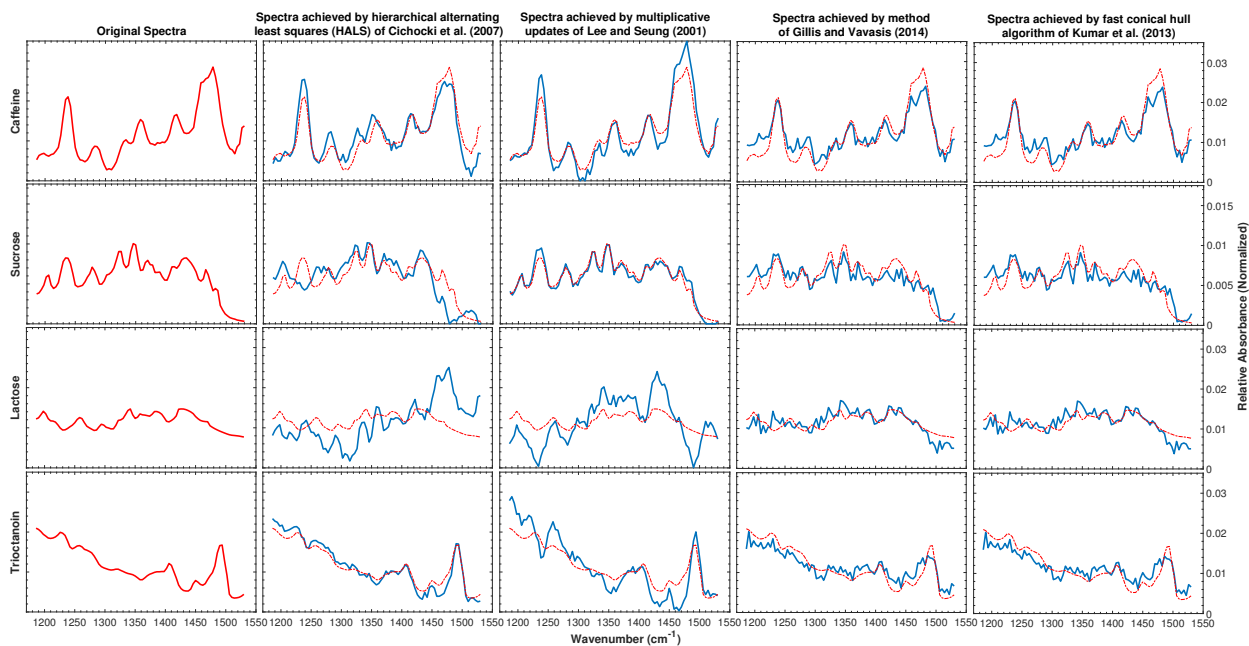


Figure 5: As in Figure 4 for four other methods.

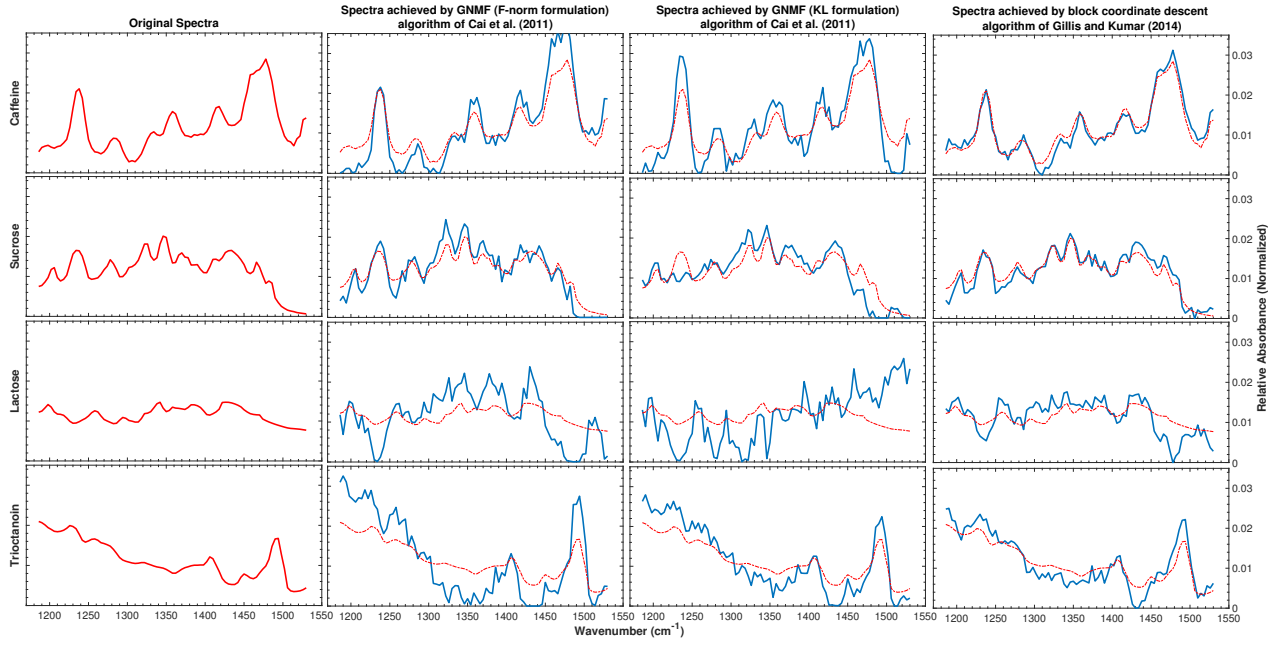


Figure 6: As in Figure 5 for three other methods.

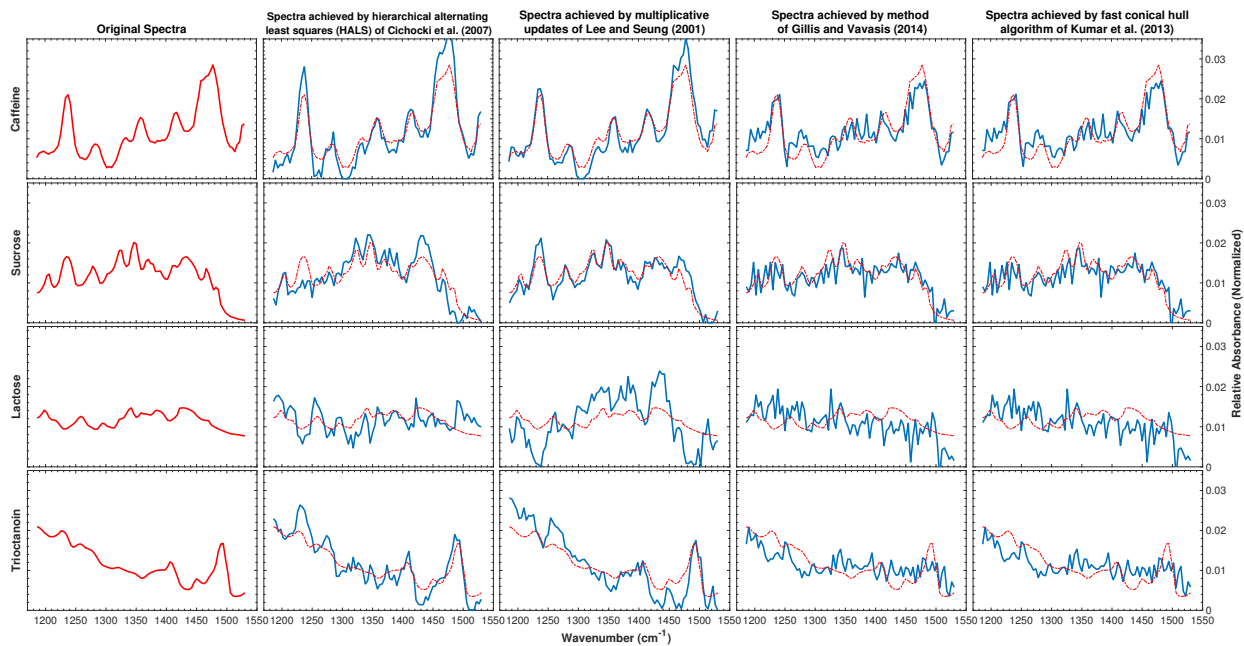


Figure 7: As in Figure 5 for four other methods.

level we run the various algorithms on 10 noise realizations for each noise level  $\sigma$ . We show in bold the smallest achieved error and the smallest error with data driven choice of the algorithm parameters.

$\sigma$	0	0.001	0.002	0.003	0.004	0.005	0.006	0.007	0.008	0.009	0.01
Projected gradient [Lin07]	0.062	0.058	0.06	0.07	0.089	0.101	0.122	0.133	0.142	0.143	0.151
Multiplicative update[LS01]	0.069	0.062	0.069	0.078	0.092	0.114	0.125	0.136	0.145	0.143	0.15
Fast Anchor Words [AGH <sup>+</sup> 13]	0.041	0.047	0.055	0.075	0.093	0.113	0.134	0.156	0.176	0.195	0.22
Block coordinate descent [GK15]	0.067	0.068	0.067	0.069	<b>0.077</b>	<b>0.086</b>	<b>0.092</b>	<b>0.095</b>	<b>0.096</b>	<b>0.098</b>	<b>0.099</b>
HALS [CZPA09]	0.073	0.077	0.074	0.077	0.95	0.112	0.117	0.131	0.140	0.145	0.151
GNMF [CHHH11] (Frobenius)	0.065	0.081	0.095	0.102	0.111	0.121	0.131	0.140	0.143	0.143	0.15
GNMF [CHHH11] (KL)	0.066	0.075	0.081	0.087	0.099	0.121	0.128	0.138	0.141	0.140	0.151
Recursive method [GV14]	0.034	0.04	0.053	0.068	0.089	0.111	0.13	0.15	0.17	0.19	0.21
Conical hull [KSK13]	0.034	0.04	0.052	0.068	0.088	0.111	0.13	0.15	0.17	0.19	0.21
Our method ( oracle $\lambda$ )	<b>0.005</b>	<b>0.015</b>	<b>0.039</b>	<b>0.06</b>	0.081	<b>0.091</b>	0.102	0.113	0.122	0.132	0.14
Our method (data driven $\lambda$ )	<b>0.006</b>	<b>0.021</b>	<b>0.041</b>	<b>0.067</b>	0.094	0.108	0.124	0.134	0.147	0.161	0.178

Table 1: Risk  $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}$  for reconstruction of the 4 spectra in Figure 1 under dependent, non-gaussian noise using some construction methods in different noise magnitudes. The trivial estimator  $\widehat{\mathbf{H}} = 0$  achieves  $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} = 0.231$ . For this table I have generated the data as  $\mathbf{X} = \mathbf{W}_0 \mathbf{H}_0 + \mu \mathbf{Q} \mathbf{Z}$  where  $Z_{ij}$ ,  $1 \leq i, j \leq n$  are i.i.d Laplace  $(0, 1)$  random variables and  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is a circulant matrix with first row equal to  $[1, \mathbf{0}_{n-(n_1+1)/2}, \rho \mathbf{1}_{(n_1-1)/2}]$  and  $\mu = (\sigma/\sqrt{2})(1 + (n_1 - 1)\rho^2/2)^{-1/2}$ . In these simulations, I have taken  $\rho = 0.5$ ,  $n_1 = 21$ . For the data driven row, parameter  $\lambda$  is chosen as in Section 4.3 with  $c_0 = 1.2$ .

### F.3 Simulations with other ground truth signals

We repeated the synthetic datasets experiments of Section 4 with using four other spectra as the ground truth signals. In these experiments we have used the reflection spectra of chalk, Maltose, Acetaminophen and baking soda from the NIST Chemistry WebBook dataset [LM]. For these signals we have  $d = 107$  and we generate  $n = 250$  data points as in Appendix A. We have also used different parameters to generate weight vectors in

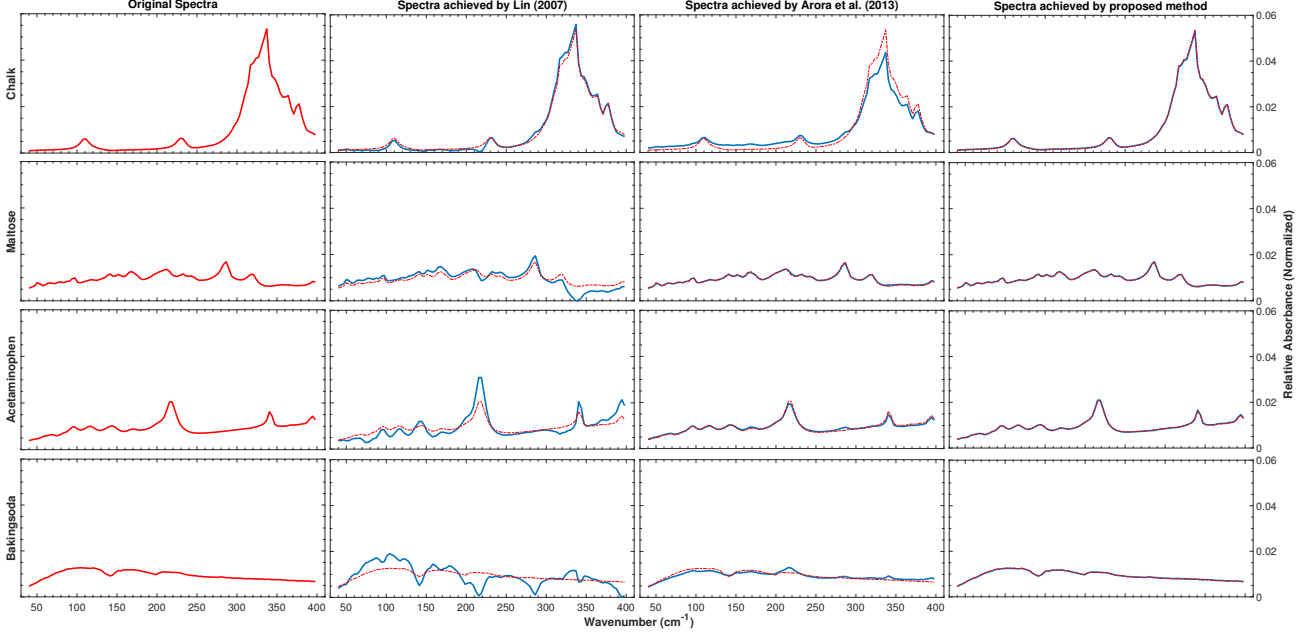


Figure 8: As in Figure 1 for four other ground truth spectra.

this case. Dirichlet parameters are chosen to be equal to one (instead of 5), the number of weight vectors with cardinality equal to 2 is equal to 12 and the number of weight vectors with cardinality equal to 3 is equal to 8. Other weight vectors have cardinality equal to 4. The recovered spectra of different algorithms in the noiseless and noisy settings are reported in Figures 8 to 16.

In Table 2, we report the average reconstruction error achieved by the same nine algorithms for non-negative matrix factorization on this dataset. For each noise level we use 10 noise realizations for each noise level  $\sigma$ .

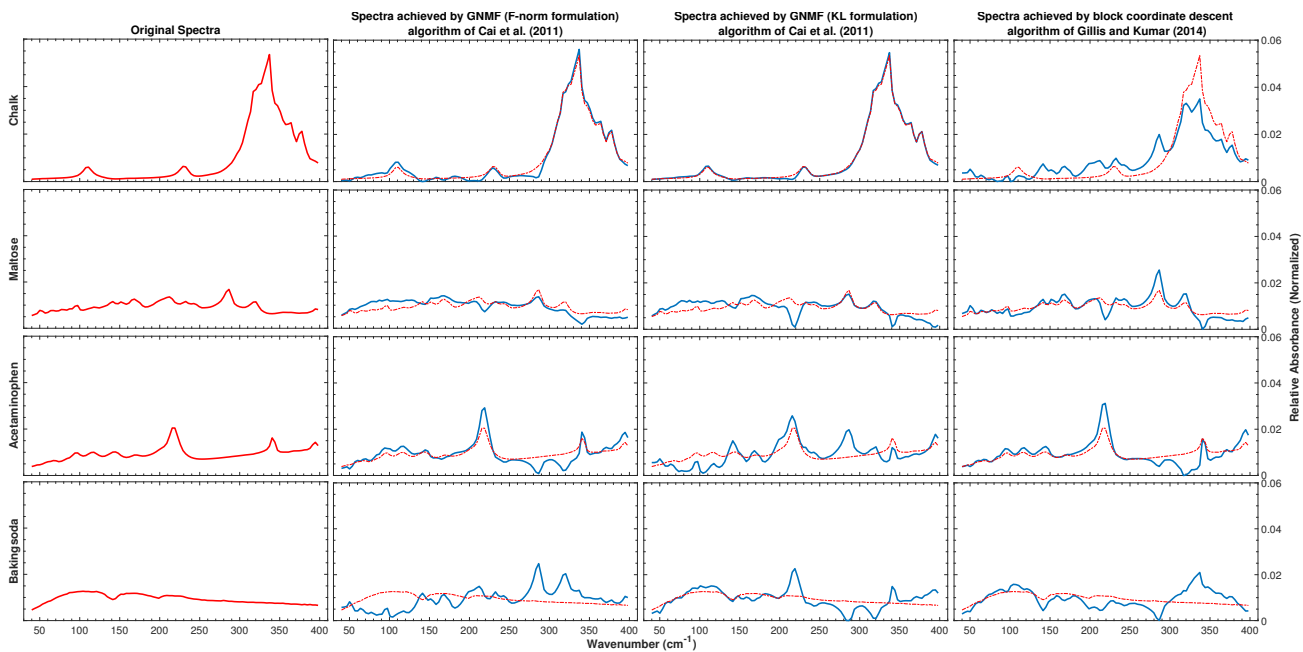


Figure 9: As in Figure 2 for four other ground truth spectra.

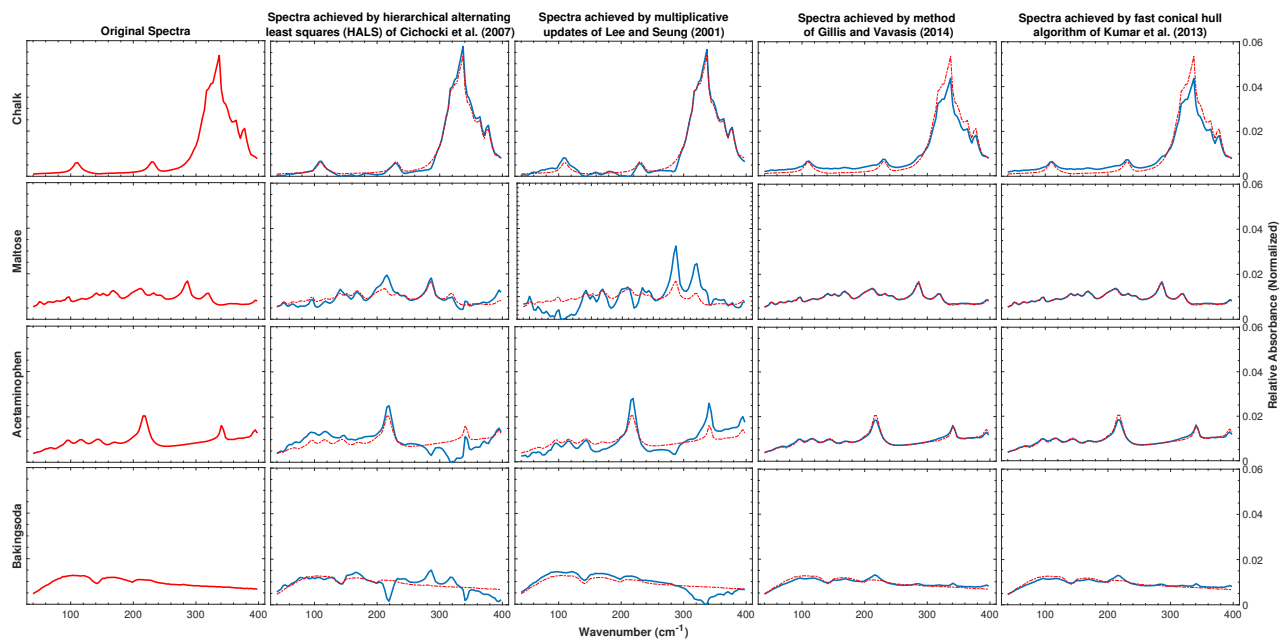


Figure 10: As in Figure 3 for four other ground truth spectra.

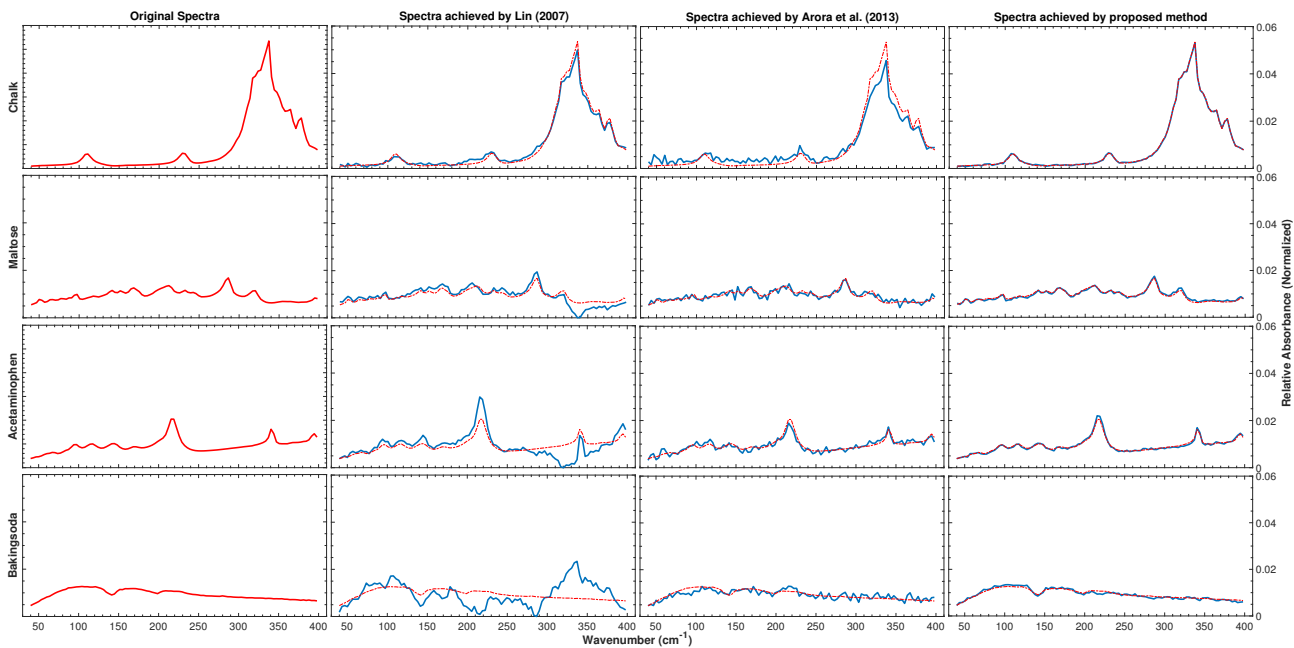


Figure 11: As in Figure 4 for four other ground truth spectra.

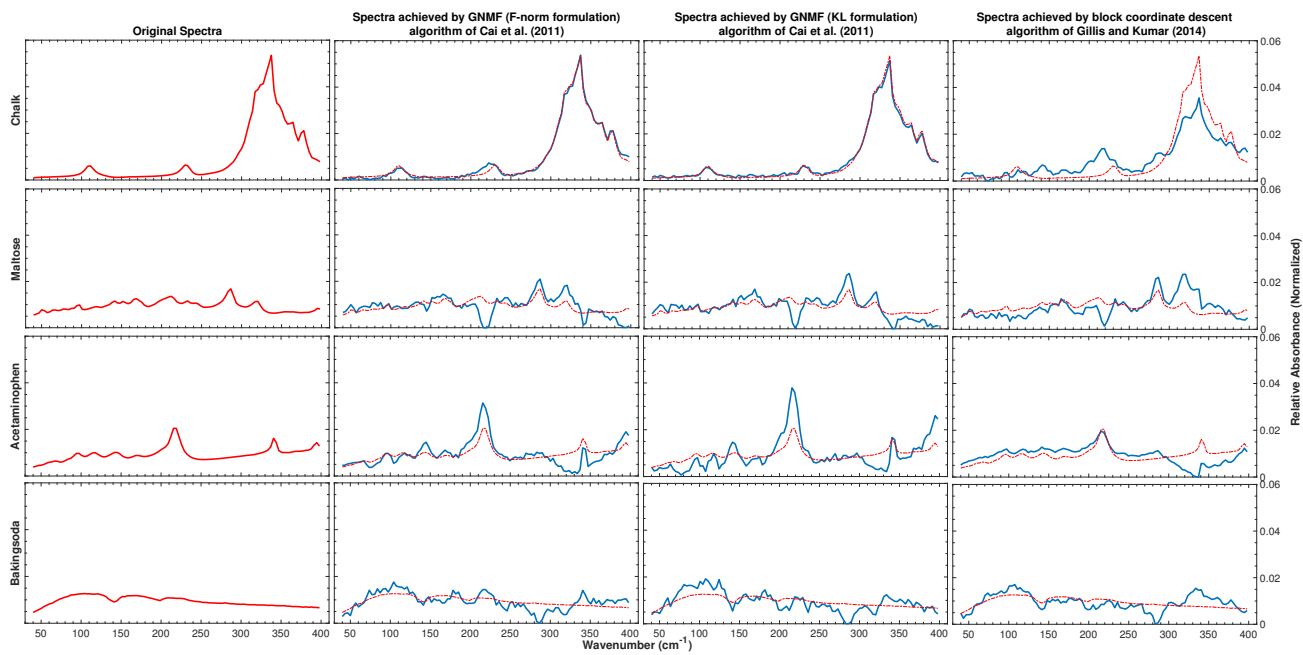


Figure 12: As in Figure 4 for four other ground truth spectra.

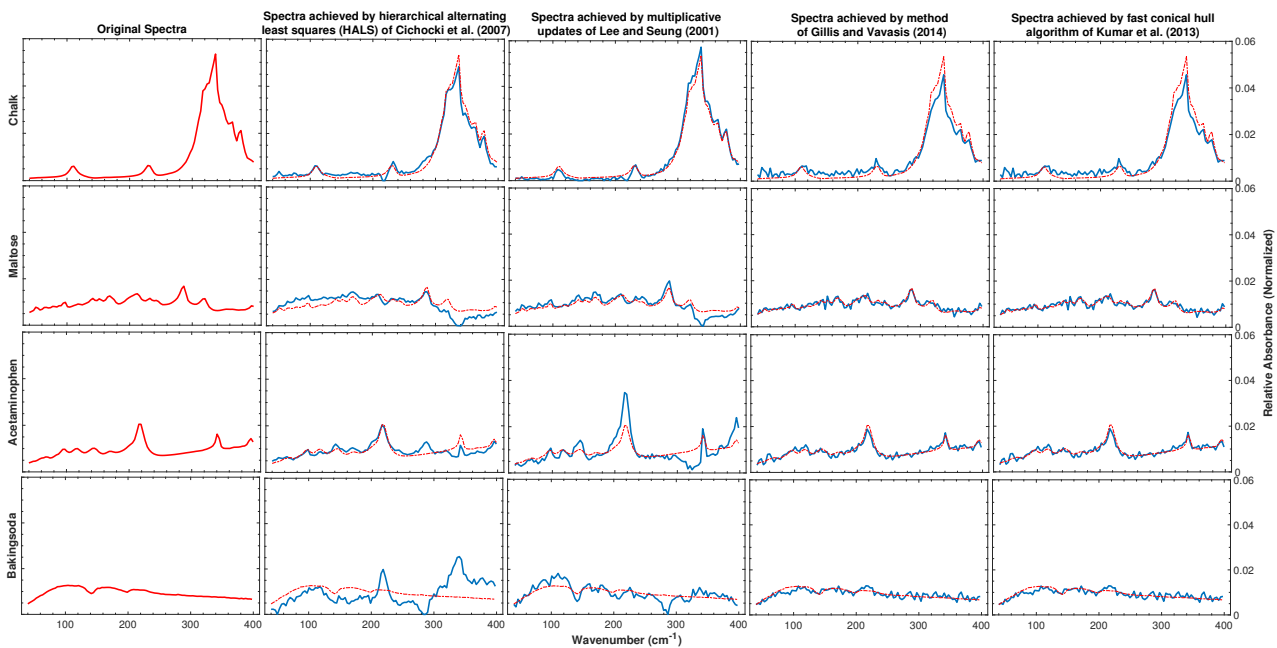


Figure 13: As in Figure 5 for four other ground truth spectra.

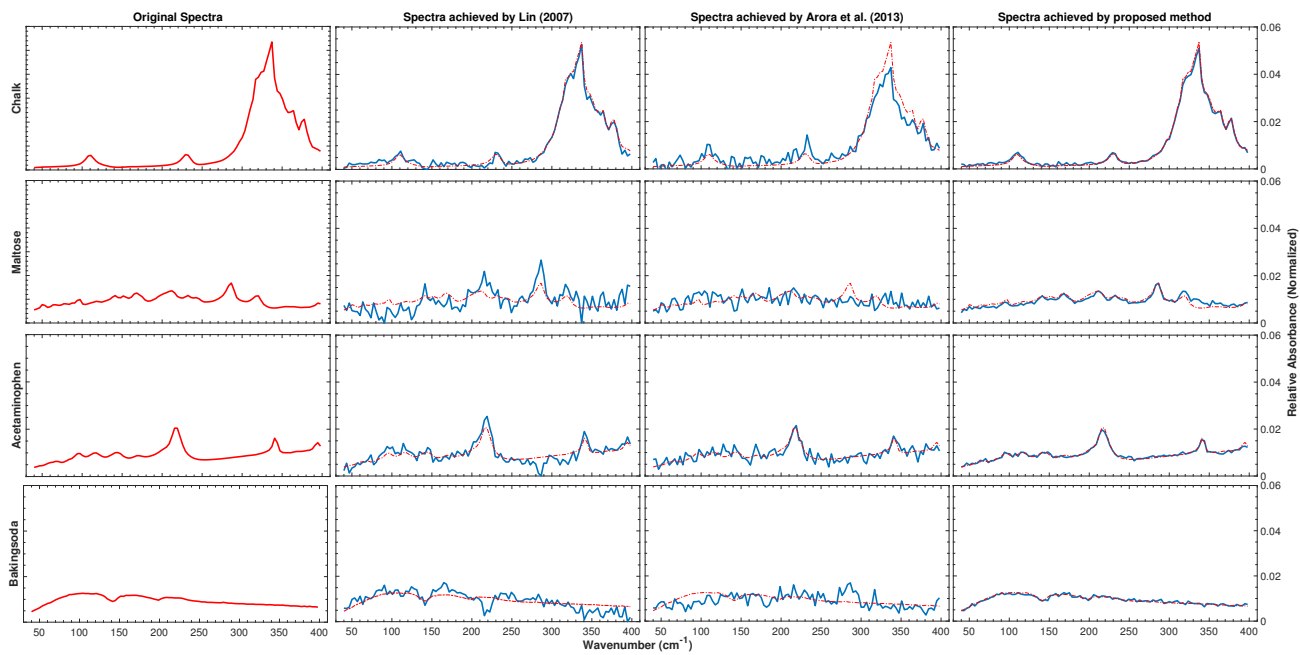


Figure 14: As in Figure 5 for four other ground truth spectra.

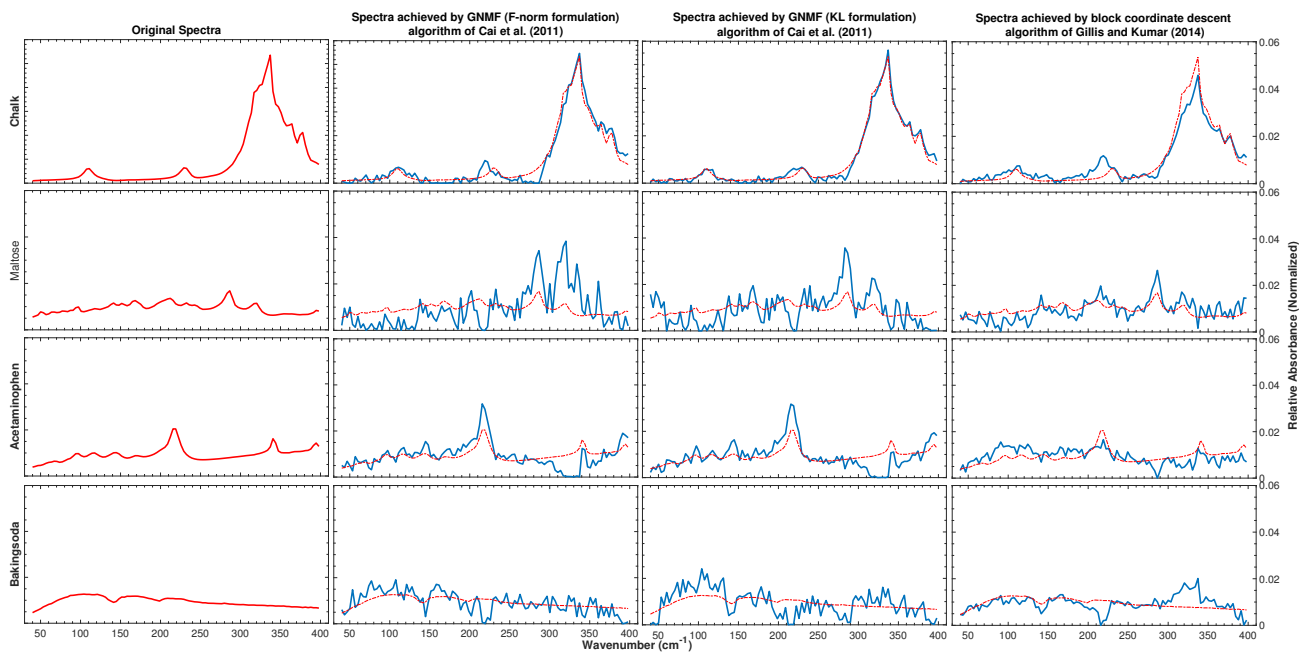


Figure 15: As in Figure 6 for four other ground truth spectra.

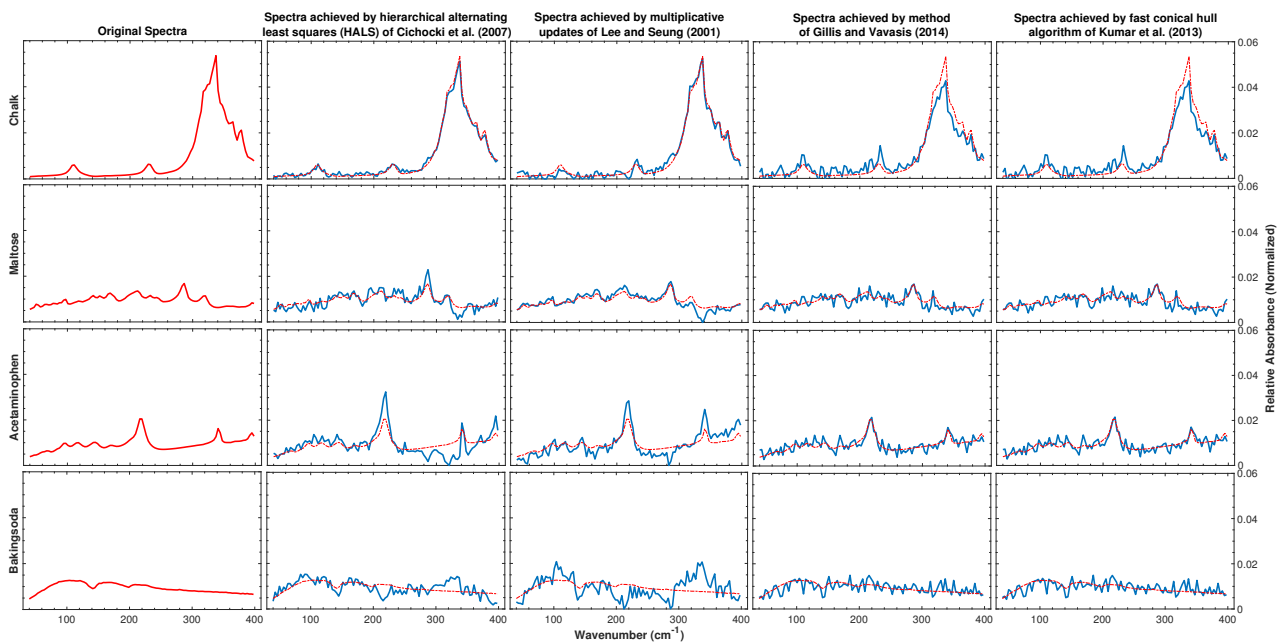


Figure 16: As in Figure 7 for four other ground truth spectra.

$\sigma$	0	0.001	0.002	0.003	0.004	0.005	0.006	0.007	0.008	0.009	0.01
Projected gradient [Lin07]	0.061	0.062	0.061	0.063	0.075	0.085	0.097	0.100	0.106	0.114	0.119
Multiplicative update [LS01]	0.064	0.064	0.064	0.069	0.075	0.086	0.092	0.099	0.110	0.115	0.119
Fast Anchor Words [AGH <sup>+</sup> 13]	0.045	0.041	0.056	0.075	0.094	0.119	0.140	0.161	0.183	0.208	0.231
Block coordinate descent [GK15]	0.089	0.089	0.085	0.086	0.083	0.091	0.095	0.088	<b>0.092</b>	<b>0.092</b>	<b>0.093</b>
HALS [CZPA09]	0.0601	0.059	0.064	0.067	0.073	0.089	0.095	0.100	0.109	0.115	0.120
GNMF [CHHH11] (Frobenius)	0.057	0.081	0.089	0.109	0.121	0.133	0.123	0.125	0.123	0.125	0.123
GNMF [CHHH11] (KL)	0.066	0.078	0.092	0.097	0.107	0.115	0.123	0.126	0.132	0.131	0.136
Recursive method [GV14]	0.031	0.039	0.054	0.074	0.093	0.117	0.138	0.159	0.183	0.205	0.229
Conical hull [KSK13]	0.031	0.039	0.054	0.074	0.093	0.117	0.138	0.159	0.183	0.205	0.229
Our method (oracle $\lambda$ )	<b>0.002</b>	<b>0.007</b>	<b>0.014</b>	<b>0.023</b>	<b>0.040</b>	<b>0.049</b>	<b>0.068</b>	<b>0.064</b>	<b>0.068</b>	<b>0.072</b>	<b>0.076</b>
Our method (Data driven $\lambda$ )	<b>0.003</b>	<b>0.009</b>	<b>0.014</b>	<b>0.027</b>	<b>0.052</b>	<b>0.06</b>	<b>0.075</b>	<b>0.087</b>	0.095	0.114	0.133

Table 2: Risk  $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2}$  for reconstruction of the 4 spectra in Figure 8 using some reconstruction methods in different noise magnitudes. The trivial estimator  $\widehat{\mathbf{H}} = 0$  achieves  $\mathcal{L}(\mathbf{H}_0, \widehat{\mathbf{H}})^{1/2} = 0.237$ . For the data driven row, parameter  $\lambda$  is chosen as in subsection 4.3 with  $c = 1.2$  in  $\{0.001, 0.002, 0.005, 0.007, 0.01, 0.03, 0.05, 0.07, 0.1, 0.2, 0.3, 0.5, 0.75, 1, 2, 5\}$ .

## References

- [AGH<sup>+</sup>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.
- [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.
- [CHHH11] Deng Cai, Xiaofei He, Jiawei Han, and Thomas S Huang, *Graph regularized nonnegative matrix factorization for data representation*, IEEE Transactions on Pattern Analysis and Machine Intelligence **33** (2011), no. 8, 1548–1560.

- [CZPA09] Andrzej Cichocki, Rafal Zdunek, Anh Huy Phan, and Shun-ichi Amari, *Non-negative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*, John Wiley & Sons, 2009.
- [CB94] Adele Cutler and Leo Breiman, *Archetypal analysis*, Technometrics **36** (1994), no. 4, 338–347.
- [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.
- [GK15] Nicolas Gillis and Abhishek Kumar, *Exact and heuristic algorithms for semi-nonnegative matrix factorization*, SIAM Journal on Matrix Analysis and Applications **36** (2015), no. 4, 1404–1424.
- [KSK13] Abhishek Kumar, Vikas Sindhwani, and Prabhanjan Kambadur, *Fast conical hull algorithms for near-separable non-negative matrix factorization*, International Conference on Machine Learning, 2013, pp. 231–239.
- [LS01] ———, *Algorithms for non-negative matrix factorization*, Advances in neural information processing systems, 2001, pp. 556–562.
- [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).

- [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.
- [Zie12] Günter M Ziegler, *Lectures on polytopes*, vol. 152, Springer Science & Business Media, 2012.